Institute of Mathematical Statistics COLLECTIONS Volume 00

# Optimality

The Third Erich L. Lehmann Symposium

Javier Rojo, Editor



Institute of Mathematical Statistics Collections

Series Editor: Anirban DasGupta

The production of the *Institute of Mathematical Statistics Collections* is managed by the IMS Office: Rong Chen, Treasurer and Elyse Gustafson, Executive Director.

Library of Congress Control Number: 0000000000 International Standard Book Number 0-000000-00-0 International Standard Serial Number 0749-2170 Copyright © 2009 Institute of Mathematical Statistics All rights reserved Printed in the United States of America

# Contents

| Preface<br>Javier Rojo               |
|--------------------------------------|
| ontributors to this volume           |
| viii                                 |
| cientific program                    |
| ix                                   |
| ist of Participants                  |
| xvii                                 |
| cknowledgement of referees' services |
| xxii                                 |

#### PAPERS

#### Honoring Erich Leo Lehmann

| Erich L. Lehmann, The Lehmann Symposia, and November 20 <sup>th</sup> 1917<br>Javier Rojo             | 1 |
|---|---|
| The Honorable Erich L. Lehmann         Stephen Stigler  | 3 |
| Optimality  |   |
| Some History of Optimality<br>Erich L. Lehmann  | 1 |
| An Optimality Property of Bayes' Test Statistics<br>Raghu Raj Bahadur and Peter J. Bickel 11          | 3 |
| On the Non-Optimality of Optimal Procedures<br>Peter J. Huber   | 1 |
| Semi-parametric Inference   |   |
| Proportional Hazards Regression with Unknown Link Function<br>Wei Wang, Jane-Ling Wang and Qihua Wang | 3 |
| Semiparametric Models and Likelihood - The Power of Ranks<br>Kjell Doksum and Akichika Ozeki          | 3 |
| Bootstrap   |   |
| On Bootstrap Tests of Hypotheses<br>Wei-Yin Loh and Wei Zheng   | 3 |
| Nonparametric Inference   |   |
| Nonparametric Estimation for Lévy Models Based on Discrete-Sampling<br>José E. Figueroa-López         | 7 |
| On the Estimation of Symmetric Distributions under Peakedness Order Con-<br>straints                  | _ |
| Javier Rojo ana Jose Batun-Cutz   | ſ |

#### Contents

iv

| Functional Data Analysis  |   |
|---|---|
|   |   |
| a runctional Generalized Linear Wodel with Curve Selection in Cervical Pre-<br>cancer Diagnosis using Fluorescence Spectroscopy |   |
| Hongxiao Zhu and Dennis D. Cox  | l |
| Nonparametric Estimation of Hamadynamic Despanse Eurotion, A Everyonay  |   |
| Domain Approach   |   |
| Ping Bai, Young Truong and Xuemei Huang   |   |
| Mixed Models, Posterior Means and Penalized Least-Squares   |   |
| Yolanda Muñoz Maldonado   | 7 |
|   |   |
| $\mathbf{Probability}$  |   |
| From Charged Polymers to Random Walk in Random Scenery  |   |
| Xia Chen and Davar Khoshnevisan   | 7 |
| Becovery of Distributions via Moments   |   |
| Robert M. Mnatsakanov and Artak S. Hakobyan   | 2 |
|   |   |
| Asymptotic Theory   |   |
| Asymptotic Efficiency of Simple Decisions for the Compound Decision Problem   |   |
| Eitan Greenshtein and Ya'acov Ritov   | 5 |
| Large Sample Statistical Inference for Skew-Symmetric Families on the Real  |   |
| Line  |   |
| Rolando Cavazos–Cadena and Graciela González–Farías   | 6 |
|   |   |
| Multiple Comparison   |   |
| Parametric Mixture Models for Estimating the Proportion of True Null Hy-  |   |
| potheses and Adaptive Control of FDR  |   |
| Ajit C. Tamhane ana Jiaxiao Shi   | Ł |
| Bayesian Decision Theory for Multiple Comparisons   |   |
| Charles Lewis and Dorothy T. Thayer 326   | ) |
| Ad-hoc Networks   |   |
| The Challennes of Madel Objection Calenting and Fatimation for Addres Net   |   |
| work Data Sets  |   |
| Farinaz Koushanfar and Davood Shamsi  | 3 |
|   |   |
| Finance   |   |
| A Note on the Investment Proportions of a Minimum-Variance Equity Portfolio   |   |
| Wilhelmine von Turk   | 7 |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |
|   |   |

| v   |
|---|
| Preface   |
| The venue for the $3^{rd}$ Lehmann Symposium was the School of Engineering at               |
| Rice University from May $16^{th}$ through May $19^{th}$ , 2007. The collection of refereed |
| papers included in this volume represents a selection of the papers submitted for           |
| publication. Most of the work was presented at the Symposium but there are some             |
| contributions that were submitted by participants who did not present their work            |
| during the conference.  |
| All activities of the Symposium, except for a banquet held at the student center,           |
| were held in Duncan Hall – home of the Statistics Department. Duncan Hall's floor           |
| plan, with its open atrium, its main auditorium, and several conveniently located           |
| meeting rooms, allows for, and facilitates, interaction among the participants.             |
| As it has been the tradition of the Symposia, the event opens with a session for            |
| young investigators. The purpose of initiating the Symposia in this way is to free          |
| the young investigators from this activity, and introduce them to other more senior         |
|   |

young inve ree the young i ior investigators with the goal that the young investigators may more easily mingle with the group. For the third Lehmann Symposium the four young investigators were Yolanda Muñoz Maldonado, Brisa Sánchez, Farinaz Koushanfar, and José Enrique Figueroa-López. At the end, due to unforeseen circumstances, José Enrique was moved to the probability session. All four young investigators provided motivating talks and three of them submitted their work for this volume. All four have a bright future ahead of them. 

It is also the tradition of the Symposia that the young investigators session is immediately followed by the first Plenary session and this spot has always been filled by Erich L. Lehmann. Erich provided a great lecture on the history of optimality. The rest of the program, I hope that the reader will agree with me, was excellent.

The papers presented here cover several areas: some of the works consider classical aspects of the discipline and others deal with contemporary aspects of the theory and applications of statistics. Thus, the reader will find a fascinating section dedicated to the subject of optimality. Lehmann, Bahadur and Bickel, and Huber

v

vi

provide excellent discussions on various aspects of optimality. Semi-parametric and non-parametric inference, bootstrap tests of hypotheses, functional data analysis, asymptotic theory, ad-hoc networks, and finance are some of the areas represented in the volume. Intentionally, I left probability to the end. It has been a goal of the Symposium to have a probability component. It is felt that the perceived distancing of probability and statistics, even at the level of Ph.D. programs, cannot be healthy. Future Lehmann Symposia will continue to encourage a closer relationship between the two subjects. The Symposium could not occur without the financial support of several generous contributors. Several institutions have provided support for the series of Symposia. I want, however, to acknowledge the explicit support and work of the individuals within those institutions responsible for securing the funding. Demissie Alemayehu of Pfizer has been a constant and faithful supporter of the Symposia and has been responsible for Pfizer's generous contributions to the  $2^{nd}$  and  $3^{rd}$  Lehmann Sym-posia. I also want to acknowledge the support of the National Science Foundation. Shulamith Gross, Grace Yang, and Gábor Székely have been instrumental in sup-porting the  $2^{nd}$  and  $3^{rd}$  Lehmann Symposia. Gary Rosner, from the U.T. MD Anderson Cancer Center, has worked to obtain MD Anderson Cancer support for the last two Symposia. Kathleen O'Hara at the Mathematical Sciences Research

Institute (MSRI) and James O. Berger at the Statistics and Applied Mathemat-ics Sciences Institute (SAMSI) provided financial support for the  $3^{rd}$  Symposium. MSRI supported a proposal to hold the event at their facilities but plans changed. You can read the details in the first article of the volume. The efforts of Robert Hardy of the University of Texas School of Public Health, and Rudy Guerra of the Gulf Coast Consortia, are also gratefully acknowledged. Victor Pérez-Abreu, of the Centro de Investigación en Matemáticas (CIMAT) provided support for the  $1^{st}$  and  $2^{nd}$  Symposia. 

Last, but not least, I want to acknowledge the support of Rice University through

indefatigable efforts and continued help in the preparation of this volume.

I also want to give special recognition to my student Tuan S. Nguyen for his

| 1  | the School of Engineering and its Department of Statistics for allowing me to engage | 1  |
|----|--|----|
| 2  | in these activities.   | 2  |
| 3  |  | 3  |
| 4  |  | 4  |
| 5  | Javier Rojo January 30, 2009   | 5  |
| 6  | Statistics Department<br>Rice University   | 6  |
| 7  | Houston, Texas   | 7  |
| 8  |  | 8  |
| 9  |  | 9  |
| 10 |  | 10 |
| 11 |  | 11 |
| 12 |  | 12 |
| 13 |  | 13 |
| 14 |  | 14 |
| 15 |  | 15 |
| 16 |  | 16 |
| 17 |  | 17 |
| 18 |  | 18 |
| 19 |  | 19 |
| 20 |  | 20 |
| 21 |  | 21 |
| 22 |  | 22 |
| 23 |  | 23 |
| 24 |  | 24 |
| 25 |  | 25 |
| 26 |  | 26 |
| 27 |  | 27 |
| 28 |  | 28 |
| 29 |  | 29 |
| 30 |  | 30 |
| 31 |  | 31 |
| 32 |  | 32 |
|    |  |    |

imsart-coll ver. 2008/08/29 file: Rojo\_Preface3.tex date: March 26, 2009

vii

# Contributors to this volume

Bahadur, R. R., University of Chicago

Bai, P., University of North Carolina at Chapel Hill
Batún-Cutz, J., Universidad Autónoma de Yucatán, Mérida, México
Bickel, P. J., University of California, Berkeley

Cavazos–Cadena, R., Universidad Autónoma Agraria Antonio Narro Chen, X., University of Tennessee Cox, D. D., Rice University

Doksum, K., University of Wisconsin, Madison

Figueroa-López, J., Purdue University

González–Farías, G., Centro de Investigación en Matemáticas A. C. Greenshtein, E., Duke University

Hakobyan, A. S., West Virginia University Huang, X., University of North Carolina at Chapel Hill Huber, P. J., Switzerland

Khoshnevisan, D., University of Utah Koushanfar, F., Rice University

Lehmann, E. L., University of California, Berkeley Lewis, C., Fordham University Loh, W., University of Wisconsin, Madison

Mnatsakanov, R. M., West Virginia University Muñoz Maldonado, Y., Michigan Technological University

Ozeki, A., University of Wisconsin, Madison

Ritov, Y., Jerusalem, Israel Rojo, J., Rice University

Shamsi, D., *Rice University* Shi, J., *Northwestern University* Stigler, S., *University of Chicago* 

Tamhane, A. C., Northwestern UniversityThayer, D. T., Educational Testing ServiceTruong, Y., University of North Carolina at Chapel Hill

Von Turk, W., American Century Investmens

Wang, J., University of California, DavisWang, Q., Chinese Academy of Science and The University of Hong KongWang, W., Harvard Medical School and Brigham and Women's Hospital

Zheng, W., University of Wisconsin, Madison Zhu, H., Rice University

# SCIENTIFIC PROGRAM

# The Third Erich L. Lehmann Symposium May 16 - 19, 2007 Rice University

#### Symposium Chair and Organizer

Javier Rojo Statistics Department, MS-138 Rice University 6100 Main Street Houston, TX 77005

# **Plenary Speakers**

| Erich L. Lehmann                   |  |
|------------------------------------|--|
| University of California, Berkeley |  |

Some history of optimality

Lawrence D. Brown The Wharton School University of Pennsylvania

James O. Berger Duke University

**Rodrigo Bañuelos** Purdue University

Bayes, hierarchical Bayes, and random effects

A unified view of regression, shrinkage, empirical

selection

Some recent developments in Bayesian model

Isoperimetric bounds for Lévy processes

**Peter J. Bickel** University of California, Berkeley

**Stephen M. Stigler** University of Chicago

Peter J. Huber

Willem R. van Zwet University of Leiden The collapse of particle filters

Karl Pearson and testing statistical hypotheses

On the non-optimality of optimal procedures

*Statistics and the law: the case of the nonchalant nurse* 

# **Invited Session Scientific Committee**

| Javier Rojo, Chair | <b>Rice University</b>             |
|--------------------|------------------------------------|
| Jane Ling Wang     | University of California, Davis    |
| Rudy Guerra        | Rice University                    |
| Juliet P. Shaffer  | University of California, Berkeley |
| Wei-Yin Loh        | University of Wisconsin, Madison   |
| Peter J. Bickel    | University of California, Berkeley |
| Kjell Doskum       | University of Wisconsin, Madison   |
| Yongzhao Shao      | New York University                |
| Demissie Alemayehu | Pfizer and Columbia University     |
| James O. Berger    | Duke University and SAMSI          |

# **Invited Sessions**

# Young Investigators

Javier Rojo, Organizer Yolanda Muñoz Maldonado, Chair

| <b>Brisa N. Sánchez</b><br>University of Michigan             | Residual-based diagnostics for structural equation models                 |
|---|---|
| <b>Yolanda Muñoz Maldonado</b><br>UT Health Sc Center Houston | Penalized least squares and frequentist and bayesian mixed-effects models |
| Farinaz Koushanfar<br>Rice University                         | How challenging is the data set?  |
| Chair and Dural Laws in the Are                               |   |

## <u>Statistical Problems in the Analysis of Genomic and Magnetic Resonance Imaging</u> <u>Data</u>

Wei-Yin Loh, Chair

| Sunduz Keles                     | Statistical issues arising in the study of |
|----------------------------------|--|
| University of Wisconsin, Madison | transcription regulation                   |

| Shaw-Hwa Lo         | A method toward mapping of common tra | aits |
|---------------------|---------------------------------------|------|
| Columbia University |                                       |      |

Young K. Truong UNC, Chapel Hill Spatio-temporal modeling for fMRI data

## **Optimality in Bioinformatics: Theory vs Practice**

Rudy Guerra, Chair

| <b>David Dahl</b><br>Texas A&M University      | Simultaneous inference for multiple testing and clustering via Dirichlet process mixture models |  |
|--|---|--|
| <b>Chad Shaw</b><br>Baylor College of Medicine | Using annotations in the analysis of genome scale data  |  |
| <i>Rudy Guerra</i><br>Rice University          | Incorporating biological knowledge in gene expression analysis                                  |  |

## <u>Regularized Methods of Classification and Estimation of Nonparametric Regression</u> <u>and Covariance Matrices When Data is High Dimensional</u>

Peter J. Bickel, Chair

Wei-Yin Loh University of Wisconsin *Regression and variable selection in large p, small n problems* 

**Debashis Paul** University of California, Davis Principal component analysis for structured highdimensional data

**Ya'acov Ritov** Hebrew University *Consistent learning methods are approximately local* 

## **Probability, Levy Process, and Applications**

Javier Rojo, Organizer Rodrigo Bañuelos, Chair

| <b>Dennis Cox</b><br>Rice University            | Multiscale models for chemical reaction processes |
|---|---|
| <b>Davar Khoshnevisan</b><br>University of Utah | On some applications of stable processes          |
| José Enrique Figueroa                           | Non-parametric estimation for some models driven  |
| Purdue University                               | by Levy processes                                 |

### **Multiplicity: Developments and Current Issues**

Juliet P. Shaffer, Chair

Charles Lewis Fordham University

Helmut Finner Deutsches Diabetes-Zentrum

**Ajit C. Tamhane** Northwestern University Bayesian decision theory for multiple comparisons

Testing for equivalence in k sample models

A mixture model approach to estimating the number of true null hypotheses and adaptive control of FDR

## **Recent Advances in Non- and Semi-parametric Modeling**

Jane-Ling Wang, Chair

**Kjell Doksum** University of Wisconsin Semi-parametric models based on transformations and extremes

Xihong Lin Harvard University

**Naisyin Wang** Texas A&M University Statistical challenges in analyzing mass spectrometry proteomic data

Analysis of hierarchical biomedical data using semiparametric models

# <u>Statistical Inference for Population Substructures via Clustering, Mixture Models and</u> <u>Other Approaches</u>

Demissie Alemayehu, Yongzhao Shao, Organizers Yongzhao Shao, Chair

**Bruce G. Lindsay** Pennsylvania State University

**Yuewu Xu** Fordham University

**Yongzhao Shao** New York University Modal inference: halfway between clustering and mixture analyses

*A limit theory for likelihood ratio test under unidentifiability for general dependent processes* 

Recent developments in likelihood theory with applications to testing homogeneity in finite mixture models and other models

## Semiparametric Models, Longitudinal Survival Data, False Discovery Rate, And Brain fMRI

Kjell Doksum, Chair

| Jane-Ling Wang<br>University of California at Davis              | Semi-parametric analysis of longitudinal data truncated by event-time     |
|--|---|
| Kam-Wah Tsui and Shijie Tang<br>University of Wisconsin, Madison | Simultaneous testing of multiple hypotheses<br>using generalized p-values |
| <b>Chunming Zhang</b><br>University of Wisconsin, Madison        | Semi-parametric detection of significant activation for brain fMRI        |

## **Multiple Testing and Subgroup Analysis**

James O. Berger, Chair

**Juliet P. Shaffer** University of California Multiplicity and subgroup analysis

**Peter Mueller** U.T. M.D. Anderson Cancer Center The optimal discovery procedure and Bayesian decision rules

**M. J. Bayarri** Valencia, Duke U. and SAMSI *Bayesian and frequentist handling of multiple U. testing* 

### **CONTRIBUTED PAPERS**

Nancy L. Glenn, University of South Carolina, Columbia: The GEM Algorithm

**Richard C. Ott**, Mesa State College: On the Operating Characteristics of Some Nonparametric Methodologies for the Classification of Distributions by Tail Behavior

**Xiaohui Wang**, University of Texas-Pan American: *Classifications of Proteomic Mass* Spectra and Other Curve Data

**Xiaohu Li**, School of Mathematics and Statistics, Lanzhou University, People's Republic of China: *Stochastic Comparison on Conditional Order Statistics - Some New Results* 

**Robert Mnatsakanov**, West Virginia University, Some Asymptotic Properties of Varying Kernel Density Estimator

**Monnie McGee**, Southern Methodist University, A Distribution Free Summarization Method for Affymetrix GeneChip® Arrays

**Changxiang Rui**, University of Arkansas, *Point and Block Prediction in Log-Gaussian random Fields: The Non-constant Mean Case* 

Qiang Zhao, Texas State University, San Marcos, Survival Analysis of Microarray Gene Expression Data Using Correlation Principal Component Regression

Suhasini Subba Rao, Texas A&M University, Normalised Least-squares estimation in time-varying arch models

**Santanu Chakraborty**, UT Pan-American, *Parametric Inference on Zero-Inflated Poisson distribution and its variants* 

**Hongxiao Zhu**, Rice University, A Functional Generalized Linear Model with application to Cervical Pre-cancer Diagnosis using Fluorescence spectroscopy

**Xiaowei Wu**, Rice University, Some Estimation and Hypothesis Testing Problems in Fluctuation Analysis

**John Fresen**, University of Missouri - Columbia, On the Definition of Weak Convergence of a Sequence of Random Elements

Victor De Oliveira, UT San Antonio, Objective Bayesian Analysis of Spatial Data with Measurement Error

Pang Du, Virginia Tech, Smoothing spline frailty model

**Graciela Gonzalez**, CIMAT, Some important issues in inference under certain types of singularities

The Third Erich L Lehmann Symposium Javier Rojo, Organizer/Chair Invited Program – All talks are in McMurtry Auditorium Duncan Hall 1055

| Saturday May 19  | Breakfast<br>Dincan Hall             | d Multiple Testing and<br>Subgroup analysis                               | Juliet P Shaffer<br>Peter Mueller<br>Susie Bayarri | Coffee Break                                   | Plenary Speaker 10:50<br>Peter J Huber        | Plenary Speaker 11:50<br>Willem van Zwet | CLOSING<br>LUNCHEON<br>DUNCAN HALL                                       | 1:00 – 2:30                                 |  |  |                                  |  |
|------------------|--------------------------------------|---|--|--|---|--|--|---|--|--|----------------------------------|--|
| Friday May 18    | Breakfast<br>Duncan Hall             | <u>Recent Advances in Non an</u><br><u>Semiparametric Modeling</u>        | Kjell Doksum<br>Xihong Lin<br>Naisyin Wang         | Coffee Break                                   | Plenary Speaker 10:50<br>Lawrence D Brown     | Lunch 11:50 – 1:15<br>Duncan Hall        | Inference for Substructures<br>Clustering, Mixtures,<br>Bruce G. Lindsay | Yuewu Xu<br>Yongzhao Shao                   | Coffee Break<br>Plenary Speaker<br>James O Berger    | Semiparametrics, longitudin<br>survival, false discovery, fM<br>Jane-Ling Wang     | Kam-Wah Tsui<br>Chunming Zhang   | CONTRIBUTED TALKS<br>6:05 – 7:05   |
| Thursday May 17  | Breakfast<br>Duncan Hall             | <u>Classification, nonparametric</u><br>regression: high dimensional data | Wei-Yin Loh<br>Debashis Paul<br>Ya'acov Ritov      | Coffee Break 10:30 – 10:50                     | Plenary Speaker 10:50<br>Rodrigo Bañuelos     | Lunch 11:50 – 1:15<br>Duncan Hall        | Probability. Lévy Processes and<br>Applications<br>Dennis D. Cox         | Davar Khoshnevisan<br>Jose Enrique Figueroa | Coffee Break<br>Plenary Speaker<br>Stephen M Stigler | <u>Multiplicity: Developments and</u><br><u>current issues</u><br>Charles Lewis    | Helmut Finner<br>Ajit C. Tamhane | Banquet Ballroom Ley Student Cr<br>Cash Bar/Music 6:45 – 7:30<br>Dinner/Music 7:30 – 10:00 |
| Wednesday May 16 | Breakfast 7:30 – 8:45<br>Dincan Hall | Opening Remarks   | Plenary Speaker<br>Erich L Lehmann                 | <u>Young Investigators</u><br>Brisa N. Sánchez | Yolanda Muñoz Maldonado<br>Farinaz Koushanfar | Lunch 12:00 – 1:15<br>Duncan Hall        | Genomics and magnetic resonance<br>Sunduz Keles                          | Shaw-Hwa Lo<br>Young K. Truong              | Coffee Break<br>Plenary Speaker<br>Peter J Bickel    | <u>Optimality in Bioinformatics:</u><br><u>Theory vs Practice</u><br>David B. Dahl | Chad Shaw<br>Rudy Guerra         | CONTRIBUTED TALKS<br>6:05 - 7:05   |
|                  | 7:30                                 | 8:45 9:00   |  | 10:00<br>10:30 - 10:50                         |   |  | 1:15 – 3:00  |   | 3:00 – 3:20<br>3:20 4:20                             | 4:20 – 6:05  |                                  |  |



# **List of Participants**

Susan Alber Department of Biostatistics 1515 Holcomber Blvd Houston, TX 77030 salber@mdanderson.org

Veera Baladandayuthapani 1515 Holcombe Boulevard Unit 447 Houston, TX 77030 veera@mdanderson.org

Rodrigo Bañuelos Department of Mathematics Purdue University West Lafayette, IN 47906 banuelos@math.purdue.edu

Rosa Bañuelos 2630 Bissonnet St Apt 167 Houston, TX 77005 banuelos@rice.edu

M.J. Bayarri SAMSI Research Triangle Park NC 27709 <u>susie.bayarri@uv.es</u>

Dror Berel 4045 linkwood dr. #605 Houston, TX 77025 dror.berel@gmail.com

James Berger PO Box 14006 Research Triangle Park NC 27709 berger@samsi.info

Nancy Bickel Statistics 367 Evans Hall UC Berkeley Berkeley, CA 94720 bickel@stat.berkeley.edu

Peter Bickel Statistics 367 Evans Hall UC Berkeley Berkeley, CA 94720 <u>bickel@stat.berkeley.edu</u> Lawrence D Brown Statistics Dept, 400 Huntsman Hall 3730 Walnut St Philadelphia, PA 19104 <u>lbrown@wharton.upenn.edu</u>

William F Bryant 6100 Main St. Houston, TX 77005 wfb1@rice.edu

Christopher P Calderon CAAM MS#138 Rice University Houston, TX 77251 <u>calderon@rice.edu</u>

Santanu Chakraborty Department of Mathematics University of TX - Panamerican Edinburg, TX 78539 <u>schakraborty@utpa.edu</u>

Jamie Chatman 7200 Almeda Rd, #818 Houston, TX 77054 <u>jchatman@rice.edu</u>

Dunlei Cheng 8080 North Central Expressway Suite 500 Dallas, TX 75206 dunleic@baylorhealth.edu

Raj Chikara 2700 Bay Area Blvd. Houston, TX 77058 <u>chhikara@uhcl.edu</u>

Dennis D Cox Rice University MS-138 6000 Main St. Houston, TX 77005 dcox@rice.edu

David B Dahl Department of Statistics 3143 TAMU College Station, TX 77843 <u>dahl@stat.tamu.edu</u> Sujay Datta Department of Statistics 3143 TAMU College Station, TX 77843 sdatta@stat.tamu.edu

Victor De Oliveira Management Sci Statistics University of Texas San Antonio, TX 78249 victor.deoliveira@utsa.edu

Kjell A Doksum Statistics Department University of Wisconsin Madison, WI 53706 <u>doksum@stat.wisc.edu</u>

Pang Du 13000 Foxridge Ln Apt J Blacksburg, VA 24060 pangdu@vt.edu

Katherine Ensor Rice University 6100 Main St. Houston, TX 77005 ensor@rice.edu

Zhide Fang Department of Mathematics University of New Orleans New Orleans, LA 70148 <u>zfang@uno.edu</u>

Jose Enrique Figueroa-Lopez Department of Mathematics Purdue University West Lafayette, IN 47906 figueroa@stat.purdue.edu

Helmut Finner German Diabetes Center Duesseldorf, D-40225 DE finner@ddz.uni-duesseldorf.de

Mujyanama B Frank P.O BOX 117 Butare Rwanda Butare, 250 RW <u>gifrank2003@yahoo.com</u> Daniel Fresen Department of Mathematics University of Missouri Columbia, MO 65211 <u>fresenj@MO.edu</u>

John L Fresen Department of Statistics University of Missouri Columbia, MO 65211 <u>fresenj@MO.edu</u>

Nancy L Glenn 1148 Rockwood Road Columbia, SC 29209 nglenn@math.sc.edu

Graciela Gonzalez-Farias Jalisco s/n Mineral de Valenciana Guanajuato, GTO 36240 <u>farias@cimat.mx</u>

Violeta Hennessey 1515 Holcombe Boulevard Houston, TX 77030 violeta.g.hennessey@uth.h mc.edu

Christian Houdre School of Mathematics GA Institute of Technology Atlanta, GA 30332 houdre@math.gatech.edu

Jianhua Hu U.T. MD Anderson Cancer Center 1515 Holcombe Blvd Houston, TX 77030 jhu@mdanderson.org

Jianhua Huang Blocker Building TX A&M University <u>xlhuang@mdanderson.org</u>

Peter Huber P. O. Box 198 CH-7250 Klosters Klosters, CH peterj.huber@bluewin.ch Adarsh Joshi 3902 College Main Apt 1411 Bryan, TX adarsh@stat.tamu.edu

David Kahle 4102 Ascot Lane Houston, TX 77092 <u>dkahle@rice.edu</u>

Yuliya Karpievitch U.T. MD Anderson Cancer Center 1515 Holcombe Blvd Houston, TX 77030 ykarpi@mdanderson.org

Demetrios Kazakos 3100 Cleburne St. Houston, TX 77004 kazakosd@tsu.edu

Jerome Paul Keating One UTSA Circle The University of Texas San Antonio, TX 78249 jerome.keating@utsa.edu

Sunduz Keles Department of Statistics University of Madison Madison, WI 53705 keles@stat.wisc.edu

Davar Khoshnevisan University of Utah Department of Mathematics Salt Lake City, UT 84112 davar@math.UT.edu

Farinaz Koushanfar 6100 S Main Rice ECE Dept. MS 366 Houston, TX 77005 farinaz@rice.edu

Kalimuthu Krishnamoorthy Dept of Mathematics Louisiana State University Lafayette, LA 70504 <u>krishna@LA.edu</u> S.N. Lahiri Statistics Department TX A&M University College Station, TX 77843 snlahiri@stat.tamu.edu

Mike Lecocke 14818 Personality San Antonio, TX 78248 <u>mlecocke@stmarytx.edu</u>

J. Jack Lee U.T. MD Anderson Cancer Center 1515 Holcombe Blvd. Houston, TX 77030 jjlee@mdanderson.org

Jong Soo Lee 6315 Forbes Ave Apt 905 Pittsburgh, PA 15217 jslee@stat.cmu.edu

Erich L Lehmann University of California Department of Statistics 367 Evans Hall # 3860 Berkeley, CA 94720 shaffer@stat.berkeley.edu

Charles Lewis Fordham University 441 East Fordham Road Bronx, NY 10458 clewis@fordham.edu

Hua Li 7900 Cambridge ST 16-1E Houston, TX 77054 <u>hua.li@uth.tmc.edu</u>

Xiaohu Li 6061 De Zavala RD Apt 1016 San Antonio, TX 78249 mathxhli@hotmail.com

Xihong Lin Statistics Department Harvard University 655 Huntington Avenue Building 2, Office 419 Boston, MA 2115 xlin@hsph.harvard.edu Bruce George Lindsay 422 Thomas Building University Park, PA 16802 bgl@psu.edu

Peng Liu Department of Statistics Iowa State University Ames, IA 50011 <u>PLIU@iastate.edu</u>

Shaw-Hwa Lo Statistics Department Columbia University New York, NY 10027 <u>slo@stat.columbia.edu</u>

Wei-Yin Loh Department of Statistics University of Wisconsin Madison, WI 53706 loh@stat.wisc.edu

Joseph F Lucke Center for Clinical Research & Evidence-based Medicine UTHSC Medical School Houston, TX 77030 Joseph.F.Lucke@uth.tmc.edu

Denka Grudeva Markova 1806 S. 8th Street Apt 252 Waco, TX 76706 <u>denka\_markova@baylor.edu</u>

Matthias Mathaes Department of Statistics Rice University Houston, TX 77005 matze@rice.edu

Monnie McGee Department of Statistics 3225 Daniel Ave P.O. Box 750332 Dallas, TX 75275 mmcgee@smu.edu

Robert Musheg Mnatsakanov P.O. Box 6330 Morgantown, WV 26506 rmnatsak@stat.wvu.edu MaryAnn Morgan-Cox One Bear Place #97140 Waco, TX 76798 MaryAnn Morgan-Cox@baylor.edu

Peter Mueller U.T. MD Anderson Cancer Center 1515 Holcombe Blvd Houston, TX 77030 pm@wotan.mdacc.tmc.edu

Yolanda Muñoz Maldonado Mathematics Department Michigan Technological University 1400 Townsend Drive Houghton, MI 49931 <u>ymunoz@mtu.edu</u>

Michael Naaman 7315 Brompton St. APT. 368B Houston, TX 77025 <u>mikenaaman@rice.edu</u>

Tuan S Nguyen Department of Statistics Rice University 6100 Main St Houston, TX 77005 kcstevone@gmail.com

Linda Njoh 1806 S. 8th street Apt 252 Waco, TX 76706 <u>linda\_njoh@baylor.edu</u>

Richard C Ott 3505 North 12th St. Apt #D8 Grand Junction, CO 81506 rott@mesastate.edu

Haiying Pang Department of Biostatistics MD Anderson Cancer Center Houston, TX 77030 haiying.pang@uth.tmc.edu

Galen Papkov Department of Statistics Rice University Houston, TX 77005 gpapkov@rice.edu Debashis Paul Department of Statistics University of California Davis, CA 95616 <u>debashis@wald.ucdavis.edu</u>

Claudia Pedroza U. T. School of Public Health RAS E831 Houston, TX 77030 claudia.pedroza@uth.tmc.edu

Victor Pérez-Abreu CIMAT Apdo. Postal 402 Guanajuato, Gto. 36000 MX pabreu@cimat.mx

Kenneth Pietz Michael E. DeBakey VA Medical Center (152) 2002 Holcombe Blvd.

Houston, TX 77030 kpietz@houston.rr.com

Darwin Poritz 968 Southern Pass Ct Houston, TX 77062 darwin.h.poritz@nasa.gov

Jessica Pruszynski 2425 S. 21st St. Apt. 104 Waco, TX 76706 jessica\_pruszynski@baylor.edu

Lindsay Renfro 18 Barksdale Ave. Waco, TX 76705 Lindsay\_Renfro@baylor.edu

Peter Richardson 5213 Shady Maple Dr. Kingwood, TX <u>peterr@bcm.tmc.edu</u>

Rudolf H Riedi Department of Statistics Rice University 6100 Main St Houston, TX 77005 riedi@rice.edu Kristin Ritz 5440 N. Braeswood Blvd. #956 Houston, TX 77096 kristinr@rice.edu

Javier Rojo Statistics Department Rice University Houston, TX 77005 jrojo@rice.edu

Gary Rosner Department of Biostatistics UT MD Anderson Cancer Center 1515 Holcombe Boulevard Houston, TX 77030 glrosner@mdanderson.org

Changxiang Rui 7019 River Elms San Antonio, TX cxrui@uark.edu

Brisa N Sanchez 1420 Washington Heights Ann Arbor, MI 48109 <u>brisa@umich.edu</u>

David W Scott Department of Statistics Rice University Houston, TX 77005 scottdw@stat.rice.edu

John Seaman 121 Neely Rd Hewitt, TX john w seaman@baylor.edu

Juliet P Shaffer University of California Dept. of Statistics 367 Evans Hall # 3860 Berkeley, CA 94704 shaffer@stat.berkeley.edu

Yongzhao Shao 650 First Avenue Fifth Floor, #538 New York, NY 10016 shaoy01@med.nyu.edu Chad A Shaw 1 Baylor Plaza Houston, TX 77030 cashaw@bcm.tmc.edu

Yu Shen Department of Biostatistics UT MD Anderson Cancer Center Houston, TX 77030 yshen@mdanderson.org

Stephen Stigler Statistics Dept University of Chicago Chicago, IL 60637 stigler@galton.uchicago.edu

Suhasini Subba Rao Department of Statistics TX A&M University College Station, TX 77843 <u>suhasini@stat.tamu.edu</u>

Gabor J Szekely The National Science Foundation 4201 Wilson Blvd Arlington, VA 22201 gszekely@nsf.gov

Ajit C Tamhane Department of Statistics 2008 Sheridan Road Evanston, IL 60208 ajit@iems.northwestern.edu

Young K Truong Department of Biostatistics The University of North Carolina Chapel Hill, NC 27599 truong@bios.unc.edu

Kam-Wah Tsui Department of Statistics University of Madison Madison, WI 53706 <u>kwtsui@stat.wisc.edu</u>

Willem R van Zwet Dept. of Mathematics University of Leiden P.O. Box 9512 Leiden 2341CB vanzwet@math.leidenuniv.NL Marco Vilela Department of Biostatistics UT MD Anderson Cancer Center Houston, TX 77030 mvilela@mdanderson.org

Wilhelmine von Turk 43 Del Mar Avenue Berkeley, CA 94708 wvonturk@sbcglobal.net

Jane-Ling Wang Department of Statistics University of California Davis, CA 95616 wang@wald.ucdavis.edu

Naisyin Wang Department of Statistics TX A&M University College Station, TX 77843 nwang@stat.tamu.edu

Xiaohui Wang 7715 N 4th Court McAllen, TX 78504 <u>xhwang@utpa.edu</u>

Feifei Wei 8170 33rd Ave. S. MS#21111R Bloomington, MN 55425 feifei.wei@healthpartners.com

Talithia Williams 301 Platt Boulevard Claremont, CA 91711 williams@math.hmc.edu

Xiaowei Wu Department of Statistics Rice University Houston, TX 77005 <u>xwwu@stat.rice.edu</u>

Yuewu Xu Fordham School of Business 113 West 60th Street New York, NY 10023 yuxu@fordham.edu Ya'acov Ritov 10 Rabi Benyamin Jerusalem 96306 IL yaacov@mscc.huji.ac.il

Jose-Miguel Yamal 5726 Portal Dr. Houston, TX 77096 jmy@odin.mdacc.tmc.edu Po Yang Dept. of Mathematical Sciences DePaul University Chicago, IL 60614 <u>pyang3@depaul.edu</u>

Chunming Zhang 1300 University Avenue University of Wisconsin Madison, WI 53706 <u>cmzhang@stat.wisc.edu</u> Qiang Zhao 601 University Drive San Marcos, TX 78666 <u>qiang.zhao@txstate.edu</u>

Hongxiao Zhu Department of Statistics Rice University Houston, TX 77005 <u>hxzhu@stat.rice.edu</u>

# Acknowledgement of referees ' services The efforts from the following referees are gratefully acknowledged

Miguel Arcones University of Binghamton Binghamton, NY

Moulinath Banerjee University of Michigan Ann Arbor, MI

Camelia Bejan Rice University Houston, TX

Prabir Burman University of California Davis, CA

Patrick S Carmack University of Texas Southwestern Medical Ctr

Rolando Cavazos Universidad Autónoma Agraria Antonio Narro

William F. Christensen Brigham Young University Provo, UT

Dennis D Cox Rice University Houston, TX

Kjell Doksum University of Wisconsin Madison, WI

José Enrique Figueroa Purdue University West Lafayette, IN

Helmut Finner Heinrich-Heine-University Düsseldorf, DE Chong Gu Purdue University Lafayette, IN

Rudy Guerra Rice University Houston, TX

Graciela González CIMAT México

Daniel Hernández CIMAT México

Jianhua Huang Texas A&M University College Station, TX

Ruben Juárez University of Hawaii Manoa, HI

Tomasz Kozubowski University of Nevada Reno, NV

Erich L. Lehmann University of California Berkeley, CA

Xiaoyan Leng Wake Forest University Winston-Salem, NC

Ming-Ying Leung University of Texas El Paso, TX

Regina Liu Rutgers University Piscataway, NJ Wei-Yin Loh University of Wisconsin Madison, WI

Robert M. Mnatsakanov. West Virginia University Morgantown, WV

William Navidi Colorado School of Mines Golden, CO

Anna Panorska University of Nevada Reno, NV

Arthur Pewsey University of Extremadura Spain

Peihua Qiu University of Wisconsin Madison, WI

Frits Ruymgaart Texas Tech University Lubbock, TX

Sanat Sarkar (2) Temple University Philadelphia, PA

Juliet P. Shaffer University of California Berkeley

Yongzhao Shao Iowa State University Ames. IA

Stephen Stigler University of Chicago Chicago, IL Young Troung Univ of North Carolina Chapel Hill, NC

Stephan R. Sain NCAR Boulder, CO

Nicola Sartori University of Padova Italy Alex Tsodikov University of Michigan Ann Arbor, MI

Jane-Ling Wang University of California Davis

Guosheng Yin MD Anderson Cancer Ctr Houston, TX Zhiliang Ying Columbia University New York, NY

Cun-Hui Zhang Rutgers University Piscataway, NJ

Hongyu Zhao Yale University New Haven, CT

# Erich L. Lehmann, The Lehmann Symposia, and November 20<sup>th</sup> 1917

Javier Rojo??

Rice University

The Lehmann Symposia originated as a result of a conversation I had in the year 2001 with the, then, Director of the Centro de Investigación en Matemáticas (CIMAT), Victor Pérez-Abreu. We both felt that there was an urgent need to bring back into focus theoretical statistics and our proposed solution was a series of Symposia that could serve as a forum for some of the exciting theoretical work being done in statistics. The First Lehmann Symposium took place at CIMAT in May of 2002. Most of the participants were Mexican colleagues. The program can be seen at the site http://www.stat.rice.edu/lehmann/1st-Lehmann.html. The second Lehmann Symposium – http://www.stat.rice.edu/lehmann/ – was held in May of 2004 at the School of Engineering at Rice University. Initially, the venues for the Symposia would alternate between CIMAT and Rice University. However, for various reasons, some being financial, it was decided to hold the 3<sup>rd</sup> Lehmann Symposium in the United States.

The original plans for the Third Lehmann Symposium were to hold the symposium at the Mathematical Sciences Research Institute (MSRI) in Berkeley during the month of November of 2007. The Third Symposium, however, ended up being held at Rice University for a second time during May of 2007. See http: //www.stat.rice.edu/~jrojo/3rd-Lehmann/. I co-edited webcasts of the Second and Third Symposia, and these webcasts are freely available to the public. They can be found at the following sites:

33 http://webcast.rice.edu/webcast.php?action=details&event=408 — second
34 symposium, and

<sup>35</sup> http://webcast.rice.edu/webcast.php?action=details&event=1057 — third symposium.

But why was the venue for the Third Symposium changed from California back to Texas, and why was the date changed from November  $20^{th}$ , 2007 to May  $16^{th}$ , 2007? There were very good reasons for holding the opening of the Symposium on Novem-ber  $20^{th}$ , 2007. For example, November  $20^{th}$ , 2007 was the silver anniversary of the greatest big game of all time. See, for example: http://www.alumni.berkeley. edu/KCAA\_Multimedia/The\_Play\_1982.asp. Another good reason to start the Sym-posium on November  $20^{th}$  was to co-celebrate, with our Mexican counterparts, the start of the first major  $20^{th}$  century revolution. The Mexican revolution started on November 20<sup>th</sup>, 1910 to remove the dictator Porfirio Díaz who had remained in power for 30 years. This revolution led to the Constitution of 1917 and the start of the Partido Revolucionario Institucional that held power until 2000 when a candi-date from the Partido Acción Nacional, Vicente Fox, won the Presidential election. Francisco I. Madero, with the help of Francisco Villa, took over from Porfirio Díaz.

<sup>??</sup>Department of Statistics, MS-138; Rice University; 6100 Main Street; Houston, TX 77005;

 $\mathbf{2}$ 



May 17 - The paper "On the distribution of the correlation coefficient in small sam-ple. Appendix I to the papers of "Student" and R. A. Fisher. A cooperative study". by H.E. Soper, A.W. Young, B.M. Cave, A. Lee and K. Pearson, (Biometrika 1917 11: 328-413; doi:10.1093/biomet/11.4.328), and the paper "I. Tables for estimating the Probability that the Mean of a unique Sample of Observations lies between  $-\infty$ and any given Distance of the Mean of the Population from which the Sample is drawn" by "Student" (Biometrika 1917 11: 414-417; doi:10.1093/biomet/11.4.414) are published. The former would include a criticism of Fisher's maximum likelihood principle that helped ignite a feud between Fisher and Pearson. 

May 18 - The Selective Service Act passes the U.S. Congress giving the President
 the power of conscription.

**July 4** - Petrograd Street demonstration - The Bolshevik revolution looms in the horizon.



Days of revolution - barricades at the Arcenal [i.e., Arsenal], Petrograd. The photo is taken from the Library of Congress under lot 2398, reproduction number LC-USZ62-25298, 1917.

**November 7** - Bolshevik Revolution begins: The workers of St. Petersburg in Russia, led by the Bolsheviks and the Bolshevik leader Vladimir Lenin, attacked the ineffective Kerensky Provisional Government.

**November 20** - Ukraine is declared a republic.

#### November 20, 1917 and Beyond

Amidst the shadows of war and civil unrest, a small burst of light began to shine in Strasbourg, France. Erich Leo Lehmann was born November  $20^{th}$ , 1917 – a mere 7 months after The United States entered the First World War.

Some years later, at the age of 16, he and his family went to live in Switzerland to avoid the Nazis. After five years in Switzerland, and two years in Cambridge, Erich L. Lehmann arrived in the United States in 1940 with a letter of introduction from the wife of Edmund Landau. Landau had passed away a couple of years earlier from a heart attack. The letter of introduction was for Richard Courant who was in New York and had been a colleague of Landau in Göttingen. After being asked by Courant if he wanted to stay in New York or live in the United States, Lehmann responded that he wanted to live in the United States and then followed Courant's

З

З



Erich Leo Lehmann in 1919

advice to go to an "up-and-coming university" in Berkeley, California. Upon his arrival in Berkeley in 1941, Erich Lehmann met with Griffith C. Evans who had been a mathematician at the Rice Institute, now Rice University. Evans had been brought to Berkeley to develop the mathematics department that was in disarray. Evans was an excellent mathematician and many of his contributions as a mathematician and administrator have been recorded in Morrey (1983) and Lehmann (2007). As an administrator, Evans was able to attract to Berkelev some of the best mathematicians of the time. With a broad vision for the mathematics department, Evans supported a three-week visit by R. A. Fisher to Berkeley. The visit did not go well. Reid (1992) writes in her book that, despite a generous endowed lectureship whose terms required the lecturer to spend their time on campus to interact with interested faculty, Fisher spent the first five days of his visit in San Francisco and went back to England a day earlier "standing up a dinner in his honor". Reid (1992) writes that according to Raymond T. Birge, chair of the Physics department at the time, "Fisher was the most conceited man he had ever met - 'and that is saying a lot with such competitors as Millikan et al!" Birge put forth Neyman's name to Evans. Evans had never heard of Neyman but after some inquiring an offer was

З



Erich Leo Lehmann in 1924

made. Neyman accepted and a few years later, after an offer from Columbia became available, was able to negotiate with Evans for the creation of a separate statistics department. These and other fascinating details may be found in Lehmann (1993, 1996, and 2007) and Reid (1992) and other references in the bibliography.

During their first meeting, Evans offered him a probationary graduate student status and six months later a teaching assistantship in the Mathematics Department. With the advent of the Second World War, Evans suggested to change areas of study and consider a more useful subject. Either Physics or Statistics would be more useful than Mathematics. After completing the required course work, and after returning from Guam where he and Joseph Hodges had served, it was time for Erich L. Lehmann to begin work on a dissertation. A topic with a probabilistic flavor was proposed by Pao-Lu Hsu after consulting with Neyman. Progress was swift and as Lehmann prepared to write up some of the results, a reference led to other references that led to the painful discovery that the results so far obtained had been published a few decades earlier. At that time Neyman was invited as a member of a delegation to observe the Greek elections. Concerned with the disappointment of his student, and knowing that he might return until a few months

later, Neyman asked Hsu to consider providing a new problem for Lehmann's thesis. Hsu suggested a new problem that he had thought about and planned to work on, and a problem for which he already had some preliminary results. Lehmann came to know about Hsu's generosity some time later and had hoped to thank Hsu per-sonally after Hsu's return to Berkeley from Columbia, but this would never happen as Hsu opted to return to China. With Neyman in Greece and Hsu back in China, Neyman suggested George Polya as a surrogate advisor. Frequent visits to Polya at Stanford finally yielded a thesis. A new problem presented itself in that Polya was not a faculty member at Berkeley. Fortunately, Neyman was able to return in time for the thesis defense. He had been asked to return to the United States as his services were no longer needed in Greece. In effect, he had been dismissed for insubordination. He had decided to investigate on his own the possibility that the elections had been rigged. It thus happened that in 1946 Erich L. Lehmann received the degree of Doctor of Philosophy. The title of his thesis: "Optimum Tests of a Certain Class of Hypotheses Specifying the Value of a Correlation Coefficient".

Erich L. Lehmann stayed in Berkeley as a young faculty member and the "Rest of the story" is well known. Besides his many influential publications, he was able to produce 41-plus Ph.D. students. The following table provides the names of the students and the year of graduation.

| Colin Ross Blyth                 | 1950 | Gouri Kanta Bhattacharyya     | 1966 |
|----------------------------------|------|-------------------------------|------|
| Fred Charles Andrews             | 1953 | James Nwoye Adichie           | 1966 |
| Allan Birnbaum                   | 1954 | Dattaprabhakar V. Gokhale     | 1966 |
| Hendrik Salomom Konijn           | 1954 | Frank Rudolph Hampel          | 1968 |
| Balkrishna V. Sukhatme           | 1955 | Wilhelmine von Turk Stefansky | 1969 |
| V. J. Chacko                     | 1959 | Louis Jaeckel                 | 1969 |
| Piotr Witold Mikulski            | 1961 | Friedrich Wilhelm Scholz      | 1971 |
| Madan Lal Puri                   | 1962 | Dan Anbar                     | 1971 |
| Krishen Lal Mehra                | 1962 | Michael Denis Stuart          | 1972 |
| Subha Bhuchongkul Sutchritpongsa | 1962 | Claude L. Guillier            | 1972 |
| Shishirkumar Jogdeo              | 1962 | Sherali Mavjibhai Makani      | 1972 |
| Peter J. Bickel                  | 1963 | Howard J. M. D'Abrera         | 1973 |
| Arnljot Høyland                  | 1963 | Hyun-Ju Yoo Jin               | 1974 |
| R. Murty Ponnapalli              | 1964 | Amy Poon Davis                | 1977 |
| Milan Kumar Gupta                | 1964 | Jan F. Bjørnstad              | 1978 |
| Madabhushi Raghavachari          | 1964 | William Paul Carmichael       | 1981 |
| Vida Greenberg                   | 1964 | David Draper                  | 1981 |
| Kjell Andreas Doksum             | 1965 | Wei-Yin Loh                   | 1982 |
| William Harvey Lawton            | 1965 | Marc J. Sobel                 | 1983 |
| Shulamith Gross                  | 1966 | Javier Rojo                   | 1984 |
| Bruce Hoadley                    | 1966 |                               |      |
|                                  |      |                               |      |

So why, then, was the date changed from November  $20^{th}$  to May  $16^{th}$ ? After all, it would have been a great way of celebrating Erich's wonderful 90 years of life. But that is precisely the issue. To the reader who does not know Erich L. Lehmann personally, holding a conference on his birthday, a conference that is named after him, would seem only natural. However, those close to him know very well that he is very modest and an event like the Symposium held on his birthday, would be rather uncomfortable for him. He thought that the meeting would turn into a birthday celebration and he would not have it that way. The Lehmann Symposia should be true to its beginnings: A meeting to showcase good theoretical work. Thus, it came to be that the venue and the date for the  $3^{rd}$  Lehmann Symposium were changed. As it was not possible to celebrate his  $90^{th}$  birthday with the symposium, this volume is dedicated to Erich's  $90^{th}$  birthday. I am sure that, given the opportunity.

| 1<br>2<br>3<br>4   | all his Ph.D. students listed above and colleagues around the world would join me<br>in wishing Erich Leo Lehmann many more wonderful years! Our lives have been<br>greatly enriched through our interactions, professional and social, with him.   | 1<br>2<br>3<br>4   |
|--|---|--|
| 5<br>6   | Acknowledgements  | 5<br>6   |
| 7<br>8<br>9<br>10  | The work of the author was partially supported by NSF Grant DMS-053246, NSA Grant H98230-06-1-0099, and NSF REU Grant DMS-0552590.  | 7<br>8<br>9<br>10  |
| 11   | References  | 11   |
| 11<br>12<br>13<br>14<br>15<br>16<br>17<br>18<br>19<br>20<br>21<br>22<br>23<br>24<br>25<br>26<br>27<br>28<br>29<br>30<br>31<br>32<br>33<br>34<br>35<br>36<br>37<br>38<br>39<br>40<br>41<br>42<br>43<br>44<br>45<br>44<br>45<br>44<br>45<br>44<br>45<br>44<br>45<br>44<br>45<br>44<br>45<br>44<br>45<br>46<br>47<br>46<br>47<br>47<br>47<br>47<br>47<br>47<br>47<br>47<br>47<br>47 | <ul> <li>KREPERCES</li> <li>SCR, REGINALD C. (1912). Mexico: Its Ancient and Modern Civilisation, History, Political Conditions, Topography, Natural Resources, Industries and General Development. Hume, N. (2). London, T. Fisher Unwin.</li> <li>DEGROOT, MORRIS H. (1986). A conversation with Erich L. Lehmann. Statist. Sci., 12, 243–254.</li> <li>LEMMANN, E. L. (1993). Mentors and early collaborators: reminiscences from the years 1940-1966 with an epilogue. Statist. Sci., 83, 331–341.</li> <li>LEMMANN, E. L. (1993). The treation and early history of the Berkeley statistics department. Statist. Sci., 12, 12, 83–331.</li> <li>LEMMANN, E. L. (1997). Testing statistical hypotheses: the story of a book. Statist. Sci., 12, 14, 8–52.</li> <li>LEMMANN, E. L. (2004). Optimality and symposia: some history. In The First Erich L. Lehmann Symposium - Optimality J. Rojo and V. Pérez-Abreu (Eds.) IMS LNNS, 304, 41, 1-01.</li> <li>LEMANN, E. L. (2007). Reminiscences of a Statisticia: the Company I Kept. Springer, New York.</li> <li>Mexer, C. E. (1983). Ortfifth Conrad Evans 1887–1973: A Biographical Memoir. National Academy of Sciences, Washington D.C.</li> <li>Rein, C. (1982). Neyman from Life. Springer-Verlag, Berlin-Heidelberg-New York-Tokyo.</li> </ul> | 11<br>12<br>13<br>14<br>15<br>16<br>17<br>18<br>19<br>20<br>21<br>22<br>23<br>24<br>25<br>26<br>27<br>28<br>29<br>30<br>31<br>32<br>33<br>34<br>35<br>36<br>37<br>38<br>39<br>40<br>41<br>42<br>43<br>44<br>45<br>44<br>45<br>46<br>47<br>46<br>47<br>46<br>47<br>46<br>47<br>46<br>47<br>46<br>47<br>46<br>47<br>46<br>47<br>46<br>47<br>47<br>47<br>47<br>47<br>47<br>47<br>47<br>47<br>47 |
| 17<br>18<br>19<br>50   |   | 47<br>48<br>49<br>50   |
| 51   |   | 51   |

# The Honorable Erich L. Lehmann

#### Stephen Stigler

University of Chicago

З

The year 2007 marks a concurrence of important statistical anniversaries. It is the  $350^{th}$  anniversary of the publication of the first printed work on mathematical probability, the short tract that Christian Huygens wrote following a visit to Paris, where he learned of the investigations of Fermat and Pascal. Also, 2007 is the  $150^{th}$ year since the birth of Karl Pearson, father of the Chi-square test and much else. And related to both those events, it is also the year our teacher, friend, and colleague Erich Lehmann celebrates his  $90^{th}$  birthday. Christian Huygen's tract served as the textbook on probability for over a half century, helping to form that subject. Karl Pearson inaugurated an important species of hypothesis testing. Both then have important similarities to Erich Lehmann. But rather than further explore those analogies immediately, I would like to characterize an important part of Erich's ongoing research by looking back to a more modern document. 

The University of Chicago, rare among research universities, gives honorary degrees only in recognition of scholarly contributions of the highest order. We do not use those degrees to honor movie stars, philanthropists, or even heads of state (at least not over the past 80 years). There is a partial exception: we do so honor the departing Chair of our Board of Trustees. But that is the limit of the exceptions. We do not use this device to honor work done at Chicago; our major financial sup-porters are recognized in other ways; and discreet inquiries on behalf of politicians, celebrities, popular artists, and several heads of state have been politely turned aside, however meritorious they may have been on other grounds. Scholarship is the only coin of our realm. 

One of the fields where this practice has been actively pursued is statistics. The first degree our newly formed department granted was an honorary Doctorate of Science to Ronald Fisher, June 13, 1952. This was followed over the next 31 years by degrees to Harold Hotelling (1955), Jerzy Neyman (1959), Maurice Bartlett (1966), John Tukey (1969), and Fred Mosteller (1973). In 1990, in anticipation of our University's centennial celebration beginning the following year, the Department undertook to resume this practice after several years by proposing an honorary degree for Erich Lehmann. 

The procedure for granting honorary degrees at Chicago is a bit involved. After getting departmental agreement (an easy matter in this case) it is necessary to prepare a detailed case to be submitted to a cross-university committee appointed for this task. Exacting standards are upheld and only a fraction of proposals are given the nod of approval. Local legend has it that when someone proposed the Queen of England for a degree, it was turned back with a request for her list of publications. The procedure is like that for hiring a senior scholar from outside the university: several letters of recommendation must be solicited, the evidence must be assembled and carefully presented, all with an uncertain outcome. And unlike senior hires, even if the proposal is successful and the offer accepted, you do not get to keep the candidate! 

imsart-coll ver. 2008/08/29 file: Stigler.tex date: March 25, 2009

In preparing Erich's case we solicited letters from a dozen of the top international statistical scholars over the past half-century. All of these were glowing testaments to an amazingly influential body of work, as well as to Erich's collegial role in helping to build modern mathematical statistics. To buttress the case we did a careful citation study to document quantitatively the pervasive influence of Erich's articles and books. To our eye there was no doubt, but all this needed to be presented to a demanding committee of some of the universitys best professors, regrettably none with more than a superficial knowledge of our subject. It is that memo which I will now present.

\* \* \*

23 April 1990

#### Memo to Committee on Honorary Degrees

On behalf of the Department of Statistics, and with their unanimous endorsement, I wish to nominate Erich L. Lehmann for an Honorary Degree. We believe he would be an exceptionally appropriate candidate for a Centennial Honorary Degree at the convocation of October 1991; alternatively he could be considered for an earlier convocation, such as June 1991. Prior to his retirement in 1988, Lehmann was Professor of Statistics at the University of California at Berkeley, where he remains extremely active.

Great scientists can be roughly classified as one of two types: those whose pri-mary achievement was a single brilliant discovery, and those who have over the course of a life's work constructed a discipline, a school. Lehmann has made many important discoveries, but he is more of the second type than the first. In the years after the Second World War, the dominant paradigm in American mathematical statistics was the decision theoretic school of Jerzy Neyman (Sc.D., Chicago, 1959) and Abraham Wald. Wald died tragically in 1950 and Neyman's attention was absorbed by other matters after the mid-1950's. While many hands played a role in the prospering of this approach over 1950-1980, it is arguable that the chief architect of the expansion of the paradigm, and the one most responsible for its immense international influence over that period, was Erich Lehmann.

There is an interesting historical parallel to Lehmann's role in this school of thought. In the 1650's Pascal and Fermat founded probability theory. Yet they taught few and published little. It was Christian Huygens's 1657 tract that was responsible for the form in which these ideas were disseminated for half a century, and their widespread application in mathematics, philosophy, and science. Similarly, it was Lehmann's mimeographed lecture notes of estimation (1950, widely circulated and reproduced, but only published as a book in 1983) and Lehmann's book of hypothesis testing (1959) that provided the form in which the Neyman-Wald ap-proach came to dominate a major portion of the mathematical world. Not all of the results are Lehmann's (though many are), but the arrangement, the elegant seamless presentation, the coherency of the whole, are his in a way that has not been true in many sciences. To a large degree, Erich Lehmann created the curricu-lum of the world's graduate programs in mathematical statistics over the period 1960-1980. Lehmann has personally supervised over 50 Ph.D. dissertations, and he counts among his students some of the most influential statisticians in the United States, Europe, and Asia.

The focus of the school that is so associated with Lehmann is the application

imsart-coll ver. 2008/08/29 file: Stigler.tex date: March 25, 2009

З

of decision theory to statistical problems, the construction of a calculus of optimal statistical procedures. In one simple setting, parametric estimation, an objective criterion ("a loss function") is defined and an attempt is made within a framework of stochastic models to seek the best procedure, or at any rate to determine an order for the available procedures. The pitfalls of this approach are many: In all but the simplest problems, the mathematical difficulties can be immense, and the specification of models and objective criteria that capture the essence of the scien-tific problem is rarely straightforward. The success of Lehmann and his school has been due to the balance they have maintained in creating a system of mathematical structures of sufficient richness to encompass a large range of practical problems, vet not so amorphously vague that a true discipline could not be constructed around them. Mathematics flourishes in the detailed exploration of constrained spaces with widely accepted rules; statistics flourishes with the flexibility to treat the infinite variety of problems in the real world. Lehmann's genius has been his ability to reconcile these divergent goals and build a school that has enriched both sides. \* \* \* To this was appended the letters and the citation study. The proposal was successful; we were delighted when Erich accepted and both he and Julie attended the special Convocation on October 3, 1991, celebrating the University's Centennial. Their visit was the occasion for several parties, dinners, toasts, and culminated in a grand ceremony in Rockefeller Chapel [yes, we do recognize donors in other ways], attended by dignitaries representing the great universities of the world. The citation on Erich's degree read: Your research on the application of decision theory to statistical problems has helped create and organize modern mathematical statistics; your elegant treatises have guided the curricula of a majority of the nation's graduate programs and given shape to the discipline, and your teaching has inspired a generation of scholars. Erich has received many other honors, of course, including an earlier honorary degree at the University of Leiden. And we have since 1991 given more honorary degrees; to Charles Stein (1992), Ulf Grenander (1994), Bradley Efron (1995), David Aldous (2000), Persi Diaconis (2003), and Grace Wahba (2007). The occasion of Erich's University of Chicago degree retains a special place in our memories. So too, to those of us who took his classes (in my case some 40 years ago) does the memory of his careful and supremely clear lectures, and his papers and books that helped shape the modern statistical world. For all that and more, we may toast this extraordinary scholar on the  $90^{th}$  anniversary of his birth! 

**IMS Collections** Vol. 0 (2009) 11-17 © Institute of Mathematical Statistics, 2009 arXiv: math.PR/000009 З Some History of Optimality Erich L. Lehmann University of California, Berkeley Contents Combination of Observations  $\mathbf{2}$ Maximum Likelihood Estimation The Neyman-Pearson Theory of Hypothesis Testing  $\mathbf{5}$ Wald's Optimality Criteria The Hunt-Stein Theorem Large Sample Optimality of Testing 10 A Culture Clash 1. Combination of Observations

З

#### The earliest optimality considerations appear to be those of Laplace and Gauss at the beginning of the 19<sup>th</sup> Century concerned with determining the best linear estimates of parameters in linear structures. Laplace calls these optimal estimates "the most advantageous" while Gauss refers to them as "the most plausible values." Various aspects of this problem were discussed throughout the 19<sup>th</sup> Century under the heading "Combination of Observations." The version of the principal result generally accepted today is the so-called Gauss-Markov theorem. It states (in modern language) that the least squares estimates are the linear unbiased estimates with minimum variance. While restricted to linear estimates, the result is nonparametric in that it makes no assumptions about the distribution of the errors. For an account of this work see, for example, Stigler (1986), Hald (1998), or Chatterjee (2003).

### 2. Maximum Likelihood Estimation

Optimality next played an important part in Fisher's fundamental paper of 1922. In this paper (followed by a clarifying paper in 1925), Fisher considers estimation in parametric models, proposes the maximum likelihood estimator (MLE) as a generally applicable solution, and claims (but does not prove) that such estimators are consistent and asymptotically efficient efficient (i.e., that they minimize the asymptotic variance). Note that unlike the Gauss-Markov theorem, maximum likelihood estimation does assume that the distribution of the variables belongs to a given parametric family. Maximum likelihood has become the most widely used method of estimation, and there has been an enormous amount of work connected with Fisher's claims concerning it. It has led to the discovery of superefficiency on the one hand and of second order efficiency on the other. Many counterexamples have been found (even to consistency), and a number of modifications (bias correction and replacement of the MLE by a consistent root of the likelihood equation, for example) have been proposed. The situation is complex, but under suitable restrictions Fisher's conjecture is essentially correct when the likelihood equation has a unique root. For a more precise statement, see Shao (1999), and for further discussion, for example, Efron (1982), Le Cam (1990), and Barndorff-Nielsen and Cox (1994).

#### 3. The Neyman-Pearson Program

Least squares and maximum likelihood were first proposed on intuitive grounds and then justified by showing that they possessed certain optimality properties. Optimality as a deliberate program for determining good procedures was introduced in 1933 by Neyman and Pearson in a paper (on testing rather than estimation) appropriately called, "On the problem of the most efficient tests of statistical hypotheses." As they explain in the introduction:

The main purpose of the present paper is to find a general method of determining tests which ... would be most efficient" [in the sense of minimizing the probability of erroneous conclusions].

In a certain sense this is the true start of optimality theory.

З

#### 4. The Neyman-Pearson Theory of Hypothesis Testing

Neyman and Pearson (1933) implemented the above program by seeking, for any given situation, the test which, among all those controlling the probability of false rejection at a given level  $\alpha$ , has the maximum power (and hence the minimum probability of false acceptance).

For testing a simple hypothesis against a simple alternative, they found the solution to this problem to be the likelihood ratio test. This result, which is mathematically quite elementary but has crucial statistical consequences, is known as the Neyman-Pearson Lemma.

It turns out that in some (very rare) cases the same test is most powerful against all alternatives under consideration. Such a uniformly most powerful (UMP) test is then the optimal solution to the given testing problem. Where a UMP test does not exist, additional criteria must be invoked.

For example, when nuisance parameters are present, Neyman and Pearson require that under the hypothesis the rejection probability be  $\alpha$  for all values of the nuisance parameters. They call such rejection regions <u>similar regions</u>. As an important example, they show that the one-sided t-test is UMP among all similar regions.

For two-sided alternatives, one would not expect UMP tests to exist even without nuisance parameters. For such cases, Neyman and Pearson then impose the additional condition of <u>unbiasedness</u>, i.e., that the power of the test is  $\geq \alpha$  for all alternatives. In follow-ups to their 1933 paper (1936 and 1938), they show, for example, that the two-sided t-test is UMP among all unbiased tests. UMP similar or unbiased tests exist for important classes of testing problems concerning a single real-valued parameter (in the presence or not of nuisance parameters) but not for hypotheses such as

$$H: \theta_1 = \cdots = \theta_s$$

concerning several parameters.

A different condition, the <u>principle of invariance</u> (suggested by Hunt and Stein, unpublished), is successful in a number of important such multiparameter situations. If both the hypothesis and the class of alternatives remain invariant under a group G of transformations of the sample space, there does in these cases exist a UMP test among all tests invariant under G.

#### 5. Wald's Optimality Criteria

A quite-different approach to optimality was initiated by Wald in his 1939 paper, "Contributions to the theory of statistical estimation and testing hypotheses," and was then developed further in a series of publications culminating in his 1950 book, "Statistical Decision Functions."

This approach was part of Wald's formulation of a general theory of decision procedures. Instead of seeking procedures that are uniformly optimal among some suitably restricted class of decision procedures, Wald proposes to minimize some global feature of the performance.

Specialized to hypothesis testing, these proposals reduce to:

- i Maximize the average power, averaged with respect to some suitable weight function over the alternatives. For obvious reasons, Wald called such maximizing procedures Bayes solutions.
  - ii Maximize the minimum power over alternatives bounded away from the hypothesis.

#### 6. The Hunt-Stein Theorem

An important connection between the earlier invariance approach and that of Wald is the Hunt-Stein theorem. It states that if a UMP invariant test exists under a group G satisfying certain conditions (called amenability), then this test also maximizes the minimum power over any invariant set of alternatives.<sup>1</sup>

To illustrate this theorem, consider the univariate linear hypothesis. Under the combined action of the groups of location, scale and orthogonal transformation, a UMP invariant test exists. Since these groups satisfy the Hunt-Stein conditions, the resulting test therefore maximizes the minimum power against all invariant sets of alternatives.

As a second example, consider the multivariate one-sample problem. Hotelling's  $T^2$ -test is UMP among all tests that are invariant under the group G of all nonsingular linear transformations. Since G is not amenable, the Hunt-Stein theorem does not apply. The resulting maximin problem poses considerable difficulties. З

<sup>&</sup>lt;sup>49</sup> <sup>1</sup>The 1946 paper by Hunt and Stein, "Most stringent tests of statistical hypotheses," containing
<sup>50</sup> this theorem was never published. The theorem appeared in print for the first time in Lehmann,
<sup>51</sup> "Testing Statistical Hypotheses," John Wiley, 1959.
#### $E.L.\ Lehmann$

| 1<br>2   | 7. Some Extension of the Neyman-Pearson Theory   | 1<br>2   |
|----------|--|----------|
| 3        | The Neyman-Pearson theory has been extended in a number of directions. The             | 3        |
| 4        | following are two extensions of the Neyman-Pearson Lemma, which is so basic to         | 4        |
| 5        | this theory.   | 5        |
| 6        |  | 6        |
| 7        | (i) Sequential Analysis  | 7        |
| 8        |  | 8        |
| 9        | During World War II, Wald proposed sequential procedures, as a way of obtaining        | 9        |
| 10       | good power with fewer observations. In particular, he suggested the probability ratio  | 10       |
| 11       | test for testing a simple hypothesis against a simple alternative. This test continues | 11       |
| 12       | observation as long as the likelihood ratio remains between two fixed limits and       | 12       |
| 13       | takes the indicated decision (acceptance or rejection) as soon as it falls outside     | 13       |
| 14       | these limits.  | 14       |
| 15       | In 1948, Wald and Wolfowitz proved the remarkable result that for testing a            | 15       |
| 16       | simple hypothesis against a simple alternative, the sequential probability ratio test, | 16       |
| 10       | among all tests with the same (or smaller) probabilities of error, minimizes the       | 10       |
| 10       | expected number of observations both under the hypothesis and the alternative.         | 10       |
| 20       |  | 20       |
| 20       | (11) <u>Robust Inference</u>   | 20       |
| 21       | All the ment of a few (ment few the Course Menters theorem) are considered             | 21       |
| 23       | All the work reported so far (except for the Gauss-Markov theorem) was carried         | 23       |
| 24       | out under the assumption of an underlying parametric model. In practice, such an       | 24       |
| 25       | mulation Huber (1064) suggested replacing the assumption of a parametric model         | 25       |
| 26       | by that of a neighborhood of such a model  | 26       |
| 27       | In the following year, he obtained the analog of the Neyman-Pearson Lemma              | 27       |
| 28       | For testing the neighborhood of a distribution $P_0$ the test maximizing the mini-     | 28       |
| 29       | mum power over the neighborhood of an alternative $P_1$ is a censored version of the   | 29       |
| 30       | likelihood ratio test of $P_0$ against $P_1$ .   | 30       |
| 31       |  | 31       |
| 32       | (iii) Multiple Testing   | 32       |
| 33       |  | 33       |
| 34       | A very different extension of the Neyman-Pearson theory that is of great practical     | 34       |
| 35       | importance deals with the situation in which a number of hypotheses (sometimes a       | 35       |
| 36       | very large number) are being tested rather than just one. Unlike the work discussed    | 36       |
| 37       | so far, which is classical, the theory of multiple testing is an area of very active   | 37       |
| 38       | ongoing research.  | 38       |
| 39       | The first problem here (before optimization) is to find a suitable generalization      | 39       |
| 40       | of the concept of significance level that would provide satisfactory control of the    | 40       |
| 41       | probability of false rejections. After this, maximin tests have been obtained, which   | 41       |
| 42<br>43 | require, however, not only unbiasedness and invariance but also a condition of         | 42       |
| 40       | monotonicity. Surveys of the current state of this work are provided by Shaffer        | 43<br>47 |
| 45       | (2004, 2006).  | 44<br>45 |
| ±0<br>46 |  | 40       |
| 47       | 8. Large Sample Optimality of Testing  | 47       |
| 48       | General Francisco Account  | 48       |
| 49       | A small-sample theory of optimal estimation parallels that of optimal testing          | 49       |
| 50       | sketched in Sections 4-7, with concepts such as unbiasedness, equivariance (instead    | 50       |
| 51       | of invariance), and minimax variance (instead of maximin power), and will not be       | 51       |
|          |  |          |

imsart-coll ver. 2008/08/29 file: Lehmann\_hist\_opt.tex date: March 25, 2009

discussed here. Asymptotic optimality for estimation goes back to Fisher (1922),
as mentioned in Section 2, and is defined as minimum asymptotic variance.

For testing, asymptotic optimality is considerably more complex, both conceptually and technically. It was first studied by Wald in 1941. Consider testing a simple hypothesis  $\theta = \theta_0$  against a simple alternative  $\theta = \theta_1$ . If we keep both  $\theta_0$ and  $\theta_1$  fixed, and carry out the tests at a fixed level  $\alpha$ , the power of any reasonable test sequence will tend to 1. Thus any such test sequence will in a trivial sense be asymptotically UMP.

A more useful approach is obtained by considering a sequence of alternatives

(8.1) 
$$\theta_n = \theta_0 + h/\sqrt{n}$$

For a sequence with fixed h, the power will typically tend to a limit between 0 and 1 as  $n \to \infty$ . As h varies from 0 to  $\infty$ , the limiting power will be an increasing function of h, going from  $\alpha$  to 1; we shall call this the asymptotic power function. A sequence of tests can then be defined as being asymptotically most powerful (AUMP) if it maximizes the asymptotic power for all h.

Unlike the finite sample situation where UMP tests exist only rarely and then are unique, it turns out that AUMP tests exist under very weak assumptions, and that in fact many different AUMP tests exist for the same situation, among them the likelihood ratio test, the Wald test, and the locally most powerful (Rao) test. To distinguish them, one must resort to higher order asymptotics. (See, for example, Barndorff-Nielsen and Cox (1994)). An exposition of the first order theory due to Le Cam can be found in Lehmann and Romano (2005).

# 9. Optimal Design

In generalization of minimizing the variance of an estimator, optimal design theory is concerned with determining the design that minimizes some function of the covariance matrix of the best linear estimators of the parameters in question. In particular, D-optimality minimizes the determinant, E-optimality is a minimax criterion, and so on. After some isolated earlier efforts, the problem of optimal design was studied systematically by Kiefer in more than 40 papers between 1958 and his early death in 1981. They make up vol. III of his collected papers.

### 10. A Culture Clash

Not everyone was enthusiastic about optimality as a guiding principle. Criticism was highlighted at a 1958 meeting of the Royal Statistical Society at which Kiefer presented a survey talk on "Optimum Experimental Designs." It was a discussion paper, and the reaction was nearly uniformly negative. The core of the disagreement is stated clearly when Barnard quotes Kiefer as saying of procedures proposed in a paper by Box and Wilson (1951) that they "often [are] not even well-defined rules of operation." Barnard's reply:

In the field of practical human activity, rules of operation which are not well defined may be preferable to rules which are.

The conflict is discussed by Henry Wynn in his introduction to a reprint of
Kiefers paper in "Breakthroughs in Statistics," vol. I (Kotz and Johnson, Eds.).
Wynn calls it "a clash of statistical cultures."

#### E.L. Lehmann

This clash is between the Wald school of abstraction and optimality on the one hand and what Tocher, in his discussion of Kiefer's paper, calls "the school of British experimental design – a very practical people," on the other.

#### 11. Tukey's Criticism

Criticism of optimality was not confined to England. An outspoken American critic questioning the value and importance of optimization was John Tukey. His attitude is indicated by the titles of two philosophical papers in 1961 and 1962, which are titled respectively, "The tyranny of the best" and "Dangers of optimization."

Tukey's concern with optimality had its origin in the fact that at the time optimization had become the dominant interest of mathematical statistics. In the 1962 paper, he writes:

Some [statisticians] seem to equate [optimization] to statistics an attitude which, if widely adopted, is guaranteed to produce a dried-up, encysted field with little chance of real growth.

#### 

З

# 12. Conclusion

That optimality considerations had become so dominant is explained by the historical situation. Periods of great innovation are followed by periods of consolidation, in which the new work is given a final (or, as Tukey says, encysted) form. Such a dichotomy is also discussed by Huber (1975). Thus, the revolutionary work of Student and Fisher was followed by the optimization approach of Neyman, Pearson, Wald and their students described in this paper.

Today we are faced with the opposite situation. We are confronted with new problems arising from immense data sets and involving problems of great complexity. Ad hoc solutions are proposed and tried out on a few examples. This is a natural first step, but eventually we will want to justify the solutions at which we have arrived intuitively and by trial and error. A theoretical underpinning will be provided and the conditions will be found under which these solutions are the best possible.

#### References

- BARNDORFF-NIELSEN, O. E. and Cox, D. R. (1994). Inference and Asymptotics. Chapman & Hall, London.
   BOX, G. E. P. and WILSON, K. B. (1951). On the experimental attainment of optimum conditions. Journal of the Royal Statistical Society (B), 13, 1–45.
   CHATTERJEE, S. K. (2003). Statistical Thought. Oxford University Press.
- [4] EFRON, B. (1982). Maximum likelihood and decision theory. Annals of Statistics, 10, 309–368.
- [4] [4] [5] FISHER, R. A. (1922). On the mathematical foundations of theoretical statistics. *Phil. Trans. Roy.* 42 Soc. London (A), 222, 401–415.
- 43 [6] FISHER, R. A. (1925). Theory of statistical estimation. Cambridge Philos. Soc., 22, 700–725.
- [7] HALD, A. (1998). A History of Mathematical Statistics from 1750 to 1930. John Wiley, New York.
  - [8] HUBER, P. (1964). Robust estimation of a location parameter. Annals of Mathematical Statistics, 35, 73-101.
- 46 35, 73-101.
  47 [9] HUBER, P. (1965). A robust version of the probability ratio test. Annals of Mathematical Statistics, 36, 1753-1758.
- [10] HUBER, P. (1975). Applications vs. abstraction: the selling out of mathematical statistics? In Proc.
   Conference on Directions of Mathematical Statistics. Suppl. Adv. Prob., 7, 84–89.
- [11] KIEFER, J. C. (1959). Optimum experimental designs (with discussion). Journal of the Royal Statistical Society (B), 21,273–319. (Reprinted in Kotz and Johnson (1992). Breakthroughs in Statistics, Vol. 1., Springer-Verlag.

# Some History of Optimality

| 1        | [12]<br>[13] | KIMBALL, G. E. (1958). A critique of operations research. J. Wash. Acad. Sci., 48, 33–37. KOTZ, S. and JOHNSON, N. L. (Eds. 1992, 1997). Breakthroughs in Statistics, 1. Springer-Verlag,                   | 1  |
|----------|--------------|---|----|
| 3        | [1.4]        | New York.   | 3  |
| 4        | [14]         | Bayes' estimates. Univ. of Calif. Publ. in Statist., 1, 277–330.  | 4  |
| 5        | [15]         | LE CAM, L. (1990). Maximum likelihood an introduction. ISI Review, 58, 153–171.   | 5  |
| 6        | [16]         | LEHMANN, E. L. and ROMANO, J. P. (2005). <i>Testing Statistical Hypotheses</i> (3rd Ed.). Springer, New York  | 6  |
| 7        | [17]         | NEYMAN, J. and PEARSON, E. S. (1933). On the problem of the most efficient tests of statistical   | 7  |
| 8        | [10]         | hypotheses. Phil. Trans. Roy. Soc. (A), 231, 289–337.   | 8  |
| 9        | [10]         | hypotheses. Statist. Res. Memoirs, 1, 1–37; 2, 25–57.   | 9  |
| 10       | [19]         | PEARSON, E. S. (1939). "Student" as a statistician. Biometrika, <b>30</b> , 210–250.  | 10 |
| 11       | [20]         | SHAFFER, J. P. (2004). Optimality results in multiple hypothesis testing. In <i>Proceedings of the First</i><br>Lehmann Symposium, J. Rojo and V. Pérez-Abreu (Eds), IMS LNMS, 44, 11–35.                   | 11 |
| 12       | [21]         | SHAFFER, J. P. (2006). Recent developments towards optimality in multiple hypothesis testing. In  | 12 |
| 13       | [22]         | Proceedings of the Second Lehmann Symposium. J. Rojo (Ed). IMS LNMS, <b>49</b> , 16–32.   | 13 |
| 14       | [23]         | STIGLER, S. (1989). The History of Statistics. Harvard University Press, Cambridge, MA.   | 14 |
| 15       | [24]         | TUKEY, J. W. (1961). Statistical and quantitative methodology. In <i>Trends in Social Science</i> , (D.   | 15 |
| 16       | [25]         | P. Ray, Ed.). Philosophical Library, New York. (Reprinted in Vol. III of Tukey's Collected Works.)<br>TUKEY J. W. (1962). The future of data analysis. <i>Annals of Mathematical Statistics</i> . Reprinted | 16 |
| 10       | [=0]         | in Vol. III of Tukey's Collected Works.   | 10 |
| 10       | [26]         | WALD, A. (1943). Tests of statistical hypotheses concerning several parameters when the number  | 10 |
| 20       | [27]         | WALD, A. (1950). Statistical Decision Functions. John Wiley, New York.  | 20 |
| 20       | [28]         | WALD, A. and WOLFOWITZ, J. (1948). Optimum character of the sequential probability ratio test.  | 20 |
| 22       |              | Annals of Mathematical Statistics, <b>19</b> , 326–339.   | 22 |
| 23       |              |   | 23 |
| 24       |              |   | 24 |
| 25       |              |   | 25 |
| 26       |              |   | 26 |
| 27       |              |   | 27 |
| 28       |              |   | 28 |
| 29       |              |   | 29 |
| 30       |              |   | 30 |
| 31       |              |   | 31 |
| 32       |              |   | 32 |
| 33       |              |   | 33 |
| 34       |              |   | 34 |
| 35       |              |   | 35 |
| 36       |              |   | 36 |
| 37       |              |   | 37 |
| 30       |              |   | 30 |
| 40       |              |   | 40 |
| 10<br>41 |              |   | 40 |
| 42       |              |   | 42 |
| 43       |              |   | 43 |
| 44       |              |   | 44 |
| 45       |              |   | 45 |
| 46       |              |   | 46 |
| 47       |              |   | 47 |
| 48       |              |   | 48 |
| 49       |              |   | 49 |
| 50       |              |   | 50 |
| 51       |              |   | 51 |

# arXiv: math.PR/0000001

З

# An Optimality Property of Bayes' Test **Statistics**

# Raghu Raj Bahadur and Peter J. Bickel<sup>1,\*</sup>

University of Chicago and Imperial College, London

# Dedicated to Erich Lehmann on his 90<sup>th</sup> Birthday

| Introduction        |                    |       | <br> | <br> |  |  |  |
|---------------------|--------------------|-------|------|------|--|--|--|
| A Generalization of | a Theorem of Ba    | hadur | <br> | <br> |  |  |  |
| General Assumption  | ns and a Useful Le | emma  | <br> | <br> |  |  |  |
| The Main Theorem    |                    |       | <br> | <br> |  |  |  |
| Optimality of Minin | max Tests          |       | <br> | <br> |  |  |  |

### **Preface:**

This paper dates back to the late 60's when I collaborated with Raj Bahadur, who is unfortunately no longer with us. The reason it has not appeared until now is that he felt it had to be accompanied by a number of multivariate examples. We both went on to other things; the examples were not worked out although we both knew of the existence of some of them. So why is this paper appearing here (with the approval of Steve Stigler, an executor of the Bahadur estate)? First, in addition to attesting to Erich's continued vital presence, it gives me the opportunity of paying a tribute to Bahadur, who was a friend of both of ours. Second, it is an interesting reminder of how writing styles have changed on the whole I think for the better from rigorous abstract formulation and mathematically rigorous presentation to more motivation and a lot of hand waving. Third, and most importantly, the result is an example of what I think both Erich and I consider an important endeavor, the reconciliation of the Bayesian and frequentist points of view (in context of now rather unfamiliar asymptotics). In an important paper in the  $5^{th}$  Berkelev Sympo-sium [4], Bahadur showed that the maximum likelihood ratio statistic possessed an optimality property from the view of a large deviation based frequentist compari-son of tests he introduced in 1960 [1]. Our paper shows that this property is shared by Bayes test statistics for reasonable priors and conjectures that a corresponding Bayesian optimality property holds for the maximum likelihood ratio statistic. If true this can be viewed as the large deviation analogue of the well-known Bern-stein von Mises' theorem – see Lehmann and Casella [10] p.489, which establishes the equivalence at the  $n^{-\frac{1}{2}}$  scale of Bayesian and maximum likelihood estimates. Establishing this conjecture is left as a challenge to the reader. 

# <sup>1</sup>On leave from University of California, Berkeley (1965-66).

\*Prepared with partial support of N.S.F. Grant G.P. 2593.

Given the historical interest I have not changed the text save for typos and only brought references up to date.

# 1. Introduction

 In [4] one of the authors established the optimality of the classical likelihood ratio test statistic in terms of a method of stochastic comparison previously introduced by him in [1], [2], and [3].

In the main theorem of this paper, Theorem 2 of Section 4, we show that this property is shared by Bayes test statistics (averages of likelihood ratios with respect to probability measures on the parameter space) under conditions which are slightly different from, and in some respects weaker than those given in [4]. These assumptions are given and discussed in Section 3. Section 5 contains a theorem establishing the asymptotic optimality of minimax tests under appropriate restrictions.

In Section 2 we give a strengthening Theorem 1 of [4], which established a lower bound for the slope of any family of tests in terms of the Kullback-Leibler information numbers. The proof given here drops assumption 1 of [4] and weakens assumption 2 considerably. This argument seems to give some insight into the necessity of an assumption such as our modification of assumption 2 of [4].

#### 2. A Generalization of a Theorem of Bahadur

Even as in [4] we let X be an abstract space,  $\mathcal{A}$  a field on X,  $P_{\theta}$ ,  $\theta \in \Theta$ , a set of probability measures on  $(X, \mathcal{A})$ , and  $\Theta_0$  a given subset of  $\Theta$ . For any  $\theta$ ,  $\theta'$  we define,

(2.1) 
$$K(\theta, \theta') = -\int_X \log \frac{dP_{\theta'}}{dP_{\theta}}(x)dP_{\theta}(x) , \qquad 26$$

where  $\frac{dP_{\theta'}}{dP_{\theta}}$  is the ratio of the Radon Nikodym derivatives of  $P_{\theta}$ ,  $P_{\theta'}$  with respect to (say)  $P_{\theta} + P_{\theta'}$  and 0/0 is by convention equal to 1. Also let,

(2.2) 
$$J(\theta) = \inf\{K(\theta, \theta') : \theta' \in \Theta_0\}.$$

It is well known that (cf. [4]),  $0 \leq K(\theta, \theta') \leq \infty$ , and necessarily the same is true of  $J(\theta)$ . Following [4], let  $T_n$  be any sequence of extended real valued measurable functions of the infinite product space  $(X^{\infty}, \mathcal{A}^{\infty})$  such that  $T_n$  is a function of the first n co-ordinates only. Denote the cumulative distribution of  $T_n$ , when  $\theta$  obtains, by  $F_n(t,\theta)$ , i.e., 

where  $P_{\theta}$  now denotes the infinite product measure extension of  $P_{\theta}$  to  $X^{\infty}$ . Finally, let.

$$L_n(s) = \sup \left\{ 1 - F_n \big( T_n(s), \theta \big) : \theta \in \Theta_0 \right\}.$$
(2.4)
$$L_n(s) = \sup \left\{ 1 - F_n \big( T_n(s), \theta \big) : \theta \in \Theta_0 \right\}.$$

We assume that  $L_n$  is measurable. This for instance holds if  $F_n(T_n, \theta)$  is a separable stochastic process in  $\theta$  for  $\theta \in \Theta_0$ .

We can now state and prove, in the above framework,

Theorem 1. If

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

З 

Bahadur and Bickel

for every  $\theta \in \Theta - \Theta_0$ ,  $\theta' \in \Theta_0$  such that  $K(\theta, \theta') < \infty$ , then  $\liminf_{n} \frac{1}{n} \log L_n(s) \ge -J(\theta)$ З (2.5)with  $P_{\theta}$  probability 1 for every  $\theta \in \Theta - \Theta_0$ . *Proof.* Fix  $\theta \in \Theta - \Theta_0$ . Assume the Theorem has been proved for  $\Theta_0$  simple. Clearly we can suppose  $J(\theta) < \infty$  and can find  $\{\theta_m\}$  with  $K(\theta, \theta_m) < \infty$  and  $K(\theta, \theta_m) \to J(\theta)$ . But then  $\frac{1}{n}\log L_n(s) \ge \frac{1}{n}\log\left(1 - F_n(T_n(s), \theta_m)\right) \,.$ (2.6)By our assumption of the theorem for  $\Theta_0$  simple, we have  $\liminf_{n} \frac{1}{n} \log \left( 1 - F_n(T_n(s), \theta_m) \right) \ge -K(\theta, \theta_m)$ (2.7)with probability 1. Inequalities (2.6) and (2.7) then imply (2.5). If  $\Theta_0 = \{\theta_0\}$ ,  $P_{\theta}\left[\liminf_{n} \frac{1}{n} \log L_n(s) \ge -J(\theta)\right] = 1$ if and only if.  $P_{\theta}[1 - F_n(T_n, \theta_0) < a^n \exp{-nK(\theta, \theta_0)}$  infinitely often ] = 0(2.8)for every  $0 \le a \le 1$ . Fix a. Let  $A_n = \left[1 - F_n(T_n, \theta_0) < a^n \exp\{-nK(\theta, \theta_0)\}\right]$ . Then,  $P_{\theta_0}(A_n) \leq a^n \exp\{-nK(\theta, \theta_0)\}$ . (2.9)By the Neyman-Pearson lemma there exists  $c_n$ , such that,  $P_{\theta_0}\left\{\sum_{i=1}^n \log \frac{dP_{\theta}}{dP_{\theta_0}}(x_j) > nc_n\right\} \le a^n \exp\{-nK(\theta, \theta_0)\}$ (2.10)and,  $P_{\theta_0}\left\{\sum_{i=1}^n \log \frac{dP_{\theta}}{dP_{\theta_0}}(x_i) > nc_n\right\} \ge P_{\theta_0}(A_n) \ .$ (2.11)We require, **Lemma 1.** If (2.10) holds for all n, then there exists an  $\varepsilon > 0$  such that  $\liminf_n c_n \ge 0$  $K(\theta, \theta_0) + \varepsilon.$ *Proof.* By a Theorem of Chernoff [7],  $\frac{1}{n}\log P_{\theta_0}\left\{\sum_{i=1}^n\log\frac{dP_{\theta_0}}{dP_{\theta_0}}(x_i) > nz\right\} \to \inf_t\log H(t,z) ,$ (2.12)for  $z \ge \int_{X} \log \frac{dP_{\theta}}{dP_{\theta_0}}(x) dP_{\theta_0}(x) \; ,$ imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

where

$$H(t,z) = e^{-tz} \int_X \left(\frac{dP_\theta}{dP_{\theta_0}}\right)^t(x) dP_{\theta_0}(x) .$$
<sup>1</sup>
<sup>2</sup>
<sup>3</sup>
<sup>3</sup>

By the theory of the Laplace transform,  $0 \le H(t, z) \le \infty$ , H(t, z) is strictly convex in t wherever it is finite and if  $\inf_t H(t, z) < 1$ , the infimum is obtained for a unique t(z) given by the solution of the equation

$$\int \log \frac{dP_{\theta}}{dP_{\theta_0}}(x) \left[ \frac{dP_{\theta}}{dP_{\theta_0}}(x) \right]^t dP_{\theta_0}(x)$$

$$= \frac{\int \left[\frac{dP_{\theta_0}}{dP_{\theta_0}}(x)\right]^t dP_{\theta_0}(x)}{\int \left[\frac{dP_{\theta_0}}{dP_{\theta_0}}(x)\right]^t dP_{\theta_0}(x)} .$$

It is easily seen that if  $z_0 = K(\theta, \theta_0)$ , then  $t(z_0) = 1$ , and,

(2.14) 
$$\inf_{t} \log H(t, z_0) = -K(\theta, \theta_0) .$$
<sup>15</sup>

From (2.10), (2.12) and (2.14) we can immediately conclude that  $\liminf_n c_n \geq 1$  $K(\theta, \theta_0)$ . But, in fact, by the implicit function theorem as  $z \to K(\theta, \theta_0)$  we have  $t(z) \to 1$ , and  $\log H(z, t(z)) \to K(\theta, \theta_0)$ , by dominated convergence. Choose  $z_1 > z_0$ such that  $H(z_1, t(z_1) > a \exp\{-K(\theta, \theta_0)\}$ . Then  $\varepsilon = z_1 - z_0$  will do for the Lemma.

It now follows from the basic assumption of the theorem, by a result of Erdös, Hsu and Robbins [9] that,

and this by (2.11) and the Borel Cantelli lemma suffices for (2.8) and the theorem to hold.

Remarks. 1. Erdös has shown in [9] that our second moment assumption is necessary as well as sufficient for (2.15) to hold. Although, of course, (2.15) is not necessary for (2.8) the relative arbitrariness of the  $A_n$  apart from condition (2.9) would suggest that the Theorem may be false if some condition such as the one imposed does not hold.

2. As in [4], if we define  $N(\varepsilon, s) = \text{least positive } m$  such that  $L_n \ge \varepsilon$  for all  $n \ge m$ and  $\infty$  otherwise, we have under the assumptions of our theorem 1,

(2.16) 
$$\lim \inf_{\varepsilon \to 0} \frac{N(\varepsilon, s)}{-\log \varepsilon} \ge \frac{1}{J(\theta)} \text{ a.s. } P_{\theta} .$$

# 3. General Assumptions and a Useful Lemma

Before giving further structural assumptions needed for Sections 4 and 5 we prove a simple general lemma already implicit in [4] stating a useful sufficient condition for a sequence  $\{T_n\}$  to be optimal. We say  $\{T_n\}$  is asymptotically optimal if,

50  
51 (3.1) 
$$\lim_{n} \frac{1}{n} \log L_n(s) = -J(\theta)$$

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

Bahadur and Bickel

with  $P_{\theta}$  probability 1 for all  $\theta \in \Theta - \Theta_0$ . Then (3.1) implies, (cf. [4])  $\lim_{\varepsilon \to 0} \frac{N(\varepsilon, s)}{-\log \varepsilon} = \frac{1}{J(\theta)} \text{ a.s. } P_{\theta} .$ З (3.2)Lemma 2. If the conclusion of Theorem 1 holds and i)  $\liminf_{n} T_n \ge J(\theta) \ a.s. \ P_{\theta}$ ii)  $\limsup_{n} \log \left(1 - G_n(t)\right) \le -t$ where  $G_n(t) = \inf \{F_n(t, \theta_0) : \theta_0 \in \Theta_0\}$  and  $\theta$  ranges over  $\Theta - \Theta_0$ , then  $\{T_n\}$  is asymptotically optimal. *Proof.* It clearly suffices to show that  $\limsup_{n} \frac{1}{n} \log L_n(s) \le -J(\theta)$ (3.3)with  $P_{\theta}$  probability 1. But since  $L_n(s) = 1 - G_n(T_n)$  and  $1 - G_n(t)$  is monotone decreasing, i) and ii) obviously imply (3.3). We begin by giving nine general assumptions which are sufficient to ensure the validity of Theorem 2 of the main section. Assumption 1. There exists a c finite measure  $\mu$  on  $(X, \mathcal{A})$  which dominates the family  $\{P_{\theta}\}$ . We denote the density of  $P_{\theta}$  with respect to  $\mu$  by  $f(x, \theta)$ . Then,  $\frac{dP_{\theta}(x)}{dP_{\theta'}} = \frac{f(x,\theta)}{f(x,\theta')} \text{ a.e. } P_{\theta} + P_{\theta'} .$ Assumption 2.  $\Theta$  is a metric space. The topological Borel field on  $\Theta$  is denoted by  $\mathcal{B}$ .  $f(x,\theta)$  is bimeasurable in  $(x,\theta)$  on  $(X \times \Theta, \mathcal{A} \times \mathcal{B})$ . Assumption 3. We are given a probability measure  $\nu$  on  $(\Theta, \mathcal{B}), \Theta_0 \in \mathcal{B}$ , and  $\nu(\Theta_0) > 0$ . Moreover, if  $S(\theta, d)$  is the open sphere of centre  $\theta$  and radius d,  $\nu \{ S(\theta, d) \cap [\Theta - \Theta_0] \} > 0 \text{ for all } \theta \in \Theta - \Theta_0 \text{ and } d > 0.$ 

Assumption 4. There exists a suitable metric compactification  $\overline{\Theta}_0$  of  $\Theta_0$  (viz [4]). That is, we first define  $S(\theta, d)$  to be the sphere of radius d and centre  $\theta$  in  $\Theta_0$  and then take,

(3.4) 
$$g_0(x,\theta,d) = \sup \left\{ f(x,\lambda) : \lambda \in \hat{S}(\theta,d) \cap \Theta_0 \right\}.$$

We assume  $g_0$  is measurable in x for d sufficiently small and define,

(3.5) 
$$g_0(x,\theta,0) = \lim_{d \to 0} g_0(x,\theta,d) .$$
<sup>46</sup>

The final assumption, (see [4]) is,

50  
51 (3.6) 
$$E_{\theta}\left(\frac{g(x,\theta',0)}{f(x,\theta)}\right) \le 1$$
, 50  
51 51

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

where  $E_{\theta}(h(x))$  denotes  $\int_X h(x) dP_{\theta}(x)$  for any integrable function h. Assumption 5. Define for all  $\theta' \in \overline{\Theta}_0, \ \theta \in \Theta - \Theta_0$ , З  $\bar{K}(\theta, \theta') = -E_{\theta} \left( \frac{\log g_0(x, \theta', 0)}{f(x, \theta)} \right) \,.$ (3.7)(3.6) and Jensen's inequality guarantee  $0 \le \bar{K} \le \infty$ . Assume  $J(\theta) = \inf\{K(\theta, \theta') : \theta' \subset \bar{\Theta}_0\}.$ (3.8)Assumption 6.  $E_{\theta}\left(\log \frac{g_0(x, \theta', d)}{f(x, \theta)}\right) < \infty$ , (3.9)for all  $\theta \in \Theta - \Theta_0$ ,  $\theta' \in \Theta_0$ . As in, [4] p.22, this is equivalent to,  $\lim_{d \to 0} E_{\theta} \left( \frac{\log g_0(x, \theta', d)}{f(x, \theta)} \right) \le -K(\theta, \theta') \; .$ (3.10)Assumption 7. Define,  $\eta(x,\theta,d) = \inf \left\{ \log \frac{f(x,\lambda)}{f(x,\theta)} : \lambda \in S(\theta,d) \cap [\Theta - \Theta_0] \right\}.$ (3.11)Assume that  $\eta$  is a measurable function of x for d sufficiently small and that,  $\lim_{d \to 0} E_{\theta} \big( \eta(x, \theta, d) \big) = 0$ (3.12)for all  $\theta \in \Theta - \Theta_0$ . Assumption 8. Define, for  $\theta' \in \Theta_0$ ,  $\gamma(x,\theta',d) = \log \inf \left\{ \frac{f(x,\lambda)}{f(x,\theta')} : \ \lambda \in S(\theta',d) \cap \Theta_0 \right\} \,,$ (3.13) $\varphi(t,\theta',d) = E_{\theta'}\left(\exp\{-t\gamma(x,\theta',d)\}\right).$ (3.14)For every  $0 < \rho < 1$ ,  $\beta > 0$ , there exists  $d(\theta', \rho, \beta)$  such that  $\inf e^{-t\beta}\varphi(t,\theta',d(\theta',\rho,\beta)) \leq \frac{\rho}{2} .$ Assumption 9.  $\inf \left\{ \nu \left[ S(\theta', d(\theta', \rho, \beta)) \cap \Theta_0 : \theta' \in \Theta_0 \right] \right\} = m(\rho, \beta) > 0 .$ We shall now examine these assumptions in turn giving where necessary stronger, but more easily checkable conditions, which we shall denote by primes. Thus, (say) assumption 4' will imply assumption 4 and the conclusion of theorem 2 will continue 

| 1        | to hold if 4 is replaced by 4'. The more important of these useful weakenings of  | 1  |
|----------|---|----|
| 2        | Theorem 2 will be isolated as a Corollary.  | 2  |
| 3        | Assumption 1 is self-explanatory and clearly cannot be weakened appreciably.  | 3  |
| 4        | The requirement that $\Theta$ be a metric space can clearly be dropped and replaced   | 4  |
| 5        | by the requirement that $\Theta$ be a topological space and $\mathcal{B}$ the topological Borel field.  | 5  |
| 6        | However, the notational convenience involved in being able to define quantities in  | 6  |
| 7        | terms of spheres of a given radius rather than neighbourhood bases seems well   | 7  |
| 8        | worth the loss of generality. On the other hand, assumption 2 is obviously satisfied  | 8  |
| 9        | if we have the usual.   | 9  |
| 10       |   | 10 |
| 11       | <b>Assumption 2</b> ' $\Theta$ is a subset of k dimensional Euclidean space with the usual  | 11 |
| 12       | metric topology $\mathcal{B}$ is the Borel $\sigma$ field and $f(r, \theta)$ is himeasurable  | 12 |
| 13       | Weakenings of assumption 3 do not fit readily into this program, but we mention   | 13 |
| 14       | that we can drop the requirement that $u$ has a probability (finite) massure if the   | 14 |
| 15       | following two conditions hold, as well as the second part of assumption 3:  | 15 |
| 16       | following two conditions hold, as well as the second part of assumption 5.  | 16 |
| 17       | (1) There exists N such that, $\int_{\Theta} \prod_{i=1}^{N} f(x_i, \lambda) \nu(d\lambda) < \infty$ a.s. $P_{\theta}$ and  | 17 |
| 18       | (2) $\int_{\Omega} \prod_{i=1}^{N} f(x_i, \lambda) \nu(d\lambda) > 0, \int_{\Omega} \prod_{i=1}^{N} f(x_i, \lambda) \nu(d\lambda) > 0$ a.s. $P_{\theta}$ for all $n > 0$  | 18 |
| 19       | N.  | 19 |
| 20       |   | 20 |
| 21       | For details of the proof of i) of Lemma 2 for $T_n$ under these assumptions we refer  | 21 |
| 22       | to [6]. The basic idea of this generalization is to consider the process of observation   | 22 |
| 23       | as really starting after N with prior distribution, the posterior distribution of $\theta$  | 23 |
| 24       | given $x_1, \ldots, x_N$ , which by (1) is a true probability distribution. In fact, we can in  | 24 |
| 25       | general make dependent $\nu$ on the observations all along if we modify the second  | 25 |
| 26       | part of assumption 3 and assumption 9 suitably. The generalization given above  | 26 |
| 27       | is of interest in the case when reasonable tests arise from improper "priors", e.g.   | 20 |
| 28       | Lebesgue measure.   | 28 |
| 20       | The most natural replacement of assumption 4 is of course assuming that $\Theta_0$ is   | 20 |
| 30       | already compact. In this case, we have, <b>Assumption 4'.</b> $\Theta_0$ is compact.  | 30 |
| 31       |   | 31 |
| 32       | We can then drop 5, but must replace 6 by its equivalent form,  | 32 |
| 33       |   | 33 |
| 34       | <b>Assumption 6'.</b> $\lim_{d\to 0} E_{\theta} \log \left( \frac{g_0(x,\theta',d)}{(x,\theta')} \right) \leq -K(\theta,\theta')$ for all $\theta \in \Theta - \Theta_0, \theta' \in \Theta_0$ .  | 34 |
| 35       | Measurability of a must still of source be involved. Assurption 6 measured equation 6   | 35 |
| 36       | Measurability of $g_0$ must still of course be invoked. Assumption 0 may replace 0<br>if $f(m, \theta)$ is continuous in $\theta$ for almost all $m$ . A loss stringent modification in some  | 36 |
| 37       | If $f(x, \theta)$ is continuous in $\theta$ for almost an $x$ . A less stringent modulication in some   | 37 |
| 38       | senses which is most useful is combining 4, 5, and 6 with 2 to give:  | 30 |
| 20       |   | 20 |
| 39<br>40 | Assumption (4,5,6)" Assumption 2' holds and   | 39 |
| 40       | (a) $\lim_{x \to \infty} E_{\theta}\left(\log^{g_0(x,\theta',d)}\right) \leq K(\theta,\theta')$ for $\theta \in \Theta$ , $\theta' \in \Theta$ .  | 40 |
| 41       | (a) $\lim_{d\to 0} \mathcal{D}_{\theta}\left(\log -\frac{f(x,\theta)}{f(x,\theta)}\right) \leq -K(\theta,\theta)$ for $\theta \in O - O_0, \theta \in O_0$ .  | 41 |
| 42       | (b) $\lim_{d\to\infty} E_{\theta} \left( \log \sup \left\{ \log \frac{f(x_{\lambda},\lambda)}{f(x_{\lambda},\lambda)} : \lambda \in \Theta_0 \ \lambda\  > d \right\} \right) \le -J(\theta)$ , for all $\theta \in I$  | 42 |
| 43       | $(\circ)  \min_{a \to \infty} \sum_{i=0}^{n} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) + i \in \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log \operatorname{sap} \left( \log f(x,\theta) \right) \right) = \operatorname{contract} \left( \log f(x,\theta) \right) = \operatorname{contract} \left( $ | 43 |
| 44       | $\Theta - \Theta_0$ , where $\  \cdot \ $ is the usual Euclidean norm.  | 44 |
| -1-J     | This assumption is clearly equivalent to 5 and 6 if in 4 we take $\bar{\Theta}_{c}$ to be the closure   | 45 |
| 40       | of $\Theta_0$ in the one point compactification of $R$  | 40 |
| 41<br>19 | Assumption 7 is most readily replaced by  | 47 |
| 40<br>10 | rissamption r is most readily replaced by,  | 48 |
| 49<br>50 | Assumption 7? $f(x, \theta)$ is continuous in $\theta$ for element all $\pi(y)$ and $F(\pi(x, \theta, d))$  | 49 |
| 50       | Assumption $i \cdot j(x, \theta)$ is continuous in $\theta$ for almost all $x(\mu)$ and, $E_{\theta}(\eta(x, \theta, a))$   | 50 |
| 51       | $\sim \infty$ for some $u \geq 0$ for each $v \in O = O_0$ .  | 51 |

Assumption 7' and the dominated convergence theorem readily imply 7. We need not require measurability of  $g_0$  in this case in view of 4 or 2' since  $\Theta_0$  being a subset of a separable metric space is separable. The same is true of  $\eta$  if 2' holds or more З generally if  $\Theta$  is a separable metric space. A useful substitute for assumption 8 is Assumption 8'. There exists an  $M < \infty$  such that for every  $0 < T < \infty$  we can find a  $d^x(\theta, T)$  with  $\varphi(T, \theta', d^x) \leq M$ . Assumption 8' easily implies 8 since, then,  $\inf_t e^{-t\beta}\varphi(T,\theta',d^x) \le e^{-T\beta}M < \frac{\rho}{2} ,$ (3.15)for T sufficiently large. In many situations 8' is most easily verified by showing that,  $\varphi(t, \theta', d) \to 1$ , (3.16)uniformly on compacts in t as  $d \to 0$ . This in turn is implied by, **Assumption 8".**  $f(x,\theta)$  is continuous in  $\theta$  for almost all  $x(\mu)$  and for every 0 < 0 $T < \infty, \varphi(T, \theta', d) < \infty$  for d sufficiently small. Assumption 8" implies (3.16) by way of the dominated convergence theorem if we remark that  $\varphi(t, \theta', d)$  is monotone increasing in t for every fixed  $\theta', d$ . Finally, we can replace assumption 9 by, **Assumption 9'.**  $d(\theta, \rho, \beta)$  is independent of  $\theta'$ , assumption 4' holds, and  $\nu [S(\theta', d) \cap$  $|\Theta_0| > 0$  for all  $\theta' \in \Theta_0$ . To show that 9' implies 9 we need only prove that,  $\inf \left\{ \nu \left[ S(\theta', d) \cap \Theta_0 \right] : \theta' \in \Theta_0 f > 0 \right\}.$ (3.17)Then, if  $\theta'_n \to \theta'$  and,  $S^*(\theta, d) = S(\theta, d) \cap \Theta_0$  $\nu \left[ S^*(\theta',d) \right] - \nu \left[ S^*(\theta'_n,d) \right] = \nu \left[ \lambda : \ \delta(\lambda,\theta') < d, \quad \delta(\lambda,\theta'_n) \ge d, \ \lambda \in \Theta_0 \right]$  $-\nu [\lambda: \ \delta(\lambda, \theta') \ge d, \quad \delta(\lambda, \theta'_n) < d, \ \lambda \in \Theta_0].$ (3.18)Clearly, if  $n \to \infty$ , the first term of the above difference tends to 0 since the set whose measure is computed tends to the empty set. Therefore, for fixed d,  $\nu[S^*(\theta',d)]$  is a lower semi-continuous function of  $\theta'$  on  $\Theta_0$  and 2', and the third part of 9', imply (3.17). This completes our roster of simplifying assumptions. Clearly if further restric-tions are put on  $f(x,\theta)$  the verification of most can be very easy. A strong form of Lemma 3,  $\tau_n \to J(\theta)$ , is given under very weak assumptions in [5] if  $f(x,\theta)$  is of exponential type, (Theorem 4.1). If the conditions 1-7 of this paper are made two sided (viz. [5] Theorem 4.2) one can obtain this strengthening of Lemma 3 in gen-eral. In fact, the conditions detailed in theorem 4.2 of [5] are somewhat restrictive 

<sup>47</sup> The assumptions required by this paper but not by [4] are 7, 8, and 9. On <sup>48</sup> the other hand, assumption 6 of that paper and the fact that  $\Theta$  can be suitably <sup>49</sup> compactified (rather than just  $\Theta_0$ ) is not required by us. Since simple hypotheses <sup>50</sup> tend to be somewhat more common than simple alternatives this seems a gain. <sup>51</sup> Otherwise the non structural assumptions of this paper and [4] coincide.

versions of our conditions 2'-7'.

4. The Main Theorem Given  $\nu$  in assumption 3 we now define, З  $\bar{T}_n = \frac{1}{n} \log \frac{\int_{\Theta - \Theta_0} \prod_{i=1}^n f(x_i, \lambda) \nu(d\lambda)}{\int_{\Omega} \prod_{i=1}^n f(x_i, \lambda) \nu(d\lambda)} .$ (4.1)By assumption 2,  $\overline{T}_n$  is well defined  $\left(\frac{\infty}{\infty} = 1, \frac{0}{0} = 1\right)$ . In fact,  $\overline{T}_n$  is a version of the test statistic a Bayesian with prior  $\nu$  would use to test  $H: \theta \in \Theta_0$ , rejecting for large values of  $\overline{T}_n$ . We can now state the principal theorem of the paper. **Theorem 2.** If assumptions 1–9 and the conclusion of Theorem 1 holds, then  $\{\overline{T}_n\}$ is asymptotically optimal. *Proof.* The proof preceds by way of some lemmas. We have first, **Lemma 3.** Under assumptions 1–7,  $\{\overline{T}_n\}$  satisfies condition i) of Lemma 2. *Proof.* Suppose that  $\theta \in \Theta - \Theta_0$  holds. Define,  $U_n(s,\theta) = \int_{\Theta - \Theta_0} \prod_{i=1}^n \frac{f(x_i,\lambda)}{f(x_i,\theta)} \nu(d\lambda)$ (4.2) $V_n(s,\theta) = \int_{\Theta_0} \prod_{i=1}^n \frac{f(x_i,\lambda)}{f(x_i,\theta)} \nu(d\lambda) \ .$ (4.3)Then,  $\bar{T}_n = \frac{1}{n} \log \frac{U_n(s,\theta)}{V_n(s,\theta)} \; .$ (4.4)We show first,  $\limsup_{n \to \infty} \frac{1}{n} \log V_n(s,\theta) \le -J(\theta)$ (4.5)a.s.  $P_{\theta}$ . Note that,  $\frac{1}{n}\log V_n(s,\theta) \le \frac{1}{n}\log\left(\sup\left\{\frac{1}{n}\sum_{i=1}^n\log\frac{f(x_i,\lambda)}{f(x_i,\theta)} : \lambda\in\Theta_0\right\}\right).$ (4.6)And the right hand side equals  $R_0(\Theta_0, \theta)$  in the notation of [4] which converges to  $-J(\theta)$  a.s.  $P_{\theta}$  by Lemma 4 of that paper. An examination of the proof of this Lemma will show that only our assumptions 1–6 are used. To establish the Lemma we need now only show,  $\liminf_{n} \frac{1}{n} \log U_n(s,\theta) \ge 1 \text{ a.s. } P_{\theta}.$ (4.7)By assumption 7 we can find  $d_1(\theta, \varepsilon)$  such that,

(4.8) $E_{\theta}(\eta(X_i, \theta, d_1)) \geq -\varepsilon$ . 

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

But,

$$\frac{1}{n}\log U_n \ge \log \nu \left[ S(\theta, d_1) \cap [\Theta - \Theta_0] \right]$$

$$+\inf\left\{\frac{1}{n}\sum_{i=1}^{n}\frac{f(X_{i},\lambda)}{f(X_{i},\theta)}: \lambda \in S(\theta,d_{1}) \cap [\Theta - \Theta_{0}]\right\}$$

(4.9) 
$$\geq \frac{1}{n} \sum_{i=1}^{n} \eta(x_1, \theta, d_1) + \log \nu \left[ S(\theta, d_1) \cap \left[ \Theta - \Theta_0 \right] \right] \,.$$

Letting  $n \to \infty$  and then  $\varepsilon \to 0$ , (4.9), assumption 3, and the strong law of large numbers imply (4.7). The Lemma follows.

We complete the proof of the Theorem by way of two further Lemmas.

**Lemma 4.** Under the first part of assumption 3, for all 
$$n, t, \theta' \in \Theta_0$$
,

(4.10) 
$$P_{\theta'}\left[\frac{1}{n}\log U_n(s,\theta') \ge t\right] \le e^{-nt}.$$

Proof.

$$P_{\theta'}\left[\frac{1}{n}\log U_n(s,\theta') \ge t\right] = \int_W \prod_{i=1}^n f(x_i,\theta') \ \mu(dx_1)\dots\mu(dx_n)$$

where 
$$W = \left[s : \prod_{i=1}^{n} f(x_i, \theta') \le e^{-nt} \int_{\Theta - \Theta_0} \prod_{i=1}^{n} f(x_i, \lambda) \nu(d\lambda)\right].$$

Thus our probability is bounded above by

(4.11) 
$$e^{-nt} \int_{X^n} \int_{\Theta - \Theta_0} \prod_{i=1}^n f(x_i, \lambda) \nu(d\lambda) \ \mu(dx_1) \dots \mu(dx_n) \ .$$

But the right hand side of (4.11), by Fubini's theorem, is

$$e^{-nt} \int_{\Theta-\Theta_0} \int_{X^n} \prod_{i=1}^n f(x_i,\lambda) \ \mu(dx_1) \dots \mu(dx_n) \nu(d\lambda) \le e^{-nt} \ .$$

**Lemma 5.** Under assumptions 8 and 9 if  $\theta' \in \Theta_0$ , for every  $0 < \rho < 1$ ,  $\beta > 0$ , there exist  $N(\rho, \beta)$  such that for  $n \ge N(\rho, \beta)$ 

(4.12) 
$$P_{\theta'}\left[\frac{1}{n}\log V_n(s,\theta') \ge -\beta\right] \le \rho^n .$$

*Proof.* Choose  $d(\theta', \rho, \beta)$  as in assumptions 8 and 9. Then since

$$\frac{1}{n}\log V_n(S,\theta') \ge \log \nu \left[ S(\theta', d(\theta', \rho, \beta) \cap \Theta_0 \right]$$
<sup>44</sup>
<sup>45</sup>

$$+ \inf\left\{\frac{1}{n}\sum_{i=1}^{n}\log\frac{f(x_i,\lambda)}{f(x_i,\theta')} : \lambda \in S(\theta', d(\theta',\rho,\beta) \cap \Theta_0\right\}$$

$$\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$\sum_{51}^{50} (4.13) \geq \log \nu \left[ S(\theta', d(\theta', \rho, \beta)) \cap \Theta_0 \right] + \frac{1}{n} \sum_{i=1}^{1} \gamma \left( x_i, \theta', d(\theta', \rho, \beta) \right),$$

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

З

(4.14)  $P_{\theta'}\left[\frac{1}{n}\log V_n(s,\theta') \le -\beta\right] \le P_{\theta'}\left[\sum_{i=1}^n \gamma(x_i,\theta',d(\theta',\rho,\beta)) \le -n\beta - \log m(\rho,\beta)\right].$ By Lemma 1 of [4]  $P_{\theta'}\left[\frac{1}{n}\sum_{i=1}^n \gamma(x_i,\theta',d(\theta',\rho,\beta)) \le -n\beta - \log m(\rho,\beta)\right]$   $(4.15) \le \inf_t \left\{\exp\left[-\beta t - \log m(\rho,\beta)\frac{t}{n}\right]\varphi(t,\theta',d(\theta',\rho,\beta))\right\}^n$ for  $n \ge -\frac{\log m(\rho,\beta)}{n}$  which is finite by assumption 0. If  $\inf_{x_i=1}^{\infty} e^{-\beta t}e^{it_i \theta'_i} d(\theta',\rho,\beta)$ 

for  $n \ge -\frac{\log m(\rho,\beta)}{\beta}$  which is finite by assumption 9. If  $\inf_t e^{-\beta t} \varphi(t,\theta',d(\theta',\rho,\beta))$  is attained for  $t = t_0$  which is strictly positive by lemma 1 of [4], we can choose,

$$N(\rho,\beta) \leq \max\Big(-\frac{\log m(\rho,\beta)}{\beta},-\frac{t_0\log m(\rho,\beta)}{\log 2}\Big) < \infty \ .$$

Now,

$$P_{\theta'}\left[\frac{1}{n}\log\bar{T}_n \ge t\right] \le P_{\theta'}\left[\frac{1}{n}\log U_n(s,\theta') \ge -\beta, \frac{1}{n}\log V_n(s,\theta') \ge -\beta\right]$$

(4.16) 
$$+ P_{\theta'} \left[ \frac{1}{n} \log V_n(s, \theta') \le -\beta \right] \,.$$

By Lemmas 4 and 5 the right hand side of (4.16) is bounded above by  $e^{-n(t-\beta)} + \rho^n$  for  $n \ge N(\rho, \beta)$ . Hence,

$$\limsup_{n} \frac{1}{n} \log \sup \left\{ P_{\theta'}[\bar{T}_n \ge t] : \ \theta \in \Theta_0 \right\} \le \limsup_{n} \frac{1}{n} \log \left( e^{-n(t-\beta)} + \rho^n \right)$$

$$(4.17) = \max[-(t-\beta), \log \rho] .$$

Letting  $\rho \to 0$  first and then  $\beta \to 0$ , we find that ii) of Lemma 2 is satisfied by  $\overline{T}_n$ and the Theorem is proved. Gathering the most useful of the "prime" assumptions together we state,

**Corollary 1.** If assumptions 1, 2', 3, (456)", 7', 8", and 9 hold, then the conclusion of Theorem 2 is valid.

The most immediate field of application of Corollary 1 is when  $f(x, \theta)$  is the density of a Koopman-Darmois (exponential) family.

(4.18) 
$$f(x,\theta) = e^{\theta \cdot t(x)}$$

where  $\theta = (\theta_1, \dots, \theta_k), t(x) = (t_1(x), \dots, t_k(x))$  and the  $\theta_j, t_j$  are real. Assumptions 1, 2', (456)", 7' and 8" are then automatically satisfied and we need only impose conditions 3, and 9 on  $\nu$ . If  $\Theta_0$  is compact, 9' is automatic and 3 is all that is needed.

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

we have,

З

5. Optimality of Minimax Tests The main result of this section is an immediate consequence of the following Lemma. We retain the notation of the previous section, defining only  $\bar{F}_n(t,\theta) = P_{\theta}[\bar{T}_n < t]$ . (5.1) $\bar{G}_n(t) = \inf \left\{ \bar{F}_n(t,\theta') : \theta' \in \Theta_0 \right\} .$ (5.2)**Lemma 6.** Suppose that there exists a measurable subset S of  $\Theta_0$  such that, iii)  $1 - \bar{F}_n(t, \theta')$  is a constant on S (for fixed n, t) as a function of  $\theta'$ . iv) sup  $\{1 - \bar{F}_n(t, \theta') : \theta' \in S\} = 1 - \bar{G}_n(t).$ v)  $\nu[\Theta_0 - S] = 0.$ Then,  $\limsup_{n \to \infty} \frac{1}{n} \log \left( 1 - \bar{G}_n(t) \right) \le -t \; .$ (5.3)Proof.  $\left(1 - \bar{G}_n(t)\right) = \sup\left\{1 - \bar{F}_n(t,\theta'): \ \theta' \in S\right\} = \int_{\Theta_n} P_{\theta'}[\bar{T}_n \ge t]\nu(d\theta')$ (5.4)by iii), iv), v). Now, let  $C_n = \left[ \int_{\Theta_0} \prod_{i=1}^n f(x_i, \lambda) \nu(d\lambda) \le e^{-nt} \int_{\Theta-\Theta_0} \prod_{i=1}^n f(x_i, \lambda) \nu(d\lambda) \right].$ Then.  $\int_{\Theta_0} P_{\theta'}[\bar{T}_n < t]\nu(d\theta') = \int_{\Theta_0} \int_{C_n} \prod_{i=1}^n f(x_i, \lambda)\mu(dx_1)\dots\mu(dx_n)\nu(d\lambda)$  $= \int_{C_{-1}} \int_{\Theta_{0}} \prod_{i=1}^{n} f(x_{i}, \lambda) \mu(dx_{1}) \dots \mu(dx_{n}) \nu(d\lambda)$  $\leq e^{-nt} \int_{C_n} \int_{\Theta - \Theta_0} \prod_{i=1}^n f(x_i, \lambda) \nu(d\lambda) \mu(dx_1) \dots \mu(dx_n) \leq e^{-nt} .$ (5.5)We formulate, Assumption 10. For every n, ta)  $\nu \Big[ \theta \in \Theta - \Theta_0 : \sup \big\{ \bar{F}_n(t, \lambda) : \lambda \in \Theta - \Theta_0 \big\} > \bar{F}_n(t, \theta) \Big] = 0.$ b)  $\nu \Big[ \theta' : \bar{G}_n(t) < F_n(t, \theta') : \theta' \in \Theta_0 \Big] = 0.$ We can now state, **Theorem 3.** If assumption 10 holds the test which rejects if  $\overline{T}_n \geq t$  is minimax for  $H: \theta \in \Theta_0 \text{ vs } K: \theta \in \Theta - \Theta_0 \text{ at level } 1 - \overline{G}_n(t).$  If assumptions 1-7 and assumption 10b) hold, then the sequence of test statistics  $\{\overline{T}_n\}$  is asymptotically optimal.

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

*Proof.* The first part of the Theorem is classical (c.f. Lehmann [8] p.327). The second part is an immediate consequence of Theorem 2 and Lemma 5 since assumption 10b) is equivalent to iii), iv), v).

This Theorem is of interest in connecting the classical finite sample optimality results with stochastic comparison. The most immediate application of this result is in the one-parameter exponential family where minimax tests of (say)  $H: \theta \leq \theta_0$ vs  $K: \theta = \theta_1 > \theta$  are. Bayes tests with respect to two point distributions satisfy assumption 10 (cf. [8]). Unfortunately proving optimality directly is trivial in this case. More interesting candidates are in the normal situation the *t*-statistic and the  $S^2$  statistic used in testing  $H: \mu < \mu_0$  and  $H: \sigma \leq \sigma_0$  when  $\mu$ , and  $\sigma$  are respectively unknown. Although we are here presented with a situation which does not quite fall under Theorem 3, ( $\nu$  satisfying 10a,b depends on *n*, (cf. [8], p.94), one can easily check the conclusion of Lemma 2 directly and then apply Lemma 5 to obtain optimality.

Finally, it may be interesting to see that from a quasi Bayesian point of view, if stochastic comparison is defined in terms of the observed expected level of significance, (where the expectation is taken under the prior) then Lemma 5 and assumptions 1–7 and an analogue of Theorem 1 guarantee the optimality of  $\bar{T}_n$  in this sense. Formally for any sequence  $\{T_n\}$  we would then consider not  $L_n$  but,

$$L_n^*(s) = \left(1 - G_n^*(T_n)\right)$$

where  $G_n^*(t) = \int_{\Theta_0} F_n(t, \theta') \nu(d\theta').$ 

The analogue of the conclusion of Theorem 1 needed would be that a.s.  $P_{\theta}$ ,

imsart-coll ver. 2008/08/29 file: Bahadur.tex date: March 25, 2009

 $\liminf_n \frac{1}{n} \log L_n^* \geq -J(\theta) \ .$ References [1] BAHADUR, R. R. (1960). Stochastic comparison of tests. Ann. Math. Stat., 31 276-295. [2] BAHADUR, R. R. (1960). Asymptotic efficiency of tests and estimates. Sankhya, 22 229-252. BAHADUR, R. R. (1960). Simultaneous comparison of the optimum and sign tests of a normal mean. [3] In Contributions to Probability and Statistics. Stanford Univ. Press, 79-88. BAHADUR, R. R. (1965). An optimal property of the likelihood ratio statistic. In Proceedings of the [4]5<sup>th</sup> Berkeley Symposium, 1, 13–26. BICKEL, P. J. and YAHAV, J. A. (1965). Asymptotically pointwise optimal procedures in sequential analysis. In Proceedings of the  $5^{th}$  Berkeley Symposium, 1, 401–415. BICKEL, P. J. and YAHAV, J. A. (1969). Asymptotic theory of Bayes solution. Prob. Theory and [6] Related Fields, 11, 257-276. CHERNOFF, H. (1952). A measure of asymptotic efficiency for tests of hypothesis based on the sum [7] of observations. Ann. Math. Stat., 23, 493-507. ERDÖS, P. (1949). On a theorem of Hsu and Robbins. Ann. Math. Stat., 20, 286-291. [9] LEHMANN, E. L. (1959). Testing Statistical Hypotheses. J. Wiley and Sons. 

# On the Non-Optimality of Optimal **Procedures**

З

#### Peter J. Huber<sup>1</sup>

Abstract: This paper discusses some subtle, and largely overlooked, differences between conceptual and mathematical optimization goals in statistics, and illustrates them by examples.

| 14 | Co | ontents   | 14 |
|----|----|---|----|
| 15 |    |   | 15 |
| 16 | 1  | Introduction  | 16 |
| 17 | 2  | On Optimization and Models                                | 17 |
| 18 | 3  | Classical Mathematical Statistics and Decision Theory     | 18 |
| 19 | 4  | Tukey's 1962 Paper         32                             | 19 |
| 20 | 5  | Pitfalls of Optimality                                    | 20 |
| 21 | 6  | Examples from Classical Statistics                        | 21 |
| 22 |    | 6.1 The Fuzzy Concepts Syndrome                           | 22 |
| 23 |    | 6.2 The Straitjacket Syndrome                             | 23 |
| 24 |    | 6.3 The Scapegoat Syndrome                                | 24 |
| 25 | 7  | Problems with Optimality in Robustness                    | 25 |
| 26 |    | 7.1 Asymptotic Robustness for Finite $\varepsilon > 0$    | 26 |
| 27 |    | 7.2 Finite Sample Robustness for Finite $\varepsilon > 0$ | 27 |
| 28 |    | 7.3 Asymptotic Robustness for Infinitesimal $\varepsilon$ | 28 |
| 29 |    | 7.4 Optimal Breakdown Point                               | 29 |
| 30 | 8  | Design Issues   | 30 |
| 31 |    | 8.1 Optimal Designs                                       | 31 |
| 32 |    | 8.2 Regression Design and Breakdown                       | 32 |
| 33 | 9  | Bayesian Statistics                                       | 33 |
| 34 | 10 | Concluding Remarks  | 34 |
| 35 | Re | eferences   | 35 |
| 36 |    |   | 36 |
| 37 |    |   | 37 |
| 38 | 1. | Introduction  | 38 |

First, we shall identify those parts of statistics that rely in a crucial fashion on optimization. The most conspicuous among them are: classical mathematical statistics, decision theory, and Bayesian statistics.

Classical mathematical statistics was created by R. A. Fisher (1922), in a paper concerned with estimation, and by J. Neyman and E. S. Pearson (1933), in a paper concerned with testing. It was brought to completion by E. L. Lehmann in his lecture notes (1949, 1950); those notes later grew into two books (1959, 1983). Around the same time when Lehmann produced his lecture notes, A. Wald (1950) expanded the scope of mathematical statistics by creating statistical decision theory.

<sup>&</sup>lt;sup>1</sup>POB 198, 7250 Klosters, Switzerland, email: peterj.huber@bluewin.ch

Keywords and phrases: optimality, superefficiency, optimal robustness, breakdown point, optimal design, Bayesian robustness 

З

The central concerns of classical mathematical statistics were *efficiency* in estimation (i.e. minimum variance), and *power* in testing problems, both being optimality concerns. Decision theory confirmed the central interest in optimality, but shifted the emphasis to *admissibility* and *minimaxity*.

A heavy slant towards optimality, of a different origin, holds also for Bayesian statistics. For a given model, consisting of a prior distribution and a family of conditional distributions, the Bayes formula by definition gives the "best" procedure; it is admissible in decision theoretic terminology.

The above-mentioned three areas of statistics appear to be the only ones where optimality is central to the theory. Elsewhere, optimality seems to provide mere icing on the cake. Note that the papers of Fisher and of Neyman-Pearson imprinted subsequent generations of statisticians with an (often uncritical) love of optimality. By 1960, as a young mathematical statistician you would not dare submit a new procedure to a reputable journal, unless you could prove some optimality property. (Later, there was a reversal, and too many statistical algorithms may have slipped through the editorial gates with enthusiastic but inadequately substantiated claims.)

#### 2. On Optimization and Models

Mathematical optimization always operates on some model. Models are simplified approximations to the truth; the hope is that optimality at the model translates into approximate optimality at the true situation. In the sciences, the main purpose of models is different: they are to assist our conceptual understanding, and to help with communication of ideas.

In traditional statistics there is no methodology for assessing the *adequacy* of a model. At best, traditional statistics can reject a model through a goodness-of-fit test — and Bayesian statistics cannot even do that. A non-rejected model is not necessarily adequate, and even more embarrassing, a rejected model sometimes may provide a perfectly adequate approximation.

#### 3. Classical Mathematical Statistics and Decision Theory

Classical mathematical statistics provides a clean theory under very restrictive assumptions, such as restricting the underlying models to exponential families, or the procedures to unbiasedness or invariance.

Decision theory clarified the classical views and reduced the dependence on restrictions. It also opened new areas, in particular optimal design theory (Kiefer 1959). But on the whole, decision theory was less successful than originally hoped. The two principal success stories are the Stein estimates (James and Stein 1961), relating to admissibility, and robustness (Huber 1964), relating to minimaxity.

# 4. Tukey's 1962 Paper

In his long 1962 paper "The Future of Data Analysis", while ostensibly talking
about his personal predilections, Tukey actually redefined the field of statistics.
Tukey's central theme was his emphasis on *judgment* (Section 7). At the same time,
he played down the importance of mathematical rigor and optimality (Sections
5 and 6). Possibly the most important issue worked out in his long and multifaceted paper was that there is more to theoretical statistics than mathematical

| 1        | statistics. This reminds one of Clausewitz (1832), who castigated the theorists of         | 1  |
|----------|--|----|
| 2        | military strategy of his time because they "considered only factors that could be          | 2  |
| 3        | mathematically calculated".  | 3  |
| 4        | In his paper, Tukey eschewed models. Why? Perhaps because in traditional                   | 4  |
| 5        | statistics models erroneously are considered as <i>substitutes</i> for the truth, rather   | 5  |
| 6        | than as <i>simplified approximations</i> . Note in particular his quote of Martin Wilk at  | 6  |
| 7        | the end of Section 4: "The hallmark of good science is that it uses models and             | 7  |
| 8        | 'theory' but never believes them".   | 8  |
| 9        | Tukey of course was not the first to question the role of models and of optimality.        | 9  |
| 10       | Statistical methods based on ranks and nonparametrics had become popular pre-              | 10 |
| 11       | cisely because they avoided dependence on uncertain models and were valid under            | 11 |
| 12       | weaker assumptions, even if they lacked the flexibility and wide applicability of the      | 12 |
| 13       | parametric approaches.   | 13 |
| 14       | But the problems with models and optimality go deeper. They have less to do                | 14 |
| 15       | with the idealized models <i>per se</i> , but more with the procedures optimized for them. | 15 |
| 10       |  | 10 |
| 10       |  | 10 |
| 10       | 5. Pitfalls of Optimality  | 10 |
| 20       |  | 20 |
| 20       | There are four basic pitfalls, into which mathematically optimal procedures can be         | 20 |
| 21<br>22 | trapped:   | 21 |
| 22       |  | 22 |
| 24       | (i) the Fuzzy Concepts Syndrome:   | 20 |
| 25       | sloppy translation of concepts into mathematics,   | 25 |
| 26       | (ii) the Straitjacket Syndrome:  | 26 |
| 27       | overly restrictive side conditions,  | 27 |
| 28       | (iii) the Scapegoat Syndrome:  | 28 |
| 29       | confuse the model with the truth,  | 29 |
| 30       | (iv) the Souped-Up Car Syndrome:   | 30 |
| 31       | optimize speed and produce a delicate gas-guzzler.   | 31 |
| 32       |  | 32 |
| 33       | These pitfalls affect distinct, very different aspects of statistics, namely: (i) con-     | 33 |
| 34       | cepts, (ii) procedures, (iii) models, and (iv) target functions. The list of course is     | 34 |
| 35       | not exhaustive. The pitfalls shall be discussed with the help of ten examples:             | 35 |
| 36       |  | 36 |
| 37       | Classical:   | 37 |
| 38       | (1) superefficiency  | 38 |
| 39       | (2) unbiasedness, equivariance   | 39 |
| 40       | (3) efficiency at the model  | 40 |
| 41       | Robustness:  | 41 |
| 42       | (4) asymptotics for finite $\varepsilon$   | 42 |
| 43       | (5) finite $n$ , finite $\varepsilon$  | 43 |
| 44       | (6) asymptotics for infinitesimal $\varepsilon$  | 44 |
| 45       | (7) optimal breakdown point  | 45 |
| 46       | Design:  | 46 |
| 47       | (8) optimal designs  | 47 |
| 48       | (9) regression design and breakdown  | 48 |
| 49       | Bayesian statistics:   | 49 |
| 50       | (10) Bayesian robustness   | 50 |
| 51       |  | 51 |

З

6. Examples from Classical Statistics

The three "classical" examples (1)-(3) nearly illustrate the first three pitfalls.

# 6.1. The Fuzzy Concepts Syndrome

Problems caused by the Fuzzy Concepts Syndrome mostly are relics from earlier development stages of statistical theories. In a conference on Directions for Mathematical Statistics, I had argued (Huber 1975b): "In statistics as well as in any other field of applied mathematics [...] one can usually distinguish (at least) three phases in the development of a problem. In Phase One, there is a vague awareness of an area of open problems, one develops *ad hoc* solutions to poorly posed questions, and one gropes for the proper concepts. In Phase Two, the right concepts are found [...]. In Phase Three, the theory begins to have a life of its own, [...] and its boundaries of validity are explored by leading it *ad absurdum*; in short, it is squeezed dry." In the 1970s there had been widespread anxiety about the future of mathematical statistics. As a deeper reason for this anxiety I had proposed the diagnosis that too many of the then current activities belonged to the later phases of Stage Three.

In the groping phase, somewhat reckless heuristics can be beneficial. The concepts inevitably are fuzzy, and correspondingly, they are translated into mathematics in a sloppy fashion. But recklessness, fuzziness and sloppiness should be cut down at the latest at the beginning of the squeezing phase (the "consolidation phase", in Erich Lehmann's terminology). Though, in the later phases it is tempting to concentrate on the mathematical formalism and to neglect a re-examination of its conceptual origins. And admittedly, even in the mathematical formalism, any attempts to eliminate sloppiness in its entirety will lead to an admirable, but non-teachable theory, as already Whitehead and Russell with their monumental *Principia Mathematica* (1910–13) have demonstrated.

In mathematical statistics, asymptotics is exceptionally prone to sloppiness. Details notoriously are not adequately elaborated. Indeed, the expression "asymptotic theory" itself is used misleadingly. In standard mathematical usage asymptotic theory ordinarily is concerned with asymptotic expansions. Statistics knows such expansions too (e.g. Edgeworth expansions), but mostly, "asymptotic theory" denotes what more properly should be called "limiting theory". A few examples follow.

- Remainder terms? With asymptotic expansions, the first neglected term gives an indication of the size of the error. In statistics, asymptotic results hardly ever are complemented by remainder terms, however crude. That is, whatever the actual sample size is, we never know whether an asymptotic result is applicable.
- What kind of asymptotics is appropriate? In regression, for example, we have n observations and p parameters. Should the asymptotics be for p fixed,  $n \rightarrow \infty$ , or for  $p/n \rightarrow 0$ , or for what?
- Order of quantifiers and limits? Usually, one settles on an order that makes proofs easy.
- **Example 1.** Perhaps the most illuminating case of the Fuzzy Concepts Syndrome has to do with superefficiency. There is a famous pathological example due to Hodges (see LeCam 1953). Assume that the observations  $(x_1, \ldots, x_n)$  are i.i.d. nor-

mal  $\mathcal{N}(\theta, 1)$ . Estimate  $\theta$  by

$$T_n = \bar{x}, \quad \text{if } |\bar{x}| \ge n^{-1/4}$$

$$T_n = \bar{x}/2$$
, if  $|\bar{x}| < n^{-1/4}$ .

Then  $T_n$  is consistent for all  $\theta$ , with asymptotic variance  $n^{-1}$  for  $\theta \neq 0$ , but  $\frac{1}{4}n^{-1}$  for  $\theta = 0$ . That is, the estimate  $T_n$  is efficient everywhere, but superefficient at 0. See Lehmann (1983, p. 405–408) for a discussion of various responses to the unpleasantness caused by Hodges' example.

Informally, asymptotic efficiency means that in large samples the variance of the estimate approaches the information bound, and this for all  $\theta$ . Everyday language is notoriously ambiguous about the order of the quantifiers. For example we may spell out asymptotic efficiency as:

(1) 
$$(\forall \varepsilon > 0) \ (\forall \theta) \ (\exists n_0) \ (\forall n > n_0) \ \{T_n \text{ is } \varepsilon \text{-efficent}\},\$$

where we define  $\varepsilon$ -efficiency by, say,

(2) 
$$\{T_n \text{ is } \varepsilon\text{-efficent}\} = \{E_\theta \left(n(T_n - \theta)^2\right) < 1/I(\theta) + \varepsilon\}.$$

But then, for any fixed n,  $T_n$  might be arbitrarily bad for some  $\theta$ . Therefore, since we do not know the true value of  $\theta$ , we never will know whether an estimate satisfying (1) is any good, however large n is. In other words: while the definition of asymptotic efficiency may be technically in order, it is conceptually inacceptable.

On closer inspection we conclude that the order of quantifiers in (1) does not correspond to our intuitive concept of asymptotic efficiency. An improved version is obtained by interchanging the second and third quantifiers:

(3) 
$$(\forall \varepsilon > 0) \ (\exists n_0) \ (\forall \theta) \ (\forall n > n_0) \ \{T_n \text{ is } \varepsilon \text{-efficent}\}.$$

It turns out that this version excludes superefficiency.

But version (3) still is negligent. Conceptually, unbounded loss functions are unsatisfactory. Technically, the awkward fact is that for very long-tailed distributions, the expectation in (2) may fail to be finite for all n and all "reasonable" estimators (i.e. for all estimators  $T_n$  whose value is contained in the convex hull of the observations, cf. Huber 1972, p. 1047), while the limiting distribution exists and has a finite variance. To obtain a definition of asymptotic efficiency working more generally we might rewrite (3) to

(4)

$$(\forall c > 0)(\forall \varepsilon > 0)(\exists n_0)(\forall \theta)(\forall n > n_0) \left\{ E_{\theta} \left( \left( \left[ \sqrt{n}(T_n - \theta) \right]_{-c}^{+c} \right)^2 \right) < 1/I(\theta) + \varepsilon \right\}.$$

Here,  $[x]_a^b = max(a, min(b, x))$ . Of course, (4) is not yet the final word; for example, we might want to replace the global uniform bound by a local one.

In my opinion Hodges' example should not be considered as a local improvement of the standard estimate, comparable to the James-Stein estimate, but rather as an ingenious spotlight on a conceptual inadequacy of the traditional formalization of asymptotic efficiency. This interpretation is not new. In particular, the crucial technical result, namely that one-sided locally uniform bounds suffice to prevent superefficiency, had been published in an abstract more than 40 years ago (Huber 1966). But I never had found a congenial outlet for the philosophical side of the result; it took the present symposium to provide one. З

6.2. The Straitjacket Syndrome

**Example 2**. Classical examples of the Straitjacket Syndrome, that is of overly restrictive side conditions on the procedures, are well known and do not need a detailed discussion. One is furnished by unbiasedness: unbiased estimates may not exist, or they may be nonsensical, cf. Lehmann (1983, p. 114). Other examples occur with invariance (more properly: equivariance): equivariant estimates may be inadmissible (Stein estimation).

6.3. The Scapegoat Syndrome

This subsection is concerned with excessive reliance on idealized models. The word "scapegoat" refers to the pious belief that the gods of statistics will accept the model as a substitute for the real thing.

As statisticians, we should always remember that models are simplified approximations to the truth, not the truth itself. Sometimes they are not even that, namely when they are chosen for ease of handling rather than for adequacy of representation; typical examples are conjugate priors in Bayesian statistics. The following eye-opening example gave rise to robustness theory.

**Example 3**. In 1914, Eddington had advocated the use of mean absolute deviations, against root-mean-square (RMS) deviations, as estimates of scale. Fisher (1920) objected and showed that for normal errors RMS deviations are 12% more efficient. Tukey (1960) then pointed out that for the contaminated normal error model

(5) 
$$F(x) = (1 - \varepsilon)\Phi(x/\sigma) + \varepsilon\Phi(x/(3\sigma))$$

mean absolute deviations are more efficient for all  $0.002 < \varepsilon < 0.5$ .

The unfortunate fact is that errors in real data typically are better approximated by a contamination model with a contamination rate ("gross error rate") in the range  $0.01 < \varepsilon < 0.1$ , than by the normal model.

The main lesson to be learnt from the Eddington–Fisher–Tukey example is that the standard normal error model may be quite accurate, especially in the center of the distribution. The problem is that the tail behavior of real data, to which the traditional estimates are highly sensitive, usually is rather indeterminate and difficult to model. The mistake of Fisher (and others) had been to treat the standard model as the exact truth.

We note a few conclusions from such examples:

- Optimality results put in evidence what can (and what cannot) be achieved in an *ideal* world.
- Notoriously, optimal procedures are unstable under *small deviations* from the ideal situation.
- The task thus is to find procedures that achieve *near optimality* under the ideal situation, but that are more stable under small deviation.

In 1964, I had begun to implement a program suggested by this under the name of robustness. The guiding ideas were: 

- Keep the optimality criterion (asymptotic variance, ...).
- Formalize small deviations ( $\varepsilon$ -contamination, ...).

– Find best sub-optimal procedures (best in a minimax sense).

The robustness notion I had adopted corresponds to Tukey's 1960 version. Though, this is not the unique interpretation of robustness occurring in the literature. In the 1970's, under Tukey's influence, there was a semantic shift, adopted by many, namely that the purpose of robustness was to provide procedures with a strong performance for a *widest possible* selection of heavy-tailed distribution.

But I still prefer the original 1960 version. In particular, I hold that robustness should be classified with *parametric* procedures, and that *local stability* in a neighborhood of the parametric model is the basic, overriding requirement.

7. Problems with Optimality in Robustness

Robustness had been designed to safeguard against pitfalls of optimal procedures. But optimal robustness is vulnerable to the very same pitfalls, and there are even a few new variants. The conceptual problem mentioned below in Example 4, and its solution described in Example 5, both have received less resonance in the robustness literature than they would have deserved. While the influence function without doubt is the most useful heuristic tool of robustness, one ought to be aware that optimality results based on it are no better than heuristic (Example 6).

- 7.1. Asymptotic Robustness for Finite  $\varepsilon > 0$
- **Example 4.** In the decision theoretic formalization of my 1964 paper I had imposed an unpleasant restriction on Nature by allowing only symmetric contaminations. The reason for this was that asymmetric contamination causes a bias term of the order O(1). Asymptotically, this bias then would overpower the random variability of the estimates (which typically is of the order  $O(n^{-1/2})$ ). Automatically, this would have led to the relatively inefficient sample median as the asymptotically optimal estimate. On the other hand, for the sample sizes and contamination rates of practical interest, the random variability usually is more important. Simultaneously, the symmetry assumption had permitted to extend the parameterization to the entire  $\varepsilon$ -neighborhood and thereby had made it possible to maintain a standard point-estimation approach.

The assumption of exact symmetry is repugnant, it violates the very spirit of robustness. Though, restrictions on the distributions are much less serious straitjackets than restrictions on the procedures (such as unbiasedness). The reason is that after performing optimization under symmetry restrictions, one merely has to check that the resulting asymptotically "optimal" estimate remains nearly optimal under more realistic asymmetric contaminations, see Huber (1981, pp. 104–106).

Curiously, people have worried (and still continue to worry!) much more about the symmetry straitiacket than about a conceptually much more serious problem. That problem is that 1% contamination has entirely different effects in samples of size 10 or 1000. Thus, asymptotic optimality theory need not be relevant at all for modest sample sizes and contamination rates, where the expected number of contaminants is small and may fall below 1. Fortunately, this question could be settled through an exact finite sample theory – see the following example. This theory also put to rest the problem of asymmetric contamination. 

# 7.2. Finite Sample Robustness for Finite $\varepsilon > 0$

**Example 5.** To resolve the just mentioned conceptual problem, one needs a finite sample robustness theory valid for finite  $\varepsilon > 0$ . Rigorous such theories were developed early on, see Huber (1965) for tests and Huber (1968) for estimation. The latter covers the same ground as the original asymptotic robustness theory, namely single parameter equivariant robust estimation. Gratifyingly, it leads to procedures that are qualitatively and even quantitatively comparable to the *M*-estimators obtained with the asymptotic approach.

This finite sample approach to robustness does not make any symmetry assumptions and thus also avoids the other objections that have been raised against asymptotic robustness theory. In particular, by aiming not for point estimates, but for minimax interval estimates, it bypasses the parameterization and asymmetry problems. Despite its conceptual importance, the finite sample theory has attained much less visibility than its asymptotic and infinitesimal cousins. I suspect the reason is that the approach through an unconventional version of interval estimates did not fit into established patterns of thought. In the following I shall sketch the main ideas and results; for technical details see Huber (1968).

Just as in the original asymptotic theory, we consider the one-parameter location problem and assume that the error distribution is contained in an  $\varepsilon$ -neigborhood of the standard normal distribution. The optimally robust finite sample estimator turns out to be an *M*-estimate *T* defined by

(6) 
$$\sum \psi(x_i - T) = 0,$$

where  $\psi(x) = [x]_{-k}^{k} = max(-k, min(k, x))$  for some k > 0. But instead of minimizing the maximal asymptotic variance, this estimator is optimal in the sense that it minimizes the value  $\alpha$  for which one can guarantee

(7) 
$$P\{T < \theta - a\} \le \alpha, \quad P\{T > \theta + a\} \le \alpha$$

for all  $\theta$  and all distributions in the  $\varepsilon$ -neighborhood.

We have three free parameters, n,  $\varepsilon$  and a. Interestingly, the characteristic parameter k of the  $\psi$ -function depends only on  $\varepsilon$  and a, but not on the sample size n. In (7), instead of minimizing  $\alpha$  for fixed a, we might alternatively minimize a for fixed  $\alpha$ . The asymptotic theory can be linked to these exact finite sample optimality results in several different fashions. In particular, if we let  $n \to \infty$ , but keep both  $\alpha$  and k fixed, then a and  $\varepsilon$  of the optimally robust estimates go to 0 at the rate  $O(n^{-1/2})$ . Conceptually,  $\varepsilon$ -neighborhoods shrinking at a rate  $O(n^{-1/2})$  make eminent sense, since the standard goodness-of-fit tests are just able to detect deviations of this order. Larger deviations should be taken care of by diagnostics and modeling, while smaller ones are difficult to detect and should be covered (in the insurance sense) by robustness.

# 7.3. Asymptotic Robustness for Infinitesimal $\varepsilon$

**Example 6**. Parametric families more general than location and scale are beyond the scope of the above approaches to robustness. Hampel proposed to attack them via gross error sensitivity: minimize asymptotic variance at the model, subject to a bound on the influence function (see Hampel 1974, and Hampel *et al.* 1986). This

approach is infinitesimal in nature and stays strictly at the parametric model. In essence, it is concerned only with the limiting case  $\varepsilon = 0$ .

Heuristically, it combines two desirable properties of robust estimates: good ef-ficiency at the model, and low gross error sensitivity. However, a bound on the latter at the model does not guarantee robustness (local stability in a neighborhood of the model), there are counter examples with L-estimates, see Huber (1981, pp. 71–72). Thus, the conceptual basis of this approach is weak. Even if it should yield robust procedures, we have no guarantee that they are approximately optimal for non-zero  $\varepsilon$ , and we thus have to pray to the statistical gods that they will accept an infinitesimal scapegoat. As a minimum, one ought to check the breakdown point of procedures constructed by this method.

There is a conceptually more satisfactory, but technically more complicated alternative approach via shrinking neighborhoods: while  $n \to \infty$ , let  $\varepsilon \to 0$  at the rate  $O(n^{-1/2})$ . This particular asymptotic theory had been motivated by the finite sample approach of Example 5. It was introduced by C. Huber-Carol in her thesis (1970) and later exploited by Rieder in several papers, culminating in his book (1994). The limiting results are comparable to those obtained with the infinitesimal approach, and like these, in the location parameter case they agree with those obtained in Example 4.

The principal heuristic appeal of the shrinking neighborhood approach is that in the location case it yields a sequence of estimates that have a well-defined optimality property for each n. We therefore can hope that in the general case it yields a sequence of estimates that are approximately optimal for non-zero  $\varepsilon$ . But to be honest, we have no way to check whether the heuristic arguments reliably carry beyond the location case. That is, we may run into a fifth pitfall: overly optimistic heuristics.

# 7.4. Optimal Breakdown Point

Hampel, at that time a student of Erich Lehmann, in his Ph.D. thesis (1968) had introduced the breakdown point by giving it an asymptotic definition. Conceptually, this may have been misleading, since the notion is most useful in small sample situations, see Donoho and Huber (1983). With large samples and high contamination rates you may have enough data to interpret the information contained in the contamination part. Therefore, rather than blindly using high breakdown point procedures, you may spend your efforts more profitably on an investigation of mixture models.

**Example 7.** All standard regression estimates, including the one based on least absolute deviations (the  $L_1$ -estimate, which generalizes the highly robust sample median), are sensitive to observations sitting at influential positions ("leverage points"). A single bad observation at an extreme leverage point may cause breakdown. Clearly, a higher breakdown point would be desirable. How large can it be made, and how large should it be? Via projection pursuit methods it is indeed possible to approach a breakdown point of 1/2 in large samples, provided the data are in general position (i.e., under the idealized, uncorrupted model no p rows of the *n*-by-p design matrix are linearly dependent, and thus any p observations give a unique determination of  $\theta$ ). This is a result of considerable theoretical interest.

Unfortunately, all estimators that try to optimize the breakdown point seem to run into the Souped-up Car Syndrome. The first among them was the LMSestimate (Hampel 1975, Rousseeuw 1984). 

The LMS- (Least Median of Squares) estimate of  $\theta$  modifies the least squares approach by minimizing the *median* instead of the *mean* of the squared residuals:

(8) 
$$\operatorname{median}\left\{(y_i - x_i^T \theta)^2\right\}.$$

If the data points are in general position, its breakdown point is  $([n/2] - p + 2)/n \rightarrow 1/2$ . But it has the following specific drawbacks:

- Its efficiency is low: the dispersion of the estimate decreases at the rate  $n^{-1/3}$ , instead of  $n^{-1/2}$ .

- Its computational complexity increases exponentially with p.

– What if the points are not in general position?

My conclusion is that an asymptotic theory for large p and n does not make much sense under such circumstances, and a small sample theory is not available.

S-estimates were introduced by Rousseeuw and Yohai (1984) to overcome some of these shortcomings. They estimate  $\theta$  by minimizing a suitable robust *M*-estimate of the scale of the residuals. Under suitable regularity conditions their breakdown point also approaches 1/2 in large samples, and they reach a high efficiency at the ideal model, with a dispersion converging at the rate  $n^{-1/2}$ . Unfortunately, *S*estimators suffer from a serious flaw which probably cannot be removed, namely that uniqueness and continuity can only be proved under certain conditions, see Davies (1993, Section 1.6).

Moreover, Davies (*ibid.*) points out that all known high breakdown point estimators of regression are inherently unstable. Paradoxically, it thus seems that in order to achieve an optimal regression breakdown point we may have to sacrifice robustness.

# 8. Design Issues

# 8.1. Optimal Designs

**Example 8.** Assume that the task is to fit the best possible straight line to data originating from an exactly linear function. Then the optimal regression design puts all observations on the extreme points of the segment where observations can be made.

However, a possibly more realistic version of this task is to fit the best straight line to an *approximately* linear function. In either case, one would want to make something like the expectation of the integrated mean square error as small as possible. Of course, usually one does not know much about the small deviations from linearity one might have to cope with (and does not care about them, so long as they are small). Already Box and Draper (1959, p. 622) had recognized the crux of the situation and had pointed out: "The optimal design in typical situations in which both variance and bias occur is very nearly the same as would be obtained if *variance were ignored completely* and the experiment designed so as to *minimize bias alone*."

In other words, the "naive" design, which distributes the observations evenly
over the accessible segment, in such a situation should be very nearly optimal,
since it minimizes the integrated squared bias of the fit. Apart from that, it has an
advantage over the optimal design since its redundancy allows to check the linearity
assumption.

These aspects have been made precise in a minimax sense by Huber (1975a). The most surprising fact emerging from this study was: there is a range where the

З

З



FIG 1. Optimal design fit (based on the theoretically optimal design) and best linear fit (minimizing the integrated squared error), to a not-quite linear function. Observational errors are neglected in this figure.

deviation from linearity is slight enough to stay below statistical detectability, yet large enough so that the "naive" design will outperform the "optimal" design and give a better linear approximation to the true function. Even though the effects are much less dramatic than in Example 3, we evidently have run here into another example of the Scapegoat Syndrome.

8.2. Regression Design and Breakdown

**Example 9.** In higher dimensions, generalizing the preceding example, an optimal linear regression design would place an equal number of m observations onto each of the (p+1) corners of a p-dimensional simplex. Technically, such a design is optimal, but again, it lacks redundance.

For such a design the best possible breakdown point is

(9) 
$$\varepsilon^* = \lceil m/2 \rceil / (m(p+1)) \approx 1/(2(p+1)).$$

This breakdown point is attained by the  $L_1$ -estimate (calculate the median at each corner). The so-called high-breakdown point LMS- and S-estimates cannot do any better.

But already an arbitrarily small jittering of the design points will bring them into general position. Then the breakdown point of LMS and S is close to 1/2. How can this happen?

On closer inspection we see that the high breakdown point of *LMS*- and *S*estimates is achieved by extrapolation: at each corner, you put more faith in the value extrapolated from the *mp* observations clustering near the far-away other *p* corners, than in the *m* local values. The fitted hyperplane thus not only loses efficiency, but becomes sensitive to small errors affecting a majority of the observations, such as rounding.

The conclusion is that high breakdown point regression is not necessarily robust.
We have a clear case of the Souped-up Car Syndrome: both extremes, optimal

design and optimal breakdown point, lead to estimates with undesirable properties, and a compromise is called for. A quantitative, *design-dependent* theory of robust regression would seem to be needed. The customary assumption underlying all high breakdown point regression theories in the published literature, namely that the regression carrier is a random sample from a suitable multi-dimensional continuous distribution, in my opinion is much too narrowly conceived.

# 9. Bayesian Statistics

**Example 10.** What is Bayesian robustness? Bayesian statistics has a built-in problem with the Scapegoat Syndrome, that is, with over-reliance on the model; this problem becomes acute in connection with robustness. By definition, Bayes procedures are optimal for the chosen model, consisting of a prior  $\alpha(\theta)$  and a family of conditional densities  $f(x, \theta)$ . Instability, and conversely robustness, thus are properties of the *model*. This was emphasized in 1978 by George Box in an illuminating, facetious but profound oral interchange with John Tukey at an ARO meeting on Robustness in Statistics (Launer and Wilkinson, 1979). Box maintained that, after all, *he* had invented the term (see Box 1953), and that he could define it as he pleased, and that in his opinion robustness was to be achieved by choosing a proper model, not by tampering with the data (by trimming or Winsorizing) as Tukey was wont to do. He did not elaborate on how to choose such a model.

The philosophical problem of Bayesian statistics is that it is congenitally unable to separate the model, the underlying true situation, and the statistical procedure. It acts as if the model were exactly true, and it then uses the corresponding optimal procedure. A fundamentalist Bayesian, for whom probabilities exist only in the mind, will not be able to see that there is a problem of the Scapegoat type; it takes a pragmatist like George Box to be aware of it.

I shall now attempt to sketch a way around this Bayesian Scapegoat Sydrome. The central question is: what is a *robust model?* Ad hoc parametric supermodels, which sometimes are advertised as the Bayesian approach to robustness, do not guarantee robustness. There are no reliable guidelines to select such models, and the resulting procedures may suffer from instabilities.

If we proceed pragmatically, then, as a minimum requirement, the statistical conclusions from the model ought to be insensitive to occasional outliers. Sensitivity studies  $\dot{a} \, la$  Berger, that is: admit that the specifications are inaccurate and find the range of implied conclusions (see Wolpert 2004, p. 212), may reveal the presence of outliers: if there are outliers, small changes in the tails of the model  $f(x,\theta)$  can produce large effects. Also, they may reveal conflicts between the prior and the observational evidence: if the observational evidence points to a  $\theta$  far from the center of the prior, small changes in the tails of the latter can produce large effects. Thus, if a sensitivity analysis shows that the range of implied conclusions is narrow, any model in the uncertainty range will do. If not, we better choose a robust model. But then, why not choose a robust model right away? Unfortunately, sensitivity studies do not help us find a robust model.

The following is a proposal for an informal portmanteau definition of robustness, covering both Bayesian and non-Bayesian statistics:

Uncertain parts of the evidence should never have overriding influence on the final conclusions.

This is supposed to apply not only to questionable data (outliers), but also to uncertainties in the model densities  $f(x, \theta)$  and to uncertainties in the prior  $\alpha(\theta)$ ,

and even to vagueness in the specification of the target loss function.

How to implement such a loose definition? The first two of the above four requirements are interconnected and tricky to separate: insensitivity to dubious data features (gross errors), and insensitivity to uncertain model specifications. I claim that the following implementation should do the job for both aspects: Choose a model  $f(x, \theta)$  within the uncertainty range, such that the conclusions are insensitive to gross errors. This has to be made precise.

The mode of the posterior density solves

(10)  $\alpha'(\theta)/\alpha(\theta) + \sum f'(x_i,\theta)/f(x_i,\theta) = 0,$ 

where the prime denotes the derivative with respect to  $\theta$ . For a flat prior, the mode of the posterior coincides with the maximum likelihood estimate.

As Freedman (1963) has expressed it, there is a "striking and mysterious fact", namely that asymptotically, Bayes and M.L. estimates behave similarly: they not only have the same asymptotic distribution, but if the true underlying distribution belongs to the parametric family, the Bayesian posterior distribution, centered at the M.L. estimate and scaled by  $n^{-1/2}$ , is asymptotically normal and coincides with the asymptotic distribution of the M.L. estimate, centered at the true  $\theta$  and also scaled by  $n^{-1/2}$ . See also LeCam (1957); the result apparently goes back to Bernstein and von Mises.

Thus, if we are willing to adopt the infinitesimal approach via gross error sen-sitivity (see Example 6), asymptotic robustness ideas should carry over from non-Bayesian *M*-estimates. Though, Hampel's approach through gross error sensitivity does not apply without some caveats, since it does not automatically lead to  $\psi$ -functions that are logarithmic derivatives of probability densities (which is a nec-essary side condition in the Bayes context — another example of a straitjacket). Finite  $\varepsilon$ -neighborhoods need somewhat more work. Assume that the M- and Bayes estimates both are calculated on the basis of the least favorable density (instead of the unknown true underlying distribution, which is supposed to lie anywhere in the given  $\varepsilon$ -neighborhood). Then, the M- and Bayes estimates still have the same asymptotically normal distribution, but the equivalence with the asymptotic pos-terior is lost. Though, in the one-dimensional location case it can be shown that the asymptotic variance of the posterior then lies between the asymptotic variance of the M-estimate and the upper bound for that variance obtained from the least favorable distribution (to be published in the forthcoming 2nd edition of my Robust Statistics). — As an amusing aside on the subject of pitfalls, I might mention that the usual applications in econometrics of one the formulas relevant in this context (the so-called "sandwich formula") go so far beyond its original intention that they deserve an honorable mention in the category of overly optimistic heuristics, see Freedman (2006). 

In short, the heuristic conclusion, deriving from hindsight based on non-Bayesian robustness, thus is that f'/f ought to be bounded. In (10) the prior acts very much like a distinguished additional observation. Thus, in analogous fashion, also  $\alpha'/\alpha$ ought to be bounded. In both cases, the bounds should be chosen as small as feasible. Ordinarily, these bounds are minimized by the least informative distributions, with Fisher information used as measure of information. Thus, a possible optimization goal can be expressed:

 A most robust Bayesian model can be found by choosing  $\alpha$  and f to be least informative within their respective (objective or subjective) uncertainty ranges.

З

For all practical purposes this is the same recipe as the one applying to the non-Bayesian case. But like there, it is difficult to implement once one wants to go beyond the location case. And if it is adopted overly literally, we might even get trapped in one the pitfalls of optimality.

10. Concluding Remarks

In the 1970s statistical theory clearly had been in the third, consolidation or "squeezing" phase of the development cycle. At present, we seem to have entered a new cycle and seem to be in the middle of a new "groping" phase, trying to get conceptual and theoretical handles on challenging problems involving extremely large and complexly structured data sets.

I hope that this time the laxness of the groping phase will be eliminated in time, and will not be cemented into place during the consolidation phase. Perhaps it may help to keep in mind the following aphorisms on optimality and optimization. They are not new, they are re-iterating sentiments already expressed by Tukey in 1962. Those sentiments subsequently had been studiously ignored by most statisticians. I hope that this time they will fare better.

- Optimality *results* are important: they show what can (and what cannot) be achieved under ideal conditions, and in particular they show whether a given procedure still has worthwhile potential for improvement.
  - Optimal *procedures* as a rule are too dangerous to be used in untempered form.
- Beware of sloppy asymptotics.
- Never confuse the idealized model with the truth.
- Do not optimize one aspect to the detriment of others.
- There are no clear-cut rules on how the tempering of optimal procedures should be done compromises are involved, and one must rely on human judgment. But if one insists on a mathematical approach, minimizing Fisher information within a subjective uncertainty range often will do a good job, both for Bayesians and non-Bayesians.

#### References

- [1] Box, G. E. P. (1953). Non-normality and tests on variances. Biometrika, 40, 318–335. [2] BOX, G. E. P. and DRAPER, N. R. (1959). A basis for the selection of a response surface design. J. Amer. Statist. Ass., 54, 622–654. CLAUSEWITZ, C. VON (1832). Vom Kriege. 19th edition (1991). Dümmler Verlag, Bonn. - (1984). On War. Edited and translated by M. HOWARD and P. PARET. Princeton University [4]Press, Princeton, NJ. DAVIES, P. L. (1993). Aspects of robust linear regression. The Annals of Statistics, 21, 1843–1899. [6] DONOHO, D. L. and HUBER, P. J. (1983). The notion of breakdown point. In A Festschrift for Erich L. Lehmann. Bickel, P. J., Doksum, K. A., and Hodges, J. L. (Eds.) Wadsworth, Belmont, CA. [7] EDDINGTON, A. S. (1914). Stellar Movements and the Structure of the Universe. Macmillan, London. [8] FISHER, R. A. (1920). A mathematical examination of the methods of determining the accuracy of an observation by the mean error and the mean square error. Monthly Not. Roy. Astron. Soc., 80, 758-770. [9] (1922). On the mathematical foundation of theoretical statistics. Philos. Trans. Roy. Soc. London, Ser. A, 222, 309-368.
- [10] FREEDMAN, D. A. (1963). On the asymptotic behavior of Bayes' estimates in the discrete case. *Annals of Mathematical Statistics*, 34, 1386–1403.

   [11] (2020) On the second with the South is Estimated in the Bayes''. The second seco
- [11] (2006). On the so-called "Huber Sandwich Estimator" and "Robust Standard Errors". The
   American Statistician, 60, 209–302.

imsart-coll ver. 2008/08/29 file: Huber.tex date: March 25, 2009

| 1        | [12]         | HAMPEL, F. R. (1968). Contributions to the theory of robust estimation. Ph. D. Thesis. University   | 1        |
|----------|--------------|---|----------|
| 2        |              | of California, Berkeley.  | 2        |
| 3        | [13]         | (1974). The influence curve and its role in robust estimation. J. Amer. Statist. Ass., 62,  | 3        |
| 4        | [1.4]        | (1075) Periord location parameters: Polyst concepts and methods Prog. (0th Session I. S.  | 4        |
| -        | [1.4]        | I. Warsaw 1975. Bull. Int. Statist. Inst., 46. Book 1, 375–382.   | 5        |
| 5        | [15]         | HAMPEL, F. R., RONCHETTI, E. M., ROUSSEEUW, P. J. and STAHEL, W. A. (1986). Robust Statistics.  | 5        |
| 6        |              | The Approach Based on Influence. John Wiley, New York.  | 6        |
| 7        | [16]         | HUBER, P. J. (1964). Robust estimation of a location parameter. Ann. Math. Statist., <b>35</b> , 73–101.  | 7        |
| 8        | [17]<br>[18] | (1965). A robust version of the probability ratio test. Ann. Math. Statist., <b>36</b> , 1753–1758.   | 8        |
| 9        | [19]         | —(1968). Robust confidence limits. Z. Wahrscheinlichkeitstheorie Verw. Gebiete, 10, 269–  | 9        |
| 10       |              | 278.  | 10       |
| 11       | [20]         | —(1972). Robust statistics: A review. Ann. Math. Statist., 43, 1041–1067.   | 11       |
| 12       | [21]         | ——(1975a). Robustness and designs. In A Survey of Statistical Design and Linear Models, J.  | 12       |
| 13       | [22]         | (1975b), Application vs. abstraction: the selling out of mathematical statistics. Suppl. Adv.   | 13       |
| 14       |              | Appl. Prob., 7, 84–89.  | 14       |
| 15       | [23]         | ——(1981). Robust Statistics. John Wiley, New York.  | 15       |
| 10       | [24]         | HUBER-CAROL, C. (1970). Etude asymptotique de tests robustes. Ph.D. Thesis. Eidgen, Technis-  | 10       |
| 10       | [25]         | Che Hochschule, Zurich.<br>JAMES W and STEIN C (1961) Estimation with quadratic loss Proc. Fourth Berkeley  | 10       |
| 17       | [20]         | Symp. Math. Statist. Prob. I, 311–319.  | 17       |
| 18       | [26]         | KIEFER, J. (1959). Optimum experimental designs. J. Roy. Statist. Soc. Ser. B 21, 272–319.  | 18       |
| 19       | [27]         | LAUNER, R. L. and WILKINSON, G. N. (Eds.) (1979). Proc. ARO Workshop on Robustness in   | 19       |
| 20       | [00]         | Statistics, April 11–12, 1978, Academic Press, New York.  | 20       |
| 21       | [28]         | LECAM, L. (1953). On some asymptotic properties of maximum likelihood estimates and related<br>Bayes' estimates Univ Calif Publ Statist 1 277–330 | 21       |
| 22       | [29]         | (1957). Locally asymptotically normal families of distributions. Univ. Calif. Publ. Statist.,   | 22       |
| 23       |              | <b>3</b> , 37–98.   | 23       |
| 24       | [30]         | LEHMANN, E. L. (1959). Testing Statistical Hypotheses. John Wiley, New York.  | 24       |
| 25       | [31]         | (1983). Theory of Point Estimation. John Wiley, New York.   | 25       |
| 20       | [32]         | hypotheses. Philos. Trans. Roy. Soc. London. Ser. A, 231, 289–337.  | 20       |
| 20       | [33]         | RIEDER, H. (1994). Robust Asymptotic Statistics. Springer, Berlin.  | 20       |
| 27       | [34]         | ROUSSEEUW, P. J. (1984). Least median of squares regression. J. Amer. Statist. Assoc., 79, 871–   | 27       |
| 28       | [95]         | 880.  | 28       |
| 29<br>30 | [35]         | and Nonlinear Time Series Analysis. Franke, J., Härdle, W. and Martin, R. D. (Eds.) Lecture   | 29<br>30 |
| 21       | [0.0]        | Notes in Statistics, <b>26</b> . Springer, New York.  | 21       |
| 31       | [36]         | TUKEY, J. W. (1960). A survey of sampling from contaminated distributions. In <i>Contributions</i>  | 31       |
| 32       | [37]         |   | 32       |
| 33       | [38]         | WALD, A. (1950). Statistical Decision Functions. Wiley, New York.   | 33       |
| 34       | [39]         | WHITEHEAD, A. N. and RUSSELL, B. (1910–13). Principia Mathematica, 3 vols. Cambridge Uni-   | 34       |
| 35       | [40]         | Versity Press.  | 35       |
| 36       | [40]         | WOLPERT, R. L. (2004). A conversation with James O. Berger. Statistical Science, 19, 205–218.   | 36       |
| 37       |              |   | 37       |
| 38       |              |   | 38       |
| 39       |              |   | 39       |
| 40       |              |   | 40       |
| <br>/1   |              |   | -10      |
| 41       |              |   | 41       |
| 42       |              |   | 42       |
| 43       |              |   | 43       |
| 44       |              |   | 44       |
| 45       |              |   | 45       |
| 46       |              |   | 46       |
| 47       |              |   | 47       |
| 48       |              |   | 48       |
| 49       |              |   | 49       |
| 50       |              |   | 50       |
| 51       |              |   | E1       |
| 01       |              |   | 51       |

imsart-coll ver. 2008/08/29 file: Huber.tex date: March 25, 2009

IMS Collections Vol. 0 (2009) 46-65 © Institute of Mathematical Statistics, 2009 arXiv: math.PR/0000018

# Proportional Hazards Regression with Unknown Link Function

# Wei Wang<sup>1</sup>, Jane-Ling Wang<sup>2</sup> and Qihua Wang<sup>3</sup>

Harvard Medical School and Brigham and Women's Hospital; University of California, Davis; Chinese Academy of Science and The University of Hong Kong Abstract: Proportional hazards regression model assumes that the covariates affect the hazard function through a link function and an index which is a linear function of the covariates. Traditional approaches, such as the Cox proportional hazards model, focus on estimating the unknown index by assuming a known link function between the log-hazard function and covariates. A linear link function is often employed for convenience without any validation. This paper provides an approach to estimate the link function, which can then be used to guide the choice of a proper parametric link function. This is accomplished through a two-step algorithm to estimate the link function and the effects of the covariates iteratively without involving the baseline hazard estimate. The link function is estimated by a smoothing method based on a local version of partial likelihood, and the index function is then estimated using a full version of partial likelihood. Asymptotic properties of the non-parametric link function estimate are derived, which facilitates model checking of the adequacy of the Cox Proportional hazards model. The approach is illustrated through a survival data and simulations. Contents 2.12.2Main Results 2.3Simulation Studies Conclusion and Future Research Acknowledgement 1. Introduction Proportional hazards regression model has played a pivotal role in survival analysis since Cox proposed it in 1972. Let T represent survival time and Z its associate <sup>1</sup>Harvard Medical School and Brigham and Women's Hospital, Boston, USA, email: wwang4@ partners.org <sup>2</sup>University of California, Davis, USA, Apartado Postal 402, email: wang@wald.ucdavis.edu <sup>3</sup>Chinese Academy of Science, Beijing, China, email: ghwang@amss.ac.cn and University of Hongkong, Hongkong, email: qhwang@hku.hk AMS 2000 subject classifications: Primary 62G05; secondary 62N02 Keywords and phrases: Partial likelihood, local partial likelihood, Nonparametric smoothing, Dimension reduction  1 covariate vector. Under the proportional hazards model, the hazard function for T, 2 given a particular value z for the covariate Z, is defined as

4 (1) 
$$\lambda\{t \mid z\} = \lambda_0(t) \exp\{\psi(\beta_0^T z)\},$$

where  $\lambda_0(t)$  is an unknown baseline hazard function corresponding to  $z = (0, \dots, 0)$ , and  $\psi(\cdot)$  is called the link function with  $\psi(0) = 0$ . With fully specified link function  $\psi$ , the partial likelihood method was introduced in [4, 5] to estimate the regression parameters,  $\beta_0$ , with the option to accommodate censored data. The most common choice for  $\psi$  is the identity function, which corresponds to the time-honored Cox model. In reality the link function is unknown and needs to be estimated. This is especially useful to validate a preferred choice, as an erroneous link function could dramatically distort risk assessment or interpretation of odds ratios. When the link function is known, such as in the Cox model, model (1) is a special case of the transformation model first proposed in [7] and subsequently studied in [6], [3] etc. Our goal in this paper is to consider model (1) with an unknown link function. This problem was first studied in an unpublished Ph.D. thesis [19]. However, the procedure there was less efficient and we propose an improved estimate, studying its asymptotic properties. 

Previous work focuses on the special case when the covariate is one-dimensional, or equivalently when  $\beta$  is known in (1). Under this special one-dimensional case, a local partial likelihood technique in [18] and a variation of the local scoring algo-rithm of [12] can be used to estimate the unknown link function in (1). Gentleman and Crowley [10] proposed a local version of the full likelihood instead of partial likelihood by alternating between estimating the baseline hazard function and estimating the covariate effects. The local likelihood methods in these papers were based on data whose covariate values fall in a neighborhood of the targeted location. Fan, Gijbels and King [8] used instead a local polynomial method to approximate the local partial likelihood, and derived rigorous asymptotic results for their ap-proach. Simulation studies there showed that the local partial likelihood method is comparable to the local full likelihood method in [10]. Two spline approaches have also been considered, with smoothing splines resulting from a penalized partial likelihood in O'Sullivan [16] and regression splines from [17]. 

While the aforementioned approaches can be easily extended to q-dimensional covariates by estimating a multivariate unknown link function  $\psi(z_1, \dots, z_q)$ , such nonparametric approaches are subject to the curse of dimensionality and may not be suitable for  $q \geq 3$ . Moreover, the resulting model would be different from model (1), which has the attractive dimension reduction feature that the covariate information is succinctly summarized in a single index and is a nonparametric extension of the Cox proportional hazards model. Model (1) could also be used as an exploratory tool to guide the choice of a suitable parametric link function. 

A two-step iterative algorithm to estimate the link function and the covariate effects is proposed in Section 2. In the first step, an initial estimate of the regression parameter  $\beta$  is plugged in model (1) so that the link function can be estimated by a smoothing method based on a local version of partial likelihood ([4, 5]). The second step involves updating the regression parameters using the full partial likelihood with the estimated link function in step 1 inserted. These two steps will be iterated until the algorithm converges. Asymptotic results for the link estimators are stated in Section 2. In particular, Theorem 2 provides the building blocks to check the link function and inference for the individual risk,  $\psi(\beta^T z)$ . It also reveals that the nonparametric estimate of the link function is as efficient as the one for model 

(1) but with a known regression parameter  $\beta$ . Thus, there is no efficiency loss to estimate the link function even if  $\beta$  is unknown in our setting. This is also reflected in the simulation studies in Section 3. The approach in Section 2 is further illustrated in Section 4 through a data set from the Worcester Heart Attach Study. All the proofs of the main results are relegated to an appendix.

We remark here that [15] also studied model (1) with a different approach. They assumed that the link function is in a finite dimensional subspace spanned by polynomial spline bases functions and the dimension of this subspace is known. This leads to a flexible parametric model where the spline coefficients corresponding to the link functions and  $\beta$  can then be estimated directly through traditional partial likelihood approaches. While this has the benefit of simplicity as everything is in the parametric framework, it tends to underestimate the standard errors of the estimates. Two sources of bias arise, one derives from the fact that in reality the number of spline bases depends on the data and is a function of the sample size, so the standard errors are underestimated by the simple parametric inference. In addition, the link estimation might be biased, as in theory an infinite number of spine bases might be required to span the unknown link function. These biases could significantly affect the asymptotic results. In contrast, our approach provides correct asymptotic theory and efficient estimation of the link function.

#### 2. Estimation Procedure and Main Results

Ν

Since there are three different unknown parameters,  $\lambda_0(\cdot)$ ,  $\psi(\cdot)$  and  $\beta$  in model (1), we need to impose some conditions to ensure identifiability. To identify  $\lambda_0$ , it suffices to set  $\psi(v) = 0$  at some point v, a common choice is v = 0. Since only the direction of  $\beta$  is identifiable if  $\psi$  is unknown, we assume that  $\|\beta\|=1$  (here  $\|\cdot\|$ represents the Euclidean norm) and that the sign of the first component of  $\beta$  is positive. As for the sampling plan, we assume an independent censoring scheme, in which the survival time T and censoring time C are conditionally independent, given the covariate vector Z. Let X = min(T, C) be the observed event-time and  $\Delta = I\{T \leq C\}$  be the censoring indicator. The data  $\{X_i, Z_i, \delta_i\}$  is an i.i.d. sample of  $\{X, Z, \Delta\}$ . We use the notation  $t_i < \cdots < t_N$  to denote the N distinctive ordered failure times, and (j) to denote the label of the item failing at time  $t_i$ . The risk set at time  $t_j$  is denoted by  $\mathcal{R}_j = \{i : X_i \ge t_j\}.$ 

For a fixed parametric value  $\beta$ , one can estimate the link function  $\psi(\cdot)$  by any smoothing method, such as those cited in Section 1 when  $\beta$  is assumed known. We adopt the local partial likelihood approach in [8] and assume, for a given point v, that the p-th order derivative of  $\psi(v)$  at point v exists. A Taylor expansion for  $\beta^T Z$ in a neighborhood of v then yields,

(2)

$$\psi(\beta^{T}Z) \approx \psi(v) + \psi'(v)(\beta^{T}Z - v) + \dots + \frac{\psi^{(p)}(v)}{p!}(\beta^{T}Z - v)$$

$$= \psi(v) + (\beta^T \mathbf{Z})^T \gamma(v),$$

where  $\gamma(v) = \{\psi'(v), \cdots, \psi^{(p)}(v)/p!\}^T$  is the *p*-dimensional vector associated with the derivatives of  $\psi$  and  $\beta^T \mathbf{Z} = \{\beta^T Z - v, \cdots, (\beta^T Z - v)^p\}^T$ .

Let K be a kernel function, h be a bandwidth, and define  $K_h(u) = h^{-1}K(u/h)$ . Applying kernel weights to the logarithm of the global partial likelihood

> ſ *'*)

$$\sum_{j=1} \psi(\beta^T Z_{(j)}) - \log\left\{\sum_{i \in \mathcal{R}_j} \exp\left\{\psi(\beta^T Z_{(j)})\right\}\right\}$$
<sup>50</sup>  
51

P

and replacing  $\psi(\beta^T z)$  by the local approximation in (2), we arrive at (similarly to [8]) the local version of the log partial likelihood:

$$\sum_{j=1}^{N} K_h \left( \beta^T Z_{(j)} - v \right) \left| \left( \beta^T \mathbf{Z}_{(j)} \right)^T \gamma(v) \right|$$

(3) 
$$-\log\left\{\sum_{i\in\mathcal{R}_{j}}\exp\left\{\left(\beta^{T}\mathbf{Z}_{i}\right)^{T}\gamma(v)\right\}K_{h}\left\{\beta^{T}Z_{i}-v\right\}\right\}\right],$$

where  $\beta^T \mathbf{Z}_i$  and  $\beta^T \mathbf{Z}_{(i)}$  are defined as  $\beta^T \mathbf{Z}$  with  $\mathbf{Z}$  replaced by  $\mathbf{Z}_i$  and  $\mathbf{Z}_{(i)}$  respectively. It can be shown that the local log partial likelihood in (3) is strictly concave with respect to  $\gamma(\cdot)$ , so for a fixed  $\beta$ , it has a unique maximizer with respect to  $\gamma$ . Let  $\hat{\gamma}(v)$  be the local partial likelihood estimate with  $\hat{\gamma}_k(v)$  denoting its k-th component, then  $\psi^{(k)}(v)$  can be estimated by  $\hat{\psi}^{(k)}(v) = k! \hat{\gamma}_k(v)$ , for  $k = 1, \dots, p$ . In principle, one could maximize (3) with respect to both  $\beta$  and  $\gamma$ , and this corresponds to maximizing the real local log likelihood. But we choose to maximize (3)only with respect to  $\gamma$  for a fixed estimated value of  $\beta$ , and this corresponds to maximizing a pseudo local log likelihood as the true  $\beta$  in (3) is replaced by an estimate. There are two reasons for our choice. First (3) is concave in  $\gamma$ , but not necessarily in  $\beta$ . Second, maximizing with respect to both parameters is probably not worth the additional computational cost, as the local likelihood procedures mainly serve as a smoother and the choice of the smoother is usually not crucial.

To estimate the link function, we use

where  $\hat{\psi}'(v)$  is the first component of  $\hat{\gamma}(v)$  at the last iteration step. There are several ways to approximate this integral, such as the trapezoidal rule or Gaussian quadrature. For computational simplicity, we apply the trapezoidal rule in the simulation studies, as suggested in [18], and this appears to be satisfactory.

# 2.1. Algorithm and Computational Issues

The procedure described in the previous subsection requires a certain choice of  $\beta$ in equation (2). This can be done either independently or iteratively as once an estimate of  $\psi$  is obtained, one can then estimate  $\beta$  through the global partial likelihood. An iterative algorithm, as shown below, can be established by alternatingly updating the estimates for  $\beta$  and  $\psi$ . Such an iteration procedure may improve the link estimate as a better estimate of  $\beta$  will lead to a better estimate of  $\psi$ . **Step 1**. (a) Assign a nonzero initial value to  $\beta$ , and call it  $\hat{\beta}$ .

(b) For a given v, plug  $\hat{\beta}$  into the pseudo log local partial likelihood and maximize

$$\sum_{i=1}^{N} K_h \{ \hat{\beta}^T Z_{(j)} - v \} \cdot \left[ [\hat{\beta}^T \mathbf{Z}_{(j)}]^T \gamma(v) \right]^{-1}$$

 $\overline{j=1}$ 

$$-\log\left\{\sum_{i\in\mathcal{R}_{i}}\exp\{[\hat{\beta}^{T}\mathbf{Z}_{i}]^{T}\gamma(v)\}K_{h}\left\{\hat{\beta}^{T}Z_{i}-v\right\}\right\}\right]$$

with respect to  $\gamma(v)$  to get the estimate  $\hat{\gamma}(v)$ .

З
(c) Obtain the values of  $\hat{\gamma}(v)$ , for  $v = \hat{\beta}^T Z_i, i = 1, \cdots, n$ .

З

(d) Apply the trapezoidal rule to obtain  $\{\hat{\psi}(\hat{\beta}^T Z_i) : i = 1, \cdots, n\}$ .

**Step 2**. Plug  $\hat{\psi}(\cdot)$  into the log (global) partial likelihood

$$l_G(\beta, \hat{\psi}) = \sum_{j=1}^{N} \bigg[ \hat{\psi}(\beta^T Z_{(j)}) - \log \bigg\{ \sum_{i \in \mathcal{R}_j} \exp \{ \hat{\psi}(\beta^T Z_i) \} \bigg\} \bigg],$$

and maximize it with respect to  $\beta$  to update the estimate  $\hat{\beta}$ . We use the angle between two estimated  $\hat{\beta}$  at two consecutive iterations as the convergence criterion.

**Remark 1.** The Newton-Raphson method is used to find the estimators in Step 1 and 2. The initial value of  $\beta$  can be set in different ways but cannot be zero as a nonzero value is needed in step 1 to estimate the link function. However, this restriction does not exclude the final  $\beta$ -estimate to be zero or close to zero. A simple choice is to fit the usual Cox model and use this estimator in the first step. To accelerate the computation, one can also use alternative estimates as described below.

**Remark 2.** It is possible to accelerate the computation by using a  $\sqrt{n}$ -consistent initial estimator, as Theorems 1 and 2 below imply that no iteration is required for the link estimate and that it will converge efficiently at the usual nonparametric rate. Namely, the link function can be estimated with the same efficiency as when  $\beta$  is known. In practice, we find that one iteration helps to improve the numerical performance but further iteration is usually not necessary. There are two choices for a  $\sqrt{n}$ -consistent initial estimator, one is the estimator in [2] that extends the sliced inverse regression (SIR) approach to censored data. Specifically, this approach requires a certain design condition as listed in (2.3) there but has the advantage that it leads to a  $\sqrt{n}$ -consistent estimator for  $\beta$  without the need to estimate the link function. Another initial estimate which does not rely on the design conditions (2.3) in [2] is provided in a Ph. D. thesis [19]. Specifically, this involves replacing the  $\psi$  function in step 2 above by its local version (2), which leads to the cancelation of the term  $\psi$  and results in a version of log (global) partial likelihood that involves only the derivative  $\gamma(v)$  of  $\psi$  but not  $\psi$  itself. Thus, Step 2 above is replaced by

**Step 2\***. Maximize the following approximate log (global) partial likelihood with respect to  $\beta$ :

$$\sum_{j=1}^{N} \left[ [\beta^T Z_{(j)}]^T \hat{\gamma} \{ \hat{\beta}^T Z_{(j)} \} - \log \left\{ \sum_{i \in \mathcal{R}_j} \exp \left( [\beta^T Z_i]^T \hat{\gamma} \{ \hat{\beta}^T Z_i \} \right) \right\} \right].$$

This approximation may result in some efficiency loss, but has computational advantages over the estimate in Step 2, since we do not need to estimate  $\psi(v)$  and thus can skip Step 1(d). The resulting estimate for  $\beta$  was shown in [19] to be  $\sqrt{n}$ -consistent, consequently an ideal choice as the initial estimate for  $\beta$ .

**Remark 3.** In step 1, the local log partial likelihood in (3) is replaced by a pseudo log partial likelihood with  $\beta$  replaced by  $\hat{\beta}$ . As this  $\hat{\beta}$  approaches  $\beta$ , the link estimate resulting from maximizing the pseudo log partial likelihood can be expected to approach the true link function at the usual nonparametric rate. This is because the parametric estimate  $\hat{\beta}$  converge to its target at the  $\sqrt{n}$ -rate, which is faster than the nonparametric link estimate. A rigorous proof is provided in Theorem 1 and Theorem 2.

imsart-coll ver. 2008/08/29 file: Wang.tex date: March 25, 2009

**Remark 4**. For large sample sizes, it is unnecessary to estimate the link function at each data point. An alternative way is to estimate the link function at equaldistance grid points, then use interpolation or smoothing methods to obtain the estimated value at each data point. Our simulation results show that this short-cut is computationally economical while retaining similar accuracy.

#### 2.2. Main Results

Let  $f(\cdot)$  be the probability density of  $\beta^T Z$ , for a given v, let  $P(t \mid v) = P(X \ge t \mid \beta^T Z = v)$ ,  $Y(t) = I\{X \ge t\}$ ,  $H = diag\{h, \cdots, h^p\}^T$  and  $\mathbf{u} = \{u, \cdots, u^p\}^T$ .

**Theorem 1.** Under conditions (C1)-(C5) in the Appendix, for any  $\sqrt{n}$  consistent estimator  $\hat{\beta}$  of the true parameter  $\beta_0$ , let  $\hat{\gamma}(\cdot)$  be the corresponding estimator for the derivatives  $\gamma_0(\cdot)$  of the true link  $\psi$  and  $\hat{\psi}(v) = \int_0^v \hat{\psi}'(w) dw$ , where  $\hat{\psi}'(\cdot)$  is the first component of  $\hat{\gamma}(\cdot)$ . If  $h \to 0$ ,  $nh/\log n \to \infty$ ,  $nh^4 \to \infty$  then

$$\sup_{v} |\hat{\gamma}(v) - \gamma_0(v)| \to_p 0,$$

and

$$\sup |\hat{\psi}(\hat{\beta}^T z) - \psi(\beta_0^T z)| \to_p 0$$

**Theorem 2.** Under the conditions in Theorem 1 and for bounded  $nh^{2p+3}$ ,

$$(a)\sqrt{nh}\left\{H(\hat{\gamma}(v)-\gamma_0(v))-\frac{\psi^{(p+1)}(v)}{(p+1)!}A^{-1}bh^{p+1}\right\}$$

$$\to_D N \left\{ 0, \frac{\sigma^2(v)}{f(v)} A^{-1} D A^{-1} \right\}.$$

Furthermore, we have

$$(b)\sqrt{nh}\left\{H(\hat{\gamma}(\hat{\beta}^{T}z) - \gamma_{0}(\beta_{0}^{T}z)) - \frac{\psi^{(p+1)}(\beta_{0}^{T}z)}{(p+1)!}A^{-1}bh^{p+1}\right\}$$

$$\to_D N \bigg\{ 0, \frac{\sigma^2(\beta_0^T z)}{f(\beta_0^T z)} A^{-1} D A^{-1} \bigg\},$$

where  $A = \int \mathbf{u} \mathbf{u}^T K(u) du - \nu_1 \nu_1^T$ ,  $b = \int u^{p+1} (\mathbf{u} - \nu_1) K(u) du$ ,  $D = \int K^2(u) (\mathbf{u} - \nu_1)^{\otimes 2} du$ ,  $\nu_1 = \int \mathbf{u} K(u) du$ , and  $\sigma^2(v) = E\{\delta | \beta^T Z = v\}^{-1}$ .

Theorem 1 establishes the uniform consistency of the local partial likelihood estimator of  $\gamma_0$  and Theorem 2 provides the joint asymptotic normality of the derivative estimators. The limiting distribution of  $\hat{\gamma}$  is identical to the one in [8], where  $\beta$  is assumed to be known. Thus, there is no efficiency loss as long as  $\beta$  can be estimated at the usual  $\sqrt{n}$ -rate.

#### 

# 2.3. Model Checking and Selection

While an estimated link function is of interest to correctly reflect the risk associated
with a covariate, a parametric link function is often preferable to a nonparametric
one to lend a parsimonious model with more interpretable results. Thus, a main
incentive to estimate the link function could be for exploratory model selection to
facilitate the choice of a proper parametric link function in the proportional hazards

model (1). If so, the  $\beta$  estimate in Step 2 only aids in the link estimation and need not be the end product. Once a suitable link function has been selected, Theorem 2 can be used for model checking. For instance, to check the identity link function under the Cox model, one can test  $H_0: \psi'(v) = 1$ . Since the first component of  $\gamma(v)$ is  $\psi'(v)$ , a local polynomial of order p = 2 is usually employed to estimate such a derivative, and the resulting asymptotic distribution of the corresponding estimate is given below.

**Corollary 1.** Under the condition of Theorem 2, and with p = 2 there, we have

$$\sqrt{nh^3} \bigg( \hat{\psi}'(v) - \psi'(v) - \frac{1}{6} \psi^{(3)}(v) h^2 \frac{\int u^4 K(u) du}{\int u^2 K(u) du} \bigg) \to_D N \bigg( 0, \frac{\sigma^2(v) \int u^2 K^2(u) du}{f(v) \int u^2 K(u) du} \bigg).$$

Corollary 1 facilitates the construction of testing procedures and asymptotic simultaneous confidence bands for the link function, but rigorous asymptotic theory requires much further work and is not available yet. In principle, one could check the appropriateness of the link function at all data points v that falls in the range of  $\beta^T Z$ . Since the true value of  $\beta$  is unknown, it is natural to replace it with an estimate. However, one must bear in mind the precision of this estimate as well as the low precision of  $\psi'(v)$  for v in the boundary region of  $\hat{\beta}^T Z$ . Here boundary region is defined as within one bandwidth of the data range, where a smoothing procedure is employed. Since the bandwidth h is usually of a higher order than  $n^{-\frac{1}{2}}$ , the anticipated rate of convergence for  $\hat{\beta}$ , we recommend to restrict inference on  $\psi'(v)$  for v that is in the interior and at least one bandwidth h away from either boundary of the range of  $\hat{\beta}$ .

Short of such a rigorous inference procedure for model checking, pointwise confidence intervals have often been used as a substitute for exploratory purposes. In the example in Section 4, we illustrate how to check the appropriateness of the Cox model, i.e. identity link function, using pointwise confidence intervals developed from Corollary 1. Readers should bear in mind that this is only an exploratory data analysis tool rather than a formal inference procedure.

#### 3. Simulation Studies

To see how the algorithm in Section 2 works for the proposed model, we conducted simulation studies where a quadratic link function  $\psi(\beta^T Z) = (\beta^T Z)^2$  with  $\beta =$  $(1,3)^T$  and a constant baseline,  $\lambda_0 = 0.005$ , were employed. The design for the two-dimensional covariate,  $Z = (Z_1, Z_2)^T$  is:  $Z_1 \sim U(-1, 1)$  and  $Z_2$  is a truncated N(0,1) with values in [-1,1]. Parameters of  $\beta$  were chosen in such a way that the simulation generates a reasonable signal to noise ratio (cf. Fig 1). If we take  $\varepsilon$  to have the standard exponential distribution,  $\exp(1)$ , the resulting hazard function will be  $\lambda_0 \exp\{\psi(\beta^T Z)\}$ , and survival times from this model can be generated as  $T = \exp\{\psi(-\beta^T Z)\}\varepsilon/\lambda_0$ . Different uniform distributions were utilized to generate three independent censoring times so that the censoring rates were 0%, and roughly 25% and 50%. The Epanechnikov kernel was adopted in the link estimation. Two sample sizes, 200 and 50, were selected to see whether the methods are flexible for moderate to small samples. For n = 200 we used 25 equal-distance grid points to estimate the link function to save computational time as elucidated in Remark 4 of Section 2. Piece-wise spline interpolation was then used to get the link estimate during each iteration of the algorithm. 

<sup>50</sup> Due to the complication from the identifiability problem, the link function can <sup>51</sup> only be identified up to a constant. Thus,  $\hat{\psi}(v)$  and  $\hat{\psi}(v) + c$ , for any constant

# Proportional Hazards Regression TADLE 1

|  |          |                       |             | $h^*$                            |                    | Optim |
|--|----------|-----------------------|-------------|----------------------------------|--------------------|-------|
| Censoring                                    | †        | 0.1                   | 0.2         | 0.3                              | 0.4                | MSE   |
| No   | 1        | 3.152/0.181           | 3.152/0.181 | 3.152/0.181                      | 3.152/0.181        | 9.970 |
|  | 4        | 0.191/0.114           | 0.191/0.114 | 0.191/0.114                      | 0.191/0.114        | 0.049 |
|  | 2        | 1.234/2.554           | 0.439/0.142 | 0.843/0.150                      | 1.417/0.151        | 0.213 |
|  | 3        | 1.210/3.522           | 0.490/0.144 | 0.907/0.168                      | 1.476/0.168        | 0.261 |
| 25%  | 1        | 3.151/0.181           | 3.151/0.181 | 3.151/0.181                      | 3.151/0.181        | 9.961 |
|  | 4        | 0.201/0.112           | 0.201/0.112 | 0.201/0.112                      | 0.201/0.112        | 0.053 |
|  | 2        | 1.256/2.548           | 0.425/0.157 | 0.684/0.172                      | 1.197/0.178        | 0.206 |
|  | 3        | 0.981/1.991           | 0.468/0.156 | 0.746/0.185                      | 1.256/0.194        | 0.244 |
| 50%  | 1        | 3.149/0.181           | 3.149/0.181 | 3.149/0.181                      | 3.149/0.181        | 9.947 |
|  | 4        | 0.210/0.117           | 0.210/0.117 | 0.210/0.117                      | 0.210/0.117        | 0.056 |
|  | 2        | 1.361/2.529           | 0.525/0.215 | 0.535/0.161                      | 0.729/0.210        | 0.322 |
|  | 3        | 1.236/2.138           | 0.536/0.201 | 0.571/0.169                      | 0.808/0.226        | 0.327 |
| <sup>†</sup> Method 1 is under identity link |          |                       |             | Method 2 is und                  | ler unknown link   |       |
| and unknown $\beta$ .                        |          |                       |             |                                  | and true $\beta$ . |       |
| Method 3                                     | is un    | der unknown link      |             | Method 4 is under quadratic link |                    |       |
| method c                                     | ) 15 UII | and unknown $\beta$ . |             | Method 4 is und                  | and unknown (      | 3.    |

c, are considered to be equivalent procedures, and any measures of performance would declare these two procedures identical. This points to selecting a measure which measures the variation instead of the real difference. We adopt a measure proposed in [10] which is the standard deviation of the differences between the fitted values  $\hat{\psi}(\hat{\beta}^T Z)$  and the true values  $\psi(\beta^T Z)$  at all data points. More specifically, this measure, denoted by d, is the standard deviation of the difference  $\{\hat{\psi}(\hat{\beta}^T Z) - \hat{\psi}(\hat{\beta}^T Z)\}$  $\psi(\beta^T Z): i = 1, \cdots, n$ . We report in Table 1 the average values for this measure and its standard deviation based on 100 simulation runs.

Since at each estimating step,  $\hat{\beta}$  was updated and the range of  $\hat{\beta}^T Z$  might be different, we used a bandwidth  $h^*$ , which took a certain portion of the range of  $\hat{\beta}^T Z$ . For instance, an  $h^* = 0.3$  means that the actual bandwidth is 0.3 times the range of the values of  $\hat{\beta}^T Z$ . Various bandwidths were explored, but we report only the results for bandwidth  $h^*$  varying from 0.1 to 0.4 (with 0.1 increment) times the data range of  $\beta^T Z$  at each iteration stage. Results for other bandwidths were inferior and are not reported here.

Four procedures were compared and the results for n = 200 are shown in Table 1. Method 1 assumes that the link is identity (which is incorrect here) and the regression coefficient estimate  $\hat{\beta}$  is therefore the Cox estimate based on the par-tial likelihood estimate. The aim is to see the effect of erroneously assuming the conventional Cox proportional hazards model. Method 2 assumes that  $\beta$  is known and estimates the unknown link function as in [8]. Method 3 is the new procedure where both the link function and regression coefficient  $\beta$  are estimated. Method 4 assumes that the true quadratic link function is known and the regression coefficient estimate,  $\hat{\beta}$ , is the partial likelihood estimate. The comparisons for the distance d are reported in Table 1. The results of the best procedures together with the corre-sponding optimal bandwidths are highlighted with boxes. It is not surprising that the best results came from the procedures with true quadratic link function and unknown  $\beta$ . Our estimators are close to those from method 2 ([8]) with known  $\beta$ , while the estimators based on the identity link model have much larger d. To demonstrate the effect of an erroneous link function on regression estimates,

W. Wang, J. Wang and Q. Wang

| $(\beta^T Z) =$ | $(\beta^T)$ | $Z)^2$ , where $\beta =$ | = $(1,3)^T$ and Z                          | $= (Z_1, Z_2)^T $             | with $Z_1 \sim U(-$ | (1,1) and Z       | $_{2} \sim N(0,1)$ |
|-----------------|-------------|--------------------------|--|-------------------------------|---------------------|-------------------|--------------------|
| (truncate       | ea at       | [-1,1]), n=200<br>devid  | ). The numbers itions of the $\hat{eta}_2$ | before and an<br>based on 100 | simulations.        | e oiases ana      | stanaara           |
|                 |             |                          | h*   | 5                             |                     | Optimal           | Optimal            |
| Censoring       | t           | 0.1                      | 0.2  | 0.3                           | 0.4                 | mse <sup>‡</sup>  | MSE§               |
| No              | 1           | 5.344/50.847             | 5.344/50.847                               | 5.344/50.847                  | 5.344/50.847        | 2613.949          | 11742.108          |
|                 | 4           | -0.015/0.142             | -0.015/0.142                               | -0.015/0.142                  | -0.015/0.142        | 0.021             | 0.703              |
|                 | 3           | 0.023/0.276              | -0.051/0.154                               | -0.096/0.158                  | -0.215/0.190        | 0.026             | 0.968              |
| 25%             | 1           | -1.229/24.742            | -1.229/24.742                              | -1.229/24.742                 | -1.229/24.742       | 613.684           | 11320.063          |
|                 | 4           | -0.013/0.154             | -0.013/0.154                               | -0.013/0.154                  | -0.013/0.154        | 0.024             | 0.801              |
|                 | 3           | 0.038/0.307              | -0.055/0.170                               | -0.089/0.171                  | -0.162/0.195        | 0.032             | 1.144              |
| 50%             | 1           | -1.472/6.486             | -1.472/6.486                               | -1.472/6.486                  | -1.472/6.486        | 44.237            | 11266.339          |
|                 | 4           | -0.006/0.170             | -0.006/0.170                               | -0.006/0.170                  | -0.006/0.170        | 0.029             | 0.968              |
|                 | 3           | 0.055/0.371              | -0.050/0.188                               | -0.075/0.181                  | -0.110/0.190        | 0.038             | 1.362              |
| † Method 1      | is un       | der identity link a      | and unknown $\beta$ .                      | Method 4 is un                | der true link and   | unknown $\beta$ . |                    |

З

we report in Table 2 the results of the various estimates for  $\beta$ . Since there is no regression parameter estimate for method 2, only three procedures are compared in Table 2. There are several ways to compare the regression estimates, one way is to set the first component of  $\beta$  to the true value, then compare the difference between  $\beta$  and  $\beta$  based on the second component. Table 2 shows the results of the difference between the true  $\beta_2$  and the estimate  $\beta_2$  for various procedures. The best procedures for the profile estimator in method 3 are shown in boxes under the optimal bandwidths. Another way to compare the different estimators is to calculate the angles between these estimators and the true parameter. To save space only the MSEs based on the optimal bandwidths are listed in the last column of Table 2 for the angle measure (degrees). Based on both optimal MSE measures reported in the last two columns of Table 2, the differences between the new profile estimators and the true parameters are way smaller than those from the identity link model, and reasonably close to those under the true link.

We can see that using the wrong link will lead to huge bias and MSE under all censoring patterns. The average angles between the  $\beta$  estimates assuming identity link and the true parameters are around 90°, which suggests that the  $\beta$  estimates with identity link are perpendicular to the true parameter space, indicating a total inability to estimate the regression parameter. This is in addition to the problem that the link function itself has been misspecified. Both underscore the importance to check the shape of the link function at the beginning of data analysis. 

Four typical simulated curves are shown in Fig 1. The procedure (method 4) with known quadratic link function and unknown  $\beta$  performed the best. The procedure (method 2) with known  $\beta$  and our procedure captured the shape of the true curve well, but the procedure (method 1) based on the Cox model failed to capture the shape of the link function. Results for n = 50 summarized in Table 3 are consistent with the findings for n = 200 in Table 1.

#### 4. Data Analysis

In this section, we illustrate the proposed model and estimation algorithm in Section 2 through the Worcester Heart Attack Study (WHAS) data. One of the goals of this study is to identify factors associated with the survival rates following hospital admission for acute myocardial infarction. The main data set has more than 11000 





admissions, but we used only a random sample of 500 patients as listed in [14]. This data set is chosen because the proportionality assumption has been carefully examined and is reasonably satisfied. Our goal here is to check the adequacy of the identity link function in the Cox proportional hazards regression model.

There were more than 10 covariates in the data set. After detailed model selection procedure, Hosmer et al. [14] included 6 variables age (AGE), initial heart rate (HR), initial diastolic blood pressure (DIASBP), body mass index (BMI), gender (GENDER), congestive heart complications (CHF), and the interaction between age and gender (AGEGENDER) in their model. After examining the linearity assumption using fractional polynomials, they decided to apply a two-term fractional polynomial model to the variable BMI. We thus begin with the univariate covariate BMI.

We tried different bandwidths and found similar patterns of the estimated link functions. In Fig 2 we report two of the results, which exhibit reasonable level of smoothness. The estimated link function in Fig 2 suggests clear nonlinearity. We then constructed a 95% point-wise approximated confidence interval of  $\gamma(v)$  (Fig 3) to see whether it would cover the constant function 1. The results suggest that the estimated link functions have some curvature and further investigation is needed.



FIG 2. The estimated link function, under two bandwidths, for the WHAS data with BMI as a

Next we applied the proposed procedure to the multivariate model with all 7 covariates. We tried different bandwidths ranging from 1/10, 1/8, 1/7, 1/6, 1/5, 1/41/3, to 1/2 of the single index range, and plotted the results of three bandwidths in Fig 4. The estimated link function for  $h^* = 1/4$  appears oversmoothed but all three estimates exhibit two bumps. The 95% point-wise approximated confidence intervals for  $\gamma(v)$  as shown in Fig 5 also reveal curvature away from the constant (=1) horizontal line. Although it is arguable that the Cox model could be rejected at a low level, the suitability of an identity link function seems questionable.

covariate.



sion model with unknown link function and multi-dimensional covariates seem to be
 reliable for moderate to large sample sizes. Once the link function and the param eters of the index have been established, one can proceed to estimate the unknown

imsart-coll ver. 2008/08/29 file: Wang.tex date: March 25, 2009

49

50



FIG 5. The estimated confidence interval of  $\gamma$  for the WHAS data at bandwidth  $h^* = 1/6$ .

baseline hazard function in model (1) using a Breslow-type estimate ([1]).

The cost of a misspecified link function has been demonstrated through the simulation studies in Section 3. As a consequence, the risk of an individual may be misinterpreted. It is thus important to at least estimate the link function in the initial model fitting stage as a model checking tool or guidance to a suitable class of parametric link functions. A rigorous test of parametric link function will be a worthwhile future project, as is the asymptotic theory for simultaneous inference of the link function and regression parameters.

The choice of automatic smoothing parameters, the bandwidth h in this case, is a challenging problem for proportional hazards model when a likelihood based smoother, such as the local partial likelihood estimate, is employed in the link es-timate. This is because the components of the partial likelihood are dependent, hence the usual automatic bandwidth selection methods for linear models are not applicable here. The usual least square cross-validation procedure in nonparametric regressions also cannot be easily adapted to hazard based models such as the pro-portional hazards model. An alternative criterion, less computational intensive than cross-validation methods, was proposed in [18], based on a variation of the Akaike's information criterion for span selection, when the nearest neighborhood method was used for smoothing. However, the interpretation of AIC is not clear here since partial likelihood involves dependent components. The authors also acknowledged that the asymptotic correctness of the AIC criterion has not been established. Thus, automatic bandwidth choice remains an open question when the link function is be-ing estimated. Meanwhile, we recommend to try several bandwidths and choose one that yields a moderately smooth link function as we did for the WHAS data. This subjective choice based on the visual degree of smoothness is commonly adopted as an ad hoc tool. 

50 While this paper deals with time-independent covariates, it would be desirable to 51 extend model (1) to time-dependent covariates as well. One complication is that the

З

entire history of the covariate process would be required or some kind of imputation needs to be performed to get even an initial estimate of  $\beta$ . Preliminary results were reported in [20] by imputing the covariate process through a functional principal components approach, and then proceeding with the estimation of the survival components at the second stage. Such a two-stage procedure is prone to bias as is well known in the joint modelling literature. Further investigations to correct the bias would be desirable, and joint modelling the longitudinal and survival process offers some hope if one can resolve the additional complication of an unknown link function. This is yet another worthwhile project to pursue in the future.

Acknowledgement

The work of Jane-Ling Wang and Wei Wang was completed with the partial sup-port of National Science Foundation grants DMS04-06430, National Institutes of Health grant R01-DK-45939, and National Institute of Allergy and Infectious Dis-eases grant AI24643. Qihua Wang's research was supported by the National Science Fund for Distinguished Young Scholars in China (10725106), the National Natural Science Foundation of China (10671198), the National Science Fund for Creative Research Groups in China and the Research Grants Council of the Hong Kong (HKU 7050/06P). The authors are grateful for the invaluable suggestions of two reviewers and the editor.

| 1  | Appendix 1   |
|----|--|
| 2  | 2  |
| 3  | Let $f(\cdot)$ be the probability density of $\beta^T Z$ , for a given $v$ , let $P(t \mid v) = P(X \geq 3)$   |
| 4  | $t \mid \beta^T Z = v), \ \Lambda(t, v) = \int_0^t P(u \mid v) \lambda_0(u) du, \ Y(t) = I\{X \ge t\}, \ Y_i(t) = I\{X_i \ge t\}, $  |
| 5  | $H = diag\{h, \cdots, h^p\}^T \text{ and } \mathbf{u} = \{u, \cdots, u^p\}.$   |
| 6  | We begin with some regularity conditions needed for the results. <sup>6</sup>  |
| 7  | (C1) $K \ge 0$ is a bounded density with compact support, and it has bounded <sup>7</sup>  |
| 8  | first and second derivative.   |
| 9  | $(C2) \psi(\cdot)$ has a continuous $(p+1)$ th derivative around v.  |
| 10 | (C3) The density $f(\cdot)$ of $\beta^T Z$ is continuous at point v and $\inf_w f(w) > 0$ .  |
| 11 | $(C4)$ The conditional probability $P(t \mid \cdot)$ is equicontinuous at $v$ .  |
| 12 | $(C5)  \int_0^\tau \lambda_0(u) du < \infty. $   |
| 13 | Denote 13  |
| 14 | $s_0(\beta, \psi, u) = E[Y(u) \exp\{\psi(\beta^T Z)\}].$ <sup>14</sup>   |
| 15 | $s_1(\beta, \psi, u) = E[Y(u) \exp\{\psi(\beta^T Z)\}],$ $s_1(\beta, \psi, u) = E[Y(u) \exp\{\psi(\beta^T Z)\}],$ <sup>15</sup>  |
| 16 | $S_1(\beta, \psi, u) = E[Y(u) \exp[\psi(\beta TZ)] \psi(\beta TZ)],$ $S_1(\beta, \psi, u) = E[Y(u) \exp[\psi(\beta TZ)] \int \psi(\beta TZ) + [u/(\beta TZ)]^2 ZZT]$ 16  |
| 17 | $s_{2}(\beta, \psi, a) = E[I(a) \exp\{\psi(\beta Z)\} \{\psi(\beta Z) + [\psi(\beta Z)] \} ZZ ],$ $s_{17}^{*}(\beta, a) = E[V(a) \exp\{\phi(\beta Z)\} [\phi(\beta Z)] \{ZZT]$   |
| 18 | $s_2(\beta, \psi, u) = E[I(u)\exp\{\psi(\beta^2 Z)\}[\psi(\beta^2 Z)]^2 Z Z^2].$ ( <i>CC</i> ) The functions are $0, 1$ and $2$ and $x^*$ are bounded and $x$ is bounded.  |
| 19 | $(C0)$ The functions $s_r, r = 0, 1, \text{and } 2, \text{ and } s_2$ are bounded and $s_0$ is bounded<br>super from 0 on $\mathcal{R} \times [0, \pi]$ , the family of functions of $(s_1, s_2)$ and $s_1^*(s_2, s_3)$ is on  |
| 20 | away from 0 on $B \times [0,7]$ ; the family of functions $s_r(\cdot, \psi, u)$ and $s_2(\cdot, \psi, u)$ is an  |
| 21 | Equicontinuous family at $p_0$ .<br>The following low is used repeatedly in the later proof. The proof will be   |
| 22 | ine following lemma is used repeatedly in the later proofs. The proof will be  |
| 23 | confitted and can be found in [19].  |
| 24 | <b>Lemma 1.</b> Let $c_n(\beta_0, t) = n^{-1} \sum_{i=1}^n Y_i(t) g(\beta_0^T Z_i) K_h(\beta_0^T Z_i - v)$ and $c(t) = 24$   |
| 25 | $f(v)g(v)P(t \mid v) \int K(u)du$ , Under conditions (C1) and (C4), if $g(\cdot)$ is continuous 25   |
| 26 | at the point v, then 26  |
| 27 | $\sup  c_n(\beta_0, t) - c(t)  \to_P 0, $ <sup>27</sup>  |
| 28 | $0 \le t \le \tau$ 28  |
| 29 | provided that $h \to 0, nh/\log n \to \infty, 0 < \tau \le +\infty.$ <sup>29</sup>   |
| 30 | If furthermore, $\hat{\beta}$ is a $\sqrt{n}$ -consistent estimate of $\beta_0$ and $nh^4 \to \infty$ , then   |
| 31 | 31   |
| 32 | $\sup  c_n(\hat{\beta}, t) - c(t)  \to_P 0.$   |
| 33 | $0 \le t \le \tau$ 33  |
| 34 | Proof of Theorem 1   |
| 35 | 35   |
| 36 | <i>Proof.</i> For notation simplicity, we will use $\gamma$ to represent $\gamma(v)$ and $\gamma_0$ for $\gamma_0(v)$ .  |
| 37 | The log local partial likelihood function at point $v$ is given as <sup>37</sup>   |
| 38 | 38   |
| 39 | $l[\beta,\gamma,\eta] = \frac{1}{N} \sum_{K} K \left(\beta^T Z_{\gamma}, \dots, \eta\right) $ <sup>39</sup>  |
| 40 | $l\{\rho,\gamma,v\} = -\sum_{n \neq j} \kappa_h(\rho \ Z_{(j)} - v) $  |
| 41 | j=1 41   |
| 42 | $\left[ \left[ \left( \beta^T \mathbf{Z}_{(i)} \right)^T \gamma \right] - \log \left\{ \sum \exp\{ \left( \beta^T \mathbf{Z}_i \right)^T \gamma \right\} K_h \left( \beta^T Z_i - v \right) \right\} \right]. $ <sup>42</sup>  |
| 43 | $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 2 \\ i \in \mathcal{R}_i \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ i \in \mathcal{R}_i \end{bmatrix} $   |
| 44 | 44   |
| 45 | Using counting process notation $N(t) = I\{X \le t, \delta = 1\}$ and $N_i(t) = I\{X_i \le 45\}$   |
| 46 | $t, \delta_i = 1$ , under the independent censoring, <sup>46</sup>   |
| 47 | 47   |
| 48 | $M_{i}(t) = N_{i}(t) - \int^{t} Y_{i}(y) \exp\{y b(\beta_{0}^{T} Z_{i})\} \lambda_{0}(y) dy $ <sup>48</sup>  |
| 49 | $J_0 = J_0 $ |
| 50 | 50   |
| 51 | is a martingale with respect to the filtration $\mathcal{F}_t = \sigma\{N(u), I_{\{X \le u, \delta = 0\}} : 0 \le u \le t\}.$ 51   |

$$l_n(\beta, \gamma, t, v) = \int_0^t \frac{1}{n} \sum_{i=1}^n K_h(\beta^T Z_i - v) \left[ \left[ (\beta^T \mathbf{Z}_i)^T \gamma \right] \right]$$

The empirical counterpart of  $l\{\beta,\gamma,v\}$  up to time t is

$$-\log\left\{\sum_{i}^{n} Y_{i}(u) \exp\left\{\left(\beta^{T} \mathbf{Z}_{i}\right)^{T} \gamma\right\} K_{h}\left(\beta^{T} Z_{i}-v\right)\right\} dN_{i}(u).$$

Denote 
$$S_{h,0}(\beta,\gamma,u,v) = \frac{1}{n} \sum_{i=1}^{n} K_h(\beta^T Z_i - v) Y_i(u) \exp\{(\beta^T \mathbf{Z}_i)^T \gamma\}.$$

Let  $\hat{\beta}$  be a  $\sqrt{n}$ -consistent estimate of the true parameter  $\beta_0$ , and  $\hat{\gamma}$  be the corresponding estimate of the true  $\gamma_0$ , we can write

$$l_n(\hat{\beta},\gamma,\tau,v) - l_n(\beta_0,\gamma_0,\tau,v)$$

$$= \int_{0}^{\tau} \frac{1}{n} \sum_{i=1}^{n} K_{h}(\beta_{0}^{T} Z_{i} - v)$$
<sup>15</sup>
<sup>16</sup>
<sup>17</sup>

$$\times \left[ \left[ (\beta_0^T \mathbf{Z}_i)^T \gamma - (\beta_0^T \mathbf{Z}_i)^T \gamma_0 \right] - \log \frac{S_{h,0}(\beta_0, \gamma, u, v)}{S_{h,0}(\beta_0, \gamma_0, u, v)} \right] dM_i(u)$$

$$+\int_0^\tau \frac{1}{n} \sum_{i=1}^n K_h(\hat{\beta}^T Z_i - v) \bigg[ \big[ (\hat{\beta}^T \mathbf{Z}_i)^T \gamma - (\beta_0^T \mathbf{Z}_i)^T \gamma \big] - \log \frac{S_{h,0}(\hat{\beta}, \gamma, u, v)}{S_{h,0}(\beta_0, \gamma, u, v)} \bigg] dM_i(u)$$

$$+\int_{0}^{\tau} \frac{1}{n} \sum_{i=1}^{n} \left( K_{h}(\hat{\beta}^{T} Z_{i} - v) - K_{h}(\beta_{0}^{T} Z_{i} - v) \right)$$
<sup>23</sup>
<sup>24</sup>
<sup>25</sup>

$$\times \left[ (\beta_0^T \mathbf{Z}_i)^T \gamma - \log S_{h,0}(\beta_0, \gamma, u, v) \right] dM_i(u)$$

$$\overset{26}{27}$$

$$\overset{26}{27}$$

$$+ \int_0^\tau \frac{1}{n} \sum_{i=1}^n K_h(\beta_0^T Z_i - v) \left[ \left[ (\beta_0^T \mathbf{Z}_i)^T \gamma - (\beta_0^T \mathbf{Z}_i)^T \gamma_0 \right] - \log \frac{S_{h,0}(\beta_0, \gamma, u, v)}{S_{h,0}(\beta_0, \gamma_0, u, v)} \right]$$

$$\times Y_{i}(u) \exp\{\psi(\beta_{0}^{T} Z_{i})\}\lambda_{0}(u)du$$

$$+\int_0^\tau \frac{1}{n} \sum_{i=1}^n K_h(\hat{\beta}^T Z_i - v) \left[ \left[ (\hat{\beta}^T \mathbf{Z}_i)^T \gamma - (\beta_0^T \mathbf{Z}_i)^T \gamma \right] - \log \frac{S_{h,0}(\beta, \gamma, u, v)}{S_{h,0}(\beta_0, \gamma, u, v)} \right]$$

$$\times Y_i(u) \exp\{\psi(\beta_0^T Z_i)\}\lambda_0(u)du$$

$$+ \int_{0}^{\tau} \frac{1}{n} \sum_{i=1}^{n} \left( K_{h}(\hat{\beta}^{T} Z_{i} - v) - K_{h}(\beta_{0}^{T} Z_{i} - v) \right) \left[ (\beta_{0}^{T} \mathbf{Z}_{i})^{T} \gamma - \log S_{h,0}(\beta_{0}, \gamma, u, v) \right]$$

$$\equiv X_n(\beta_0, \gamma, \tau, v) + I + II + A_n(\beta_0, \gamma, \tau, v) + III + IV.$$

Under the regularity conditions and from Lemma 1, it can be shown that

(1)  $X_n(\beta_0, \gamma, \tau, v)$  is a locally square integrable martingale with the predictable variation process

$$\langle X_n(\beta_0,\gamma,\tau,v), X_n(\beta_0,\gamma,\tau,v) \rangle$$

$$= \int_0^\tau \frac{1}{n^2} \sum_{i=1}^n K_h^2 (\beta_0^T Z_i - v) \left[ (\beta_0^T \mathbf{Z}_i)^T (\gamma - \gamma_0) - \log \frac{S_{h,0}(\beta_0, \gamma, u, v)}{S_{h,0}(\beta_0, \gamma_0, u, v)} \right]^2$$

$$\times Y_i(u) \exp\{\psi(\beta_0^T Z_i)\}\lambda_0(u) du$$

$$=O_p(\frac{1}{nh}).$$

imsart-coll ver. 2008/08/29 file: Wang.tex date: March 25, 2009

З

W. Wang, J. Wang and Q. Wang

$$\begin{split} f(v) \exp\{\psi(v)\}\Lambda(\tau, v) \\ &\times \left[ (\int \mathbf{u}K(u)du)^T H(\gamma - \gamma_0) - \log\{\int \exp\{\mathbf{u}^T H(\gamma - \gamma_0)\}K(u)du\} \right] + o_p(1) \\ &\equiv A(\beta_0, \gamma, \tau, v) + o_p(1), \\ (3) \quad I = O_p(\frac{1}{nh}), \ II = O_p(\frac{1}{nh^2}), \ III = O_p(\frac{1}{\sqrt{nh^2}}), \ \text{and} \ IV = O_p(\frac{1}{\sqrt{nh^4}}). \\ \text{This means} \ X_n(\beta_0, \gamma, \tau, v), \ \text{I, II, III and IV converge to zero at a faster rate than} \\ _n(\beta_0, \gamma, \tau, v). \ \text{By Lemma 8.2.1(2) in [9]}, \ l_n(\hat{\beta}, \gamma, \tau, v) - l_n(\beta_0, \gamma_0, \tau, v) \ \text{has the same} \\ \text{niting distribution as} \ A_n(\beta_0, \gamma, \tau, v). \ \text{Thus, we have} \\ ) \qquad \qquad l_n(\hat{\beta}, \gamma, \tau, v) - l_n(\beta_0, \gamma_0, \tau, v) \rightarrow_p A(\beta_0, \gamma, \tau, v). \\ \text{It is obvious that} \ A(\beta_0, \gamma, \tau, v) \ \text{is strictly concave, with a maximum at } \gamma = \gamma_0 \\ \text{ence, the right-hand side of (4) is maximized at } \gamma = \gamma_0. \ \text{The left-hand side of} \\ \text{) is maximized at } \gamma = \hat{\gamma}, \ \text{since } \hat{\gamma} \ \text{maximizes} \ l_n(\hat{\beta}, \gamma, \tau, v). \ \text{Therefore, } \ \sup_p |\hat{\gamma}(v) - v| = 0 \\ \hline = 0 \ \sum_{n=1}^{n} (\beta_n(\gamma, \tau, v)) = 0 \ \sum_{n=1}^{n} (\beta_n(\gamma, \tau, v)) \ \sum_{n=1}^{n} (\beta$$

 $A_n$ е lin

(4) 
$$l_n(\hat{\beta}, \gamma, \tau, v) - l_n(\beta_0, \gamma_0, \tau, v) \to_p A(\beta_0, \gamma, \tau, v).$$

He of (4) ximizes  $l_n(\beta, \gamma, \tau, v)$ .  $\gamma = \gamma$ , since  $\gamma$  $p_v | \gamma(v)$  $\gamma_0(v) | \rightarrow_p 0.$ 

By Dominated Convergence Theorem, we have

(5) 
$$\sup_{v} \left| \hat{\psi}(v) - \psi(v) \right| \to_{p} 0.$$

This implies

$$\sup_{z} \left| \hat{\psi}(\hat{\beta}^{T}z) - \psi(\beta_{0}^{T}z) \right| \leq \sup_{z} \left| \hat{\psi}(\hat{\beta}^{T}z) - \psi(\hat{\beta}^{T}z) \right| + \sup_{z} \left| \psi(\hat{\beta}^{T}z) - \psi(\beta_{0}^{T}z) \right| \to_{p} 0,$$

where the second term converges to zero by continuity of  $\psi$  and  $\sqrt{n}$ -consistency of  $\hat{\beta}$ . Theorem 1 is thus proved.

# Proof of Theorem 2

*Proof.* Let  $\eta = H\gamma$ , we can write the log local partial likelihood function in terms of  $\beta$  and  $\eta$ 

$$l_n(\beta,\eta,\tau,v) = \int_0^\tau \frac{1}{n} \sum_{i=1}^n K_h(\beta^T Z_i - v) \left[ (\beta^T \mathbf{Z}_i)^T H^{-1} \eta \right]$$

$$-\log\left\{\sum_{i}^{n}Y_{i}(u)\exp\{(\beta^{T}\mathbf{Z}_{i})^{T}H^{-1}\eta\}K_{h}(\beta^{T}Z_{i}-v)\right\}\right]dN_{i}(u).$$

Accordingly let  $S_{h,0}(\beta, \eta, u, v) = \frac{1}{n} \sum_{i=1}^{n} K_h(\beta^T Z_i - v) Y_i(u) \exp\{(\beta^T \mathbf{Z}_i)^T H^{-1} \eta\},\$ and  $S_{h,1}(\beta, \eta, u, v) = \frac{1}{n} \sum_{i=1}^{n} K_h(\beta^T Z_i - v) Y_i(u) \exp\{(\beta^T \mathbf{Z}_i)^T H^{-1} \eta\}(\beta^T \mathbf{Z}_i)^T H^{-1},\$ 

and for a 
$$\sqrt{n}$$
-consistent estimate  $\beta$  of  $\beta_0$ , Lemma 1 implies

$$\left|\frac{S_{h,1}(\hat{\beta},\eta,u,v)}{1-\nu_1}-\nu_1\right|\to 0$$

$$\sup_{0 \le u < \tau} \left| \frac{\sup}{S_{h,0}(\hat{\beta}, \eta, u, v)} - \nu_1 \right| \xrightarrow{\to}_p 0,$$

where  $\nu_1 = \int \mathbf{u} K(u) du$ .

(2)  $A_n(\beta_0, \gamma, \tau, v) \to_p$ 

З

The derivative of  $l_n(\beta, \eta, \tau, v)$  with respect to  $\eta$  evaluated at  $\hat{\beta}$  and  $\eta_0 = H\gamma_0$  is  $l'_{\pi}(\hat{\beta},\eta_0,\tau,v)$ З  $= \int_{0}^{\tau} \frac{1}{n} \sum_{k=0}^{n} K_{h}(\hat{\beta}^{T} Z_{i} - v) \bigg[ (\hat{\beta}^{T} \mathbf{Z}_{i})^{T} H^{-1} - \frac{S_{h,1}(\hat{\beta}, \eta_{0}, u, v)}{S_{h,0}(\hat{\beta}, \eta_{0}, u, v)} \bigg] dN_{i}(u)$  $= \int_0^\tau \frac{1}{n} \sum_{i=1}^n K_h(\hat{\beta}^T Z_i - v) \left[ (\hat{\beta}^T \mathbf{Z}_i)^T H^{-1} - \frac{S_{h,1}(\hat{\beta}, \eta_0, u, v)}{S_{h,0}(\hat{\beta}, \eta_0, u, v)} \right] dM_i(u)$  $+ \int_{0}^{\tau} \frac{1}{n} \sum_{i=1}^{n} K_{h}(\hat{\beta}^{T} Z_{i} - v) \left[ (\hat{\beta}^{T} \mathbf{Z}_{i})^{T} H^{-1} - \frac{S_{h,1}(\hat{\beta}, \eta_{0}, u, v)}{S_{h,0}(\hat{\beta}, \eta_{0}, u, v)} \right]$  $\times Y_i(u) \exp\{\psi(\beta_0^T Z_i)\}\lambda_0(u)du$  $\equiv U_n(\hat{\beta}, \eta_0, \tau, v) + B_n(\hat{\beta}, \eta_0, \tau, v).$ The first term  $U_n(\hat{\beta},\eta_0,\tau,v)$  $= \int_{0}^{t} \frac{1}{n} \sum_{i=1}^{n} K_{h}(\hat{\beta}^{T} Z_{i} - v) \bigg[ (\hat{\beta}^{T} \mathbf{Z}_{i})^{T} H^{-1} - \frac{S_{h,1}(\hat{\beta}, \eta_{0}, u, v)}{S_{h,0}(\hat{\beta}, \eta_{0}, u, v)} \bigg] dM_{i}(u)$  $= \int_{0}^{t} \frac{1}{n} \sum_{i=1}^{n} K_{h}(\beta_{0}^{T} Z_{i} - v) \left[ (\beta_{0}^{T} \mathbf{Z}_{i})^{T} H^{-1} - \frac{S_{h,1}(\beta_{0}, \eta_{0}, u, v)}{S_{h,0}(\beta_{0}, \eta_{0}, u, v)} \right] dM_{i}(u)$ +  $\int_{0}^{t} \frac{1}{n} \sum_{i=1}^{n} \left( K_{h}(\hat{\beta}^{T}Z_{i}-v) - K_{h}(\beta_{0}^{T}Z_{i}-v) \right)$  $\times \left[ (\beta_0^T \mathbf{Z}_i)^T H^{-1} - \frac{S_{h,1}(\beta_0, \eta_0, u, v)}{S_{h,0}(\beta_0, \eta_0, u, v)} \right] dM_i(u)$  $+\int_0^t \frac{1}{n} \sum_{i=1}^n K_h(\hat{\beta}^T Z_i - v)$  $\times \left[ (\hat{\beta}^T \mathbf{Z}_i - \beta_0^T \mathbf{Z}_i)^T H^{-1} - \frac{S_{h,1}(\hat{\beta}, \eta_0, u, v)}{S_{h,0}(\hat{\beta}, \eta_0, u, v)} - \frac{S_{h,1}(\beta_0, \eta_0, u, v)}{S_{h,0}(\beta_0, \eta_0, u, v)} \right] dM_i(u)$  $\equiv U_n(\beta_0, \eta_0, \tau, v) + V + VI.$ It is clear that  $\sqrt{nh}U_n(\beta_0, \eta_0, t, v)$  is a martingale with predictable variation  $\langle \sqrt{nh}U_n(\beta_0,\eta_0,t,v),\sqrt{nh}U_n(\beta_0,\eta_0,t,v)\rangle$  $=\frac{nh}{n^2}\sum_{i=1}^{n}\int_{0}^{\tau}K_{h}^{2}(\beta_{0}{}^{T}Z_{i}-v)\left[(\beta_{0}^{T}\mathbf{Z}_{i})^{T}H^{-1}-\frac{S_{h,1}(\beta_{0},\eta_{0},u,v)}{S_{h,0}(\beta_{0},\eta_{0},u,v)}\right]^{\otimes 2}$  $\times Y_i(u) \exp\{\psi(\beta_0^T Z_i)\}\lambda_0(u)du.$  $= f(v) \exp\{\psi(v)\} \Lambda(t, v) \int K^2(u) (\mathbf{u} - \nu_1)^{\otimes 2} du + o_p(1) \equiv \Sigma_U(t, v) + o_p(1),$ where the last step follows from Lemma 1. The Lindberg conditions are satisfied (see [19] for details), we have thus proven that (6) $\sqrt{nhU_n(\beta_0,\eta_0,\tau,v)} \rightarrow_D N(0,\Sigma_U(\tau,v)).$ 

imsart-coll ver. 2008/08/29 file: Wang.tex date: March 25, 2009

(7)and (8)(9)

As for the term V and VI, similarly to the proof of Theorem 1, we have  

$$V = \int_0^\tau \frac{1}{n} \sum_{i=1}^n \left( K_h(\hat{\beta}^T Z_i - v) - K_h(\beta_0{}^T Z_i - v) \right)$$

$$\times \left[ (\beta_0^T \mathbf{Z}_i)^T H^{-1} - \frac{S_{h,1}(\beta_0, \eta_0, u, v)}{S_{h,0}(\beta_0, \eta_0, u, v)} \right] dM_i(u)$$

$$=O_p(\frac{1}{nh^2}),$$

(b)  
$$VI = \int_0^\tau \frac{1}{n} \sum_{i=1}^n K_h(\hat{\beta}^T Z_i - v)$$
12
13
14

$$\times \left[ (\hat{\beta}^T \mathbf{Z}_i - \beta_0^T \mathbf{Z}_i)^T H^{-1} - \frac{S_{h,1}(\hat{\beta}, \eta_0, u, v)}{S_{h,0}(\hat{\beta}, \eta_0, u, v)} - \frac{S_{h,1}(\beta_0, \eta_0, u, v)}{S_{h,0}(\beta_0, \eta_0, u, v)} \right] dM_i(u)$$

$$=O_p(\frac{1}{nh}).$$

Applying Lemma 1 again and by Taylor expansion we get

$$B_n(\hat{\beta}, \eta_0, \tau, v) = f(v) \exp\{\psi(v)\} \frac{\psi^{(p+1)}(v)}{(p+1)!} \Lambda(\tau, v) \int K(u) (\mathbf{u} - \nu_1) u^{p+1} du h^{p+1} + o_p(h^{p+1}) + O_p(\frac{1}{\sqrt{n}})$$

$$=b(\tau, v) + o_p(h^{p+1}) + O_p(\frac{1}{\sqrt{n}}).$$
26  
27  
28

We have thus shown that, under (6), (7), (8) and (9),

(10) 
$$\sqrt{nh}l'_{n}(\hat{\beta},\eta_{0},\tau,v) \rightarrow_{D} N\Big(b(\tau,v),\Sigma_{U}(\tau,v)\Big).$$

Next we focus on the property of the second derivative  $l''_n(\hat{\beta}, \eta, t, v)$ . Let  $\hat{\eta} = H\hat{\gamma}$ , by Taylor expansion and Lemma 1 we have

(11) 
$$0 = l'_n(\hat{\beta}, \hat{\eta}, \tau, v) = l'_n(\hat{\beta}, \eta_0, \tau, v) + l''_n(\hat{\beta}, \eta^{**}, \tau, v)(\hat{\eta} - \eta_0),$$

where  $\eta^{**}$  lies in between  $\hat{\eta}$  and  $\eta_0$ . Theorem 1 implies  $\hat{\eta} \to_p \eta_0$ , hence  $\eta^{**} \to_p \eta_0$ . Using condition (C1) and boundedness of  $\hat{\beta}^T Z$ , we arrive at

(12) 
$$l_n''(\hat{\beta}, \eta^{**}, \tau, v) = l_n''(\hat{\beta}, \eta_0, \tau, v) + o_p(1) = \Sigma_l(\tau, v) + o_p(1).$$

By (10), (11), (12) and Slutsky's theorem,

$$\sqrt{nh}(\hat{\eta} - \eta_0) = \sqrt{nh} \bigg[ -l_n''(\hat{\beta}, \eta^{**}, \tau, v)^{-1} l_n'(\hat{\beta}, \eta_0, \tau, v) \bigg] + o_p(1)$$

$$\to_D N\Big(b(\tau,v), \Sigma_l(\tau,v)^{-1}\Sigma_U(\tau,v)\Sigma_l(\tau,v)^{-1}\Big).$$

Simple calculations lead to the result in Theorem 2.

V

З

| 1  | Re    | ferences  | 1  |
|----|-------|---|----|
| 2  |       |   | 2  |
| 3  | [1]   | BRESLOW, N. (1974). Covariance analysis of censored survival data. <i>Biometrics</i> , <b>30</b> , 89–99.   | 3  |
| 4  | [2]   | CHEN, C. H. AND LI, K. C. and WANG, J. L. (1999). Dimension reduction for censored regression data. <i>The Annals of Statistics</i> , <b>27</b> , 1–23. | 4  |
| 5  | [3]   | CHENG, S. C., WEI, L. J. and YING, Z. (1995). Analysis of transformation models with censored data.<br><i>Biometrika</i> . 82. 835–845.                 | 5  |
| 7  | [4]   | Cox, D. R. (1972). Regression models and life-table (with discussion). Journal of the Royal Statist-  | 7  |
| 8  | [5]   | Cox, D. R. (1975), Partial likelihood, <i>Biometrika</i> , <b>62</b> , 269–276.   | 8  |
| 9  | [6]   | DABROWSKA, D. M. and DOKSUM, K. A. (1988). Partial likelihood in transformation models with   | 9  |
| 10 |       | censored data. Scandinavian Journal of Statistics, 15, 1–23.  | 10 |
| 11 | [7]   | DOKSUM, K. A. (1987). An extension of partial likelihood methods for proportional hazard models   | 11 |
| 12 | [8]   | FAN, J., GIJBELS, I. and KING, M. (1997). Local likelihood and local partial likelihood in hazard   | 12 |
| 13 |       | regression. The Annals of Statistics, 25, 1661–1690.  | 13 |
| 14 | [9]   | FLEMING, T. R. and HARRINGTON, D. P. (1991). Counting Processes and Survival Analysis. John   | 14 |
| 15 | [10]  | GENTLEMAN, R. and CROWLEY, J. (1991). Local full likelihood estimation for the proportional hazards   | 15 |
| 16 |       | model. Biometrics, 47, 1283–1296.   | 16 |
| 17 | [11]  | GRAY, R. J. (1992). Flexible methods for analyzing survival data using splines, with applications   | 17 |
| 18 | [12]  | HASTIE, T. and TIBSHIRANI, R. (1986). Generalized additive models. <i>Statistical Science</i> , <b>3</b> , 297–318.                                     | 18 |
| 19 | [13]  | HASTIE, T. and TIBSHIRANI, R. (1990). Exploring the nature of covariate effects in the proportional   | 19 |
| 20 | [14]  | HAZARDS MODEL. Biometrics, 40, 1005–1010.<br>HOSMER, D. W., LEMESHOW, S. and MAY, S. (2008). Applied Survival Analysis: Regression Modeling             | 20 |
| 21 |       | of Time to Event Data: Second Edition. John Wiley & Sons, Inc., New York.   | 21 |
| 22 | [15]  | HUANG, J. and LIU, L. (2006). Polynomial spline estimation and inference of proportional hazards  | 22 |
| 23 | [16]  | O'SULLIVAN, F. (1988). Nonparametric estimation of relative risk using splines and cross-validation.  | 23 |
| 24 | L - J | SIAM Journal on Scientific and Statistical Computing, 9, 531–542.   | 24 |
| 25 | [17]  | SLEEPER, L. A. and HARRINGTON, D. P. (1990). Regression splines in the Cox model with application   | 25 |
| 26 | [18]  | TIRSHIRANL R. and HASTIE, T. (1987). Local likelihood estimation. Journal of the American Sta-  | 26 |
| 27 | L - J | tistical Association, 82, 559–567.  | 27 |
| 28 | [19]  | WANG, W (2001). Proportional hazards model with unknown link function and applications to   | 28 |
| 29 | [20]  | WANG, W (2004). Proportional hazards regression with unknown link function and time-dependent   | 29 |
| 30 |       | covariates. Statistica Sinica, 14, 885–905.   | 30 |
| 31 |       |   | 31 |
| 32 |       |   | 32 |
| 33 |       |   | 33 |
| 34 |       |   | 34 |
| 35 |       |   | 35 |
| 36 |       |   | 36 |
| 37 |       |   | 37 |
| 38 |       |   | 38 |
| 39 |       |   | 39 |
| 40 |       |   | 40 |
| 41 |       |   | 41 |
| 42 |       |   | 42 |
| 43 |       |   | 43 |
| 44 |       |   | 44 |
| 45 |       |   | 45 |
| 46 |       |   | 46 |
| 47 |       |   | 47 |
| 48 |       |   | 48 |
| 49 |       |   | 49 |
| 50 |       |   | 50 |
| 51 |       |   | 51 |

| Hiteroduction       Hiteroduction         1       Introduction         2       Sklar Type Models. Cox Regression         3       Rank, Partial and Marginal Likelihood         4       Profile NP Likelihood         5       Profile NP Likelihood         5       Profile NP Likelihood         6       Simul Algorithm for the SINAMI Model with $\theta \le 0$ 5.1       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.2       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.3       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation field Sinates         6       Simulation Results         6       Simulation Results         7       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPLIKE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         7       The MM Algorithm for the SINAMI Model with $\theta \in R$ 7  |
|--|
| University of Wisconsin, Madison Abstract: We consider classes of models related to those introduced by<br>Lehmann in 1953 and Sklar in 1959. Recently developed algorithms for finding<br>profile NP likelihood procedures are discussed, extended and implemented for<br>such models by combining them with the MM algorithm. In particular we con-<br>sider statistical procedures for a regression model with proportional expected<br>hazard rates, and for transformation models including the normal copula. A<br>variety of likelihoods introduced to deal with semiparametric models are con-<br>sidered. They all generate rank results, not only tests, but also estimates,<br>confidence regions, and optimality theory, thereby, to paraphrase Lehmann<br>(1953), demonstrating "the power of ranks".  Contents  1 Introduction 1.1 Lehmann Type Models. Cox Regression 2 Proportional Hazard and Proportional Expected Hazard Rate Models 3 Rank, Partial and Marginal Likelihood 4 Profile NP Likelihood 5 Profile NP Likelihood for the PEHR Model 5.1 The MM Algorithm 5.2 The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with<br>$\theta \ge 0$ ) 5.3 The MM Algorithm for the SINAMI Model with $\theta \in \mathbb{R}$ 5.4 The MM Algorithm for the SINAMI Model with $\theta \in \mathbb{R}$ 5.5 Profile NPLL Implementation 5.6 Estimation of the Variance of the Profile NPMLE 5.1 Dre MM Algorithm for the SINAMI Model with $\theta \in \mathbb{R}$ 5.2 Model Fit for Misspecified Model 5.3 Simulation Results 6.1 PEHR Model Estimates 6.2 Model Fit for Misspecified Model 7.1 The One Covariate Case 7.2 The Multivariate Covariate Case 7.3 The Multivariate Covariate Case 7.4 The Magnetish |
| Abstract: We consider classes of models related to those introduced by Lehmann in 1953 and Sklar in 1959. Recently developed algorithms for finding profile NP likelihood procedures are discussed, extended and implemented for such models by combining them with the MM algorithm. In particular we consider statistical procedures for a regression model with proportional expected hazard rates, and for transformation models including the normal copula. A variety of likelihoods introduced to deal with semiparametric models are considered. They all generate rank results, not only tests, but also estimates, confidence regions, and optimality theory, thereby, to paraphrase Lehmann (1953), demonstrating "the power of ranks".         Contents         1       Introduction.         1.1       Lehmann Type Models. Cox Regression         2.5       Sklar Type Models. Copula Regression         2.6       Proportional Hazard and Proportional Expected Hazard Rate Models         3.7       Rank, Partial and Marginal Likelihood         4       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm         5.2       The MM Algorithm for the SINAMI Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1  |
| Contents         1       Introduction         1.1       Lehmann Type Models. Cox Regression         1.2       Sklar Type Models. Copula Regression         2       Proportional Hazard and Proportional Expected Hazard Rate Models         3       Rank, Partial and Marginal Likelihood         4       Profile NP Likelihood         5       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 1       Introduction .         1.1       Lehmann Type Models. Cox Regression         1.2       Sklar Type Models. Copula Regression         2       Proportional Hazard and Proportional Expected Hazard Rate Models         3       Rank, Partial and Marginal Likelihood         4       Profile NP Likelihood         5       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case  |
| 1       Introduction .         1.1       Lehmann Type Models. Cox Regression .         1.2       Sklar Type Models. Copula Regression .         2       Proportional Hazard and Proportional Expected Hazard Rate Models .         3       Rank, Partial and Marginal Likelihood .         4       Profile NP Likelihood .         5       Profile NP Likelihood for the PEHR Model .         5.1       The MM Algorithm .         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ ) .         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ .         5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ .         5.5       Profile NPMLE Implementation .         5.6       Estimation of the Variance of the Profile NPMLE .         6       Simulation Results .         6.1       PEHR Model Estimates .         6.2       Model Fit for Misspecified Model .         7       Estimation in the Normal Copula Model .         7.1       The One Covariate Case .         7.2       The Multivariate Covariate Case .  |
| 1.1       Lehmann Type Models. Cox Regression         1.2       Sklar Type Models. Copula Regression         2       Proportional Hazard and Proportional Expected Hazard Rate Models         3       Rank, Partial and Marginal Likelihood         4       Profile NP Likelihood         5       Profile NP Likelihood for the PEHR Model         5       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case  |
| 1.2Sklar Type Models. Copula Regression2Proportional Hazard and Proportional Expected Hazard Rate Models3Rank, Partial and Marginal Likelihood4Profile NP Likelihood5Profile NP Likelihood for the PEHR Model515.1The MM Algorithm5.2The MM Algorithm for the PEHR Model with $\theta \ge 0$ 5.3The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5Profile NPMLE Implementation5.6Estimation of the Variance of the Profile NPMLE6Simulation Results6.1PEHR Model Estimates6.2Model Fit for Misspecified Model7.1The One Covariate Case7.2The Multivariate Covariate Case   |
| 2Proportional Hazard and Proportional Expected Hazard Rate Models3Rank, Partial and Marginal Likelihood4Profile NP Likelihood5Profile NP Likelihood for the PEHR Model5.1The MM Algorithm5.2The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with<br>$\theta \ge 0$ )5.3The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5Profile NPMLE Implementation5.6Estimation of the Variance of the Profile NPMLE6Simulation Results6.1PEHR Model Estimates6.2Model Fit for Misspecified Model7.1The One Covariate Case7.2The Multivariate Covariate Case  |
| 3       Rank, Partial and Marginal Likelihood         4       Profile NP Likelihood         5       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 4       Profile NP Likelihood         5       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.3       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         5       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 5       Profile NP Likelihood for the PEHR Model         5.1       The MM Algorithm         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         5.6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 5.1       The MM Algorithm         5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 5.2       The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )         5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| $\theta \ge 0$ )   |
| 5.3       The MM Algorithm for the SINAMI Model with $\theta \le 0$ 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 5.4       The MM Algorithm for the SINAMI Model with $\theta \in R$ 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 5.5       Profile NPMLE Implementation         5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case   |
| 5.6       Estimation of the Variance of the Profile NPMLE         6       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case  |
| b)       Simulation Results         6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         6.2       Simulation in the Normal Copula Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case         8       Transformation and NB Models  |
| 6.1       PEHR Model Estimates         6.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case         8       Transformation and NB Models  |
| 0.2       Model Fit for Misspecified Model         7       Estimation in the Normal Copula Model         7.1       The One Covariate Case         7.2       The Multivariate Covariate Case         7.3       The Multivariate Covariate Case  |
| 7.1       The One Covariate Case       7.1         7.2       The Multivariate Covariate Case       7.2         8       Thereforemetics       7.2   |
| 7.1       The One Covariate Case         7.2       The Multivariate Covariate Case         8       Transformation and NB Modela  |
| 7.2 The multivariate Covariate Case  |
|  |
| 8.1 Simulation Results   |
|  |
| <sup>1</sup> Department of Statistics, University of Wisconsin, Madison, WI 53706, email: doksum@st  |
| *Supported in part by NSF grant DMS-0505651  |

# 1. Introduction

We will focus on statistical inference for models where the distribution of the data can be expressed as a parametric function of unknown distribution functions.

# 1.1. Lehmann Type Models. Cox Regression

Suppose T is a random variable with a continuous distribution function F. For testing the null hypothesis  $H_0$ :  $F = F_0$ , Lehmann (1953) considered alternatives of the form

(1.1) 
$$F_{\theta}(\cdot) = C_{\theta}(F_0(\cdot)),$$

for some continuous distribution  $C_{\theta}(\cdot)$  on [0,1], which is known except for the parameter  $\theta$ . We consider the problem of estimating  $\theta$  when  $F_0(\cdot)$  is an unknown baseline distribution. In this case, if  $T_1, \dots, T_n$  are independent with  $T_i \sim C_{\theta_i}(F_0(\cdot))$ and we set  $U_i = F_0(T_i)$ , then  $U_i$  has distribution  $C_{\theta_i}(\cdot)$  Moreover  $R_i \equiv Rank(T_i) = Rank(U_i)$ , which implies that the distribution of any statistical method based on  $R_1, \dots, R_n$  will not depend on  $F_0$ .

For regression experiments with observations  $(T_i, \boldsymbol{x}_i), i = 1, \dots, n$ , where  $T_i$ is a response and  $\boldsymbol{x}_i$  is a vector of nonrandom covariates, Cox (1972) considered the parametrization  $\theta_i = g(\boldsymbol{\beta}^T \boldsymbol{x}_i)$  with  $g(\cdot)$  a known function and  $\boldsymbol{\beta}$  a vector of regression coefficients. He considered statistical inference procedures based on the Cox (1972, 1975) partial likelihood in very general frameworks. These procedures are based on generalized ranks and show how powerful ranks are in generating statistical inference procedures.

In this paper we consider a special case of (1.1) obtained from the Lehmann models  $[F_0(t)]^N$  and  $1 - [1 - F_0(t)]^N$  by letting N be a zero truncated Poisson variable whose parameter depends on covariates and regression coefficients. We call this model "SINAMI" after SIbuya (1968) and NAbeya and MIura (1972). For a subset of the parameter space, the model has proportional expected hazard rate (PEHR). We show that semiparametric likelihood methods for the SINAMI model give more weight to intermediate survival times than the Cox proportional hazard model which heavily weights long survival times. Recently developed algorithms for finding profile nonparametric maximum likelihood estimates (profile NPMLE's) are combined with the MM algorithm to produce estimates. In the two sample case, we carry out a Monte Carlo comparison of the NPMLE with a parametric MLE and a method of moment (MOM) estimate of the two sample parameter. The NPMLE is nearly unbiased but only about 70% as efficient in terms of root MSE as the parametric estimate if the parametric model is true. The MOM estimate is slightly less efficient than the NPMLE. 

1.2. Sklar Type Models. Copula Regression

З

Suppose X and Y are random variables with continuous joint distribution  $H(\cdot, \cdot)$ and marginals  $F_1(\cdot)$  and  $F_2(\cdot)$ . Sklar (1959) considered models that include models of the form

$${}^{2}_{3} (1.2) H_{\theta}(\cdot, \cdot) = C_{\theta}(F_{1}(\cdot), F_{2}(\cdot)), 3$$

for some continuous distribution  $C_{\theta}(\cdot, \cdot)$  on  $[0,1] \times [0,1]$ , which is known except for the parameter  $\theta$ . We consider the problem of estimating  $\theta$  when  $F_1(\cdot)$  and  $F_2(\cdot)$  are unknown baseline distributions. If we set  $U = F_1(X)$ ,  $V = F_2(Y)$ , then (U, V) has distribution  $C_{\theta}(\cdot, \cdot)$ , and  $C_{\theta}(\cdot, \cdot)$  is called a *copula*. Note that if  $(X_1, Y_1), \cdots, (X_n, Y_n)$  are independent with  $(X_i, Y_i) \sim H_{\theta_i}(\cdot, \cdot)$ , then  $R_i \equiv$  $Rank(X_i) = Rank(F_1(X_i))$  and  $S_i \equiv Rank(Y_i) = Rank(F_2(Y_i))$ , which shows that the distribution of any statistical method based on these ranks will not depend on  $(F_1, F_2)$ . This model extends in the natural way to the d-dimensional case.

In this paper we consider the bivariate normal copula model where  $C_{\theta}(u, v)$  $= \Phi_{\theta}(\Phi^{-1}(u), \Phi^{-1}(v))$  with  $C_{\theta}$  the bivariate  $N(0, 0, 1, 1, \theta)$  distribution. We also consider the multivariate normal copula model and show that in regression experiments it can be used to construct a "transform both sides regression" transformation model (copula regression model.) Klaassen and Wellner (1997) have shown that the normal scores correlation coefficient is semiparametrically efficient for the bivariate normal copula. We use simulations to compare this estimate with the profile MLE for the transform both sides Box-Cox regression model and a nonparametric estimate based on splines thereby augmenting the comparisons made by Zou and Hall (2002). The normal scores estimate is nearly as efficient as the parametric MLE for estimating median regression when the transform both sides Box-Cox model is correct. We also consider the performance of the estimates for models outside the copula regression model and find that the normal scores based estimate of median regression is remarkably robust with respect to both bias and variance. On the other hand, the profile MLE of median regression derived from the transform both sides Box-Cox model is very sensitive to deviations from the model. The nonparametric spline estimate is the best for extreme deviations from the copula regression model.

# 2. Proportional Hazard and Proportional Expected Hazard Rate Models

Interesting special cases of (1.1) are obtained by considering the distributions of  $T_1 = \min(T_{01}, \dots, T_{0k})$  and  $T_2 = \max(T_{01}, \dots, T_{0k})$  where  $T_{01}, T_{02}, \dots$  are i.i.d. as  $T_0 \sim F_0$ . Then, for  $k \geq 1$ ,

(2.1) 
$$T_1 \sim F_1(t) = 1 - [1 - F_0(t)]^k,$$

and

(2.2) 
$$T_2 \sim F_2(t) = [F_0(t)]^k,$$

with  $t > 0, k = 1, 2, \cdots$ . More general forms (Lehmann (1953); Savage (1956)) are

<sup>47</sup> (2.3) 
$$T_1 \sim F_1(t) = 1 - [1 - F_0(t)]^{\Delta},$$

and

50 50 50 50 50 51 
$$(2.4)$$
  $T_2 \sim F_2(t) = [F_0(t)]^{\Delta},$  51

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

З

with  $t > 0, \Delta > 0$ . Here (2.3) can be derived by considering two-sample models where the two samples follow distributions of the form (2.1) with different k's (Bickel and Doksum (2007), Problem 1.1.12.)

For  $T_1$ , the hazard rate is

(2.5) 
$$\lambda(t) \equiv \frac{f(t)}{1 - F(t)} = \triangle \frac{f_0(t)}{1 - F_0(t)} \equiv \triangle \lambda_0(t).$$

In regression experiments, we set  $\triangle_i = g(\boldsymbol{\beta}^T \boldsymbol{x}_i)$ , and note that (2.5) is the Cox proportional hazard (PH) model (Cox (1972)).

Nabeya and Miura (1972) proposed replacing k in (2.1) and (2.2) by a random variable. In particular, they considered  $T_1 = \min(T_{01}, \cdots, T_{0N})$ , where N is independent of  $T_{01}, T_{02}, \dots$ , and has a zero truncated Poisson( $\theta$ ) distribution with  $\theta > 0$ . They also considered  $T_2 = \max(T_{01}, \cdots, T_{0M}), T_{0i} \sim F_0$  where M is independent of  $T_{01}, T_{02}, \cdots$ , and has a zero truncated Poisson $(-\theta)$  distribution with  $\theta < 0.$ 

Using Sibuya (1968), they found

(2.6) 
$$T_1 \sim F_1(t) = \frac{1 - e^{-\theta F_0(t)}}{1 - e^{-\theta}}, \quad \theta > 0,$$

(2.7) 
$$T_2 \sim F_2(t) = \frac{1 - e^{-\theta F_0(t)}}{1 - e^{-\theta}}, \quad \theta < 0.$$
 21  
22

Combining (2.6) and (2.7), we get

(2.8) 
$$T \sim F(t) = \frac{1 - e^{-\theta}}{1 - e^{-\theta}}, \quad \theta \neq 0,$$

$$= F_0(t), \qquad \theta = 0.$$

Note that model (2.6) is a mixture of proportional hazard models for individuals with the same baseline hazard rate  $\lambda_0(\cdot)$  but different hazard factors  $\triangle$  in the factorization (2.5) of the hazard rate. Let  $\lambda(t; k)$  denote the hazard rate of  $T_1$  given N = k; then by (2.5)

(2.9) 
$$E\lambda(t;N) = \sum_{k=1}^{\infty} k\lambda_0(t)p_{\theta}(k) = \tau(\theta)\lambda_0(t), \ \theta > 0,$$

where  $p_{\theta}(x)$  is the zero truncated Poisson( $\theta$ ) probability and

(2.10) 
$$\tau(\theta) = E(N) = \frac{\theta}{1 - \exp(-\theta)}.$$

Thus (2.6) is a model with proportional expected hazard rate. Note that (2.8) does not have this property for  $\theta < 0$ . We will refer to (2.6) and (2.8) as the PEHR and SINAMI models, respectively.

**Remark 2.1**: In regression experiments, the traditional frailty models are also constructed by introducing a random element in the PH model. However, these models are different from the PEHR and SINAMI models. To see this recall that in the frailty model the conditional hazard rate given the covariate vector  $\boldsymbol{x}$  (see Oakes (1992)) for the history and interpretation of frailty models) is of the form

51 (2.11) 
$$\lambda_W(t|\boldsymbol{x}) = \lambda_0(t)W \exp[\boldsymbol{\beta}^T \boldsymbol{x}],$$

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

З

Doksum and Ozeki

where W is a random effect that incorporates potential unobservable covariates that represent frailties. Semiparametric optimality theory for model (2.11) has been developed by Kosorok, Lee, and Fine (2004).

З

Consider model (2.5) with  $\Delta = N$  and N a zero truncated Poisson( $\theta$ ) random variable with  $\theta = g(\beta^T \boldsymbol{x})$ , i.e., the conditional hazard rate given  $\boldsymbol{x}$  is

(2.12) 
$$\lambda^{(N)}(t|\boldsymbol{x}) = N\lambda_0(t).$$

Here N plays the role of  $W \exp[\beta^T x]$  in (2.11). However, (2.11) and (2.12) are different because N is an integer and  $W \exp[\beta^T x]$  is not when  $\beta \neq 0$ . In model (2.12), N represents the effect of both observed covariates and frailties. In deriving the distribution function (2.6), the unobservable covariates are averaged out, that is, we compute  $P(T \leq t) = E[P(T \leq t|N)]$ .

**Remark 2.2**: Model (2.8) was considered by Bell and Doksum (1966), Example 5.2 and Table 8.1) and Ferguson (1967, p.257, Problem 5.7.7) without any of the above interpretations. Nabeya and Miura (1972) did not use any proportional hazard or frailty interpretation. These concepts had not been invented yet.

Fig.1 gives a plot of the relative hazard rate  $\lambda(t|x = 1)/\lambda(t|x = 0)$  with  $\theta = -3, -1.5, 1.5, 3$  for model (2.8) with  $\theta = \beta x$ ,  $F_0(t) = 1 - \exp(-t)$ , t > 0, and

22  
23 (2.13) 
$$\lambda(t|x) = \frac{\theta f_0(t)}{1 - e^{-\theta(1 - F_0(t))}}.$$

In the PEHR and SINAMI models, the hazard ratio between two covariate values converge to unity as time increases. This explains why the likelihoods for these models give less weight to long survival times than the likelihood for the Cox model (see Section 3). The hazard rate is decreasing for the PEHR model for any continuous  $F_0$ .

#### 3. Rank, Partial and Marginal Likelihood

In regression experiments, we observe  $(T_i, \boldsymbol{x}_i)$ ,  $i = 1, \dots, n$ , where  $T_1, \dots, T_n$  are independent responses and  $\boldsymbol{x}_i$  is a nonrandom covariate vector. In the proportional hazard model, it is customary to use model (2.5) for  $T_i$  with  $\Delta_i = \exp(\boldsymbol{\beta}^T \boldsymbol{x}_i)$ because  $\Delta_i$  needs to be positive. In the PEHR model,  $\theta_i = \boldsymbol{\beta}^T \boldsymbol{x}_i$  is a possible parametrization, but  $\theta_i = \exp(\boldsymbol{\beta}^T \boldsymbol{x}_i)$  could also be used. Let  $\boldsymbol{R} = (R_1, \dots, R_n)$ where  $R_i = \operatorname{Rank}(T_i)$ , then  $l_r(\boldsymbol{\beta}) = P(\boldsymbol{R} = \boldsymbol{r})$  is the rank likelihood (Hoeffding (1951)).

We first consider the one covariate case. Using the rank likelihood, the locally most powerful (LMP) rank test statistic for  $H_0$ :  $\beta = 0$  versus  $H_1$ :  $\beta > 0$  is (approximately) for the Cox model (Savage (1956), 1957), Cox (1964)), Oakes and Jeong (1998)):

(3.1) 
$$\sum_{i=1}^{n} \left[-\log(1-\frac{R_i}{n+1})\right](x_i - \bar{x}) \quad \text{(Savage or log rank)},$$

and for the PEHR and SINAMI models, the LMP rank test statistics is (Bell and Doksum (1966), Ferguson (1967), Nabeya and Miura (1972)):

<sup>50</sup> (3.2) 
$$\sum_{i=1}^{n} \frac{R_i}{n+1} (x_i - \bar{x})$$
 (Wilcoxon type), <sup>50</sup> <sup>50</sup> <sub>51</sub>

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009



FIG 1. SINAMI and PEHR hazard ratios for  $\theta = -3, -1.5, 1.5, 3$ .

where  $R_i = \text{Rank}(T_i)$ . The log rank statistic gives more weights to large observations, that is, in survival analysis, to those that live longer, while the Wilcoxon statistics is even handed.

In order to compare how much relative weight is given to the small, in between, and large observed survival times for the PH and PEHR models, we next consider the rank likelihood for d covariates. Note that if  $h(\cdot)$  is decreasing, then  $\operatorname{Rank}(h(T_i)) = n + 1 - R_i$ . For the proportional hazard model, transform  $T_i$  by  $U_i = 1 - F_0(T_i)$ , then by (2.3) we have  $f_{U_i}(u) = \Delta_i u^{\Delta_i - 1}, 0 < u < 1$ . Hoeffding (1951) formula shows,

(3.3) Rank lkhd = 
$$\prod_{i=1}^{n} \bigtriangleup_{i} \int_{0 < u_{1} < u_{2} < \cdots < u_{n} < 1} \prod_{i=1}^{n} u_{i}^{\delta_{i}-1} du_{1} \cdots du_{n}$$

where  $\delta_i = \Delta_{b_i}$  and  $b_i$  = index on the T with rank n + 1 - i = reverse anti-rank. It follows that

(3.4) Rank lkhd 
$$\propto \prod_{i=1}^{n} \frac{\triangle_i}{\sum_{k:T_k \ge T_{(i)}} \triangle_k},$$

that is, the familiar Cox (1972, 1975) partial likelihood formula. Here  $\{k : T_k \geq T_{(i)}\}$  = patients at risk at time  $T_{(i)}$  where  $T_{(i)}$  is the *i*th ordered survival time. Kalbfleich and Prentice (1973, 2002) called the rank likelihood the marginal likelihood and extended it to censored data.

For the PEHR model, transform 
$$T_i$$
 by the decreasing function  
 $U_i = a^{-1} \{\exp\{-F_0(T_i)\} - b\}$ , then, we have  $f_i(u) = a\tau_i(au+b)^{\theta_i-1}, 0 < u < 1$ , 51

З

where 
$$b = e^{-1}$$
,  $a = 1 - b$ , and  $\tau_i \equiv \tau(\theta_i)$ . Let  $\gamma_i = \theta_{b_i}$ ,  $b_i$ =index on the T with rank  $n + 1 - i$ . Then,

(3.5) Rank lkhd 
$$\propto (\prod_{i=1}^{n} \tau_i) \int_{\substack{0 < u_1 < u_2 < \cdots < u_n < 1}} \prod_{i=1}^{n} (au_i + b)^{\gamma_i - 1} du_1 \cdots du_n.$$

If we perform the integration, we find that the likelihood for the PEHR model is similar to the likelihood for the Cox model except that in addition to terms involving  $\{k : T_k \geq T_{(i)}\}, i = 1, \dots, n, \text{ it includes terms involving } \{k : T_{(i)} \leq T_k \leq T_{(j)}\}, i = 1, \dots, n, j = 1, \dots, n, i \neq j$ . That is, the PEHR likelihood gives more weight to the intermediate survival times than the Cox likelihood.

Computationally, the Cox rank likelihood is easier than the PEHR rank likelihood. However, we can handle the PEHR rank likelihood with available algorithms and software (e.g. MATLAB.) More generally,  $F(x) = C_{\theta}(F_0(x))$  type models, originally considered by Lehmann (1953), can be handled effectively by considering the profile NP likelihood of the next section (e.g. Tsodikov and Gabribotti (2007)), Zeng and Lin (2007) ).

З

# 4. Profile NP Likelihood

Andersen, Borgan, Gill, and Keiding (1996), Bickel, Klaassen, Ritov, and Wellner (1993, 1996), van der Vaart (1998), Murphy and van der Vaart (2000), Tsodikov and Garibotti (2007), Zeng and Lin (2007) and many others considered the problem of finding the MLE of all the parameters in a semiparametric model. It is useful to divide the procedure into two steps by grouping parameters into two groups. Suppose the distribution function of T is of the form  $P(T \leq t) = F(\theta, \eta(t))$ , where  $\theta \in \mathbb{R}^d$  and  $\eta(\cdot)$  is a nondecreasing function. If we assume temporarily that  $\eta(\cdot)$  has a positive derivative  $\eta'(t)$  for  $t \in \{t_1, \dots, t_n\}$ , then the likelihood is

$$\prod_{i=1}^{n} \eta'(t_i) f(\theta, \eta(t_i)),$$

where  $f(\theta, \eta) = \partial F(\theta, \eta) / \partial \eta$ . The NP likelihood we consider is of the form

$$L_{NP}(\theta, \eta) = \prod_{i=1}^{n} \eta\{t_i\} f(\theta, \eta[t_i]),$$

where  $\eta[t_i] = \sum_{j \leq i} \eta\{t_i\}$  is a step function with positive jumps  $\eta\{t_i\}$  at the data points  $t_i, i = 1, \dots, n$ .

We assume that

$$af(\theta, a) \to 0 \quad as \ a \to \infty$$

Next we fix  $\theta$ , and define  $\hat{\eta}_{\theta}\{t_i\}$  as

47 (4.1) 
$$\hat{\eta}_{\theta}\{.\} = ARG \ MAX_{\eta\{.\}}L_{NP}(\theta,\eta).$$

Set

(4.2) 
$$PROF \ NPLIK = l(\theta) = MAX_{\eta\{\cdot\}}L_{NP}(\theta,\eta),$$
 51

and solve

$$\hat{\theta} = ARG MAX \ l(\theta).$$

Next estimate  $\eta\{\cdot\}$  as  $\hat{\eta}_{\hat{\theta}}\{\cdot\}$ . In the Lehmann model (1.1), the NP likelihood is

$$\prod_{i=1}^{n} F_0\{x_i\} C'_{\theta}(\sum_{j \le i} F_0\{x_j\}),$$

The method is similar to finding the empirical MLE, Owen (1988, 2001), and profile (partial) MLE's as in Andersen, et al. (1996), and Murphy and van der Vaart (2000).

**Remark 4.1**: Note that when  $P(T \leq t) = F(\theta, \eta(t))$ , (4.1), (4.2), and (4.3) do not depend on the values of  $t_1, \dots, t_n$ . In regression experiments, they will depend on the ranks of  $t_1, \dots, t_n$ . For example, see (4.4) and (5.1). This is in contrast to the Hodges and Lehmann (1963) approach that uses estimating equations based on rank test statistics to obtain estimates of parameters. In this Hodges-Lehmann "rank inversion" approach, estimates are functions of the "raw" data rather than the ranks.

As an example that will guide the algorithm for the PEHR model, consider the Cox model. Set  $\Lambda(t) = -\log(1 - F_0(t))$ , then,

$$L_{NP}(\boldsymbol{\beta}, \Lambda) = \prod_{i=1}^{n} e^{\boldsymbol{\beta}^{T} x_{i}} \Lambda\{t_{i}\} e^{-(e^{\boldsymbol{\beta}^{T} x_{i}})\Lambda[t_{i}]},$$

where

$$\Lambda[t_i] = \sum_{j:t_i \le t_i} \Lambda\{t_j\}.$$

Using calculus, we find

$$\hat{\Lambda}_{\beta}\{t_i\} = ARG \ MAX_{\Lambda\{t_i\}} \ L_{NP}(\boldsymbol{\beta}, \Lambda) = (\sum_{j:t_j \ge t_i} e^{\boldsymbol{\beta}^T x_j})^{-1},$$

and

(4.4) 
$$l(\boldsymbol{\beta}) = PROF \ NPLIK = \prod_{i=1}^{n} \frac{e^{\boldsymbol{\beta}^{T} x_{i}}}{\sum_{j:t_{i} \geq t_{i}} e^{\boldsymbol{\beta}^{T} x_{j}}}.$$

This is exactly the same as the rank, the partial, and the marginal likelihood.

#### 5. Profile NP Likelihood for the PEHR Model

Consider model (2.6) with  $\theta > 0$ . Set  $\tau(\theta) = \theta [1 - e^{-\theta}]^{-1}$ , then

(5.1) 
$$L_{NP}(\boldsymbol{\theta}, F_0) = [\prod_{i=1}^n \tau(\theta_i)] \prod_{i=1}^n F_0\{t_i\} e^{-\theta_i F_0[t_i]},$$
(5.1) (5.1)

47 where  $F_0[t_i] = \sum_{j:t_j \le t_i} F_0\{t_j\}$ . Now set  $p_i \equiv F_0\{t_i\}$  and maximize with respect 48 to  $p_1, \dots, p_n$  with  $\theta$  fixed. The maximization problem looks very similar to Cox 49 model maximization except for the constraint  $\sum p_i = 1$ . We handle this constraint 50 by writing  $F_0(t) = 1 - \exp[-\Lambda(t)]$  with  $\Lambda(t)$  unconstrained except for  $\Lambda(t) \ge 0$  and 51 by using a new approach based on the MM algorithm.

З

# 5.1. The MM Algorithm

Lang, Hunter and Yang (2000) introduced a concept called the MM algorithm. Its idea is that instead of maximizing a complicated original objective function, use a simpler surrogate function so that each iteration is faster and guarantees that the original objective function increases. Given the original objective function  $l(\mathbf{h})$  for a maximization problem, a surrogate function  $g(\mathbf{h}|\mathbf{h}_{old})$  must satisfy two properties:

$$g \qquad (5.2) \qquad \qquad l(\boldsymbol{h}_{old}) = g(\boldsymbol{h}_{old}|\boldsymbol{h}_{old})$$

$$(5.3) l(\mathbf{h}) \ge g(\mathbf{h}|\mathbf{h}_{old})$$

The EM algorithm is a special case of the MM algorithm. A practical implementation issue of the MM algorithm is that we have to find a nice surrogate function case by case.

Now we construct a surrogate function based on Tsodikov (2003). Suppose we can write  $l(\cdot)$  in the form  $l(\mathbf{h}) = B(\mathbf{h}) - A(\mathbf{h})$  for some parameter vector  $\mathbf{h} > \mathbf{0}$ , where A and B are differentiable concave functions. Then by the concavity property,

(5.4) 
$$g(\boldsymbol{h}|\boldsymbol{h}_{old}) = B(\boldsymbol{h}) - A(\boldsymbol{h}_{old}) - \nabla^T A(\boldsymbol{h}_{old})(\boldsymbol{h} - \boldsymbol{h}_{old}),$$

where  $\nabla^T A(\mathbf{h}) = \partial A / \partial \mathbf{h}$  is the gradient of A, satisfies (5.2) and (5.3), and  $g(\mathbf{h}|\mathbf{h}_{old})$ is a surrogate function for  $l(\mathbf{h})$ . Differentiating (5.4) gives

(5.5) 
$$\nabla^T B(\boldsymbol{h}_{new}) = \nabla^T A(\boldsymbol{h}_{old}).$$

Solve (5.5) for  $h_{new}$ . Iterate the procedure until there is a minimal change in h.

# 5.2. The MM Algorithm for the PEHR Model with $\theta \ge 0$ (SINAMI with $\theta \ge 0$ )

Let  $\theta_i = g(\boldsymbol{\beta}^T \boldsymbol{x}_i)$ , for some known function  $g(\cdot) \ge 0$ . When  $\theta_i = 0$ , the distribution function of  $T_i$  is  $F_0(t)$ . Let  $F_0(t) = 1 - \exp[-\Lambda(t)]$ ,  $h_k = \Lambda\{t_k\}$ , and  $\Lambda[t_i] = \sum_{k=1}^i h_k$  with  $\Lambda(t) \ge 0$  and  $h_k \ge 0$ . Then for a temporarily fixed numerical vector  $\boldsymbol{\beta}$ ,

$$l(\boldsymbol{h}) = log[L_{NP}(\boldsymbol{eta}, \boldsymbol{h})] = \sum_{i=1}^{n} \log \tau(\theta_i) + \sum_{i=1}^{n} \log h_i$$

(5.6)

$$-\sum_{i=1}^{n} [\sum_{k=1}^{i} h_k + \theta_i (1 - \exp(-\sum_{k=1}^{i} h_k))].$$

42 Now we can write  $l(\mathbf{h}) = B(\mathbf{h}) - A(\mathbf{h})$  with

(5.7) 
$$B(h) = \sum_{i=1}^{n} \log h_i,$$
 43  
44  
45

47  
48  
49
(5.8)
$$A(\boldsymbol{h}) = \sum_{i=1}^{n} [\sum_{k=1}^{i} h_k + \theta_i (1 - \exp(-\sum_{k=1}^{i} h_k))].$$
47  
48  
49  
49

Here we may ignore  $\Sigma \log \tau(\theta_i)$  because we maximize (5.6) w.r.t. **h**.

 $B(\mathbf{h})$  and  $A(\mathbf{h})$  are concave, because for  $0 \leq t \leq 1$ ,  $B(t\mathbf{h}_a + (1-t)\mathbf{h}_b) \geq 0$  $tB(\mathbf{h}_a) + (1-t)B(\mathbf{h}_b)$  and by mathematical induction,  $A(t\mathbf{h}_a + (1-t)\mathbf{h}_b) \ge tA(\mathbf{h}_a) + tB(\mathbf{h}_a) + tB(\mathbf{h$  $(1-t)A(\mathbf{h}_b)$  hold. Note that

(5.9) 
$$\partial B/\partial h_j = 1/h_j, \ j = 1, ..., n,$$

(5.10) 
$$\partial A/\partial h_j = \sum_{i=1}^n (1 + \theta_i (1 - \exp(-\sum_{k=1}^i h_k))) 1(j \le i).$$

Using (5.5), (5.9), and (5.10), update  $h_i$ ,  $j = 1, \dots, n$ , at the same time,

(5.11) 
$$h_{j,new} = \left[\sum_{i=1}^{n} (1 + \theta_i (1 - \exp(-\sum_{k=1}^{i} h_{k,old}))) 1(j \le i)\right]^{-1}.$$

Iterate (5.11) until there is a minimal change in  $l(\hat{h}_{new})$ ; call the result  $\hat{h}_A$ (Approximated profile NPMLE). Note that we call  $\hat{h}_A$  approximated profile NPMLE because  $\hat{h}_A$  is obtained by fixing  $\beta$ . This approximation is necessary because there is no closed form  $\hat{h}$  w.r.t.  $\beta$ . Next set  $l(\beta) = log[L_{NP}(\beta, \hat{h}_A)]$  and maximize w.r.t. β.

# 5.3. The MM Algorithm for the SINAMI Model with $\theta < 0$

Consider model (2.7) with  $\theta_i = g(\boldsymbol{\beta}^T \boldsymbol{x}_i)$ , for some known function  $g(\cdot) \leq 0$ . In this case we can use the algorithm of Section 5.2. To see this, suppose T satisfies model (2.7) with parameter  $\theta_2 < 0$ . Set  $V = 1 - F_0(T)$ , then V satisfies model (2.6) with parameter  $\theta_1 = -\theta_2$ . Moreover, the rank of  $1 - F_0(T_i)$  is  $n + 1 - R_i$ .

# 5.4. The MM Algorithm for the SINAMI Model with $\theta \in R$

Consider model (2.7) with  $\theta_i = g(\beta^T x_i)$ , for some known function  $g(\cdot) \in R$ . In this case we can not use the transformation in Section 5.3 because it changes the likelihood and the monotonicity of the likelihood as a function of  $\theta$  does not necessarily hold. Instead, we modify the algorithm as follows: If the value  $\theta_j$  in the jth iteration is positive, use the MM algorithm in Section 5.2. to find  $\hat{h}_i$ . If  $\hat{\theta}_i < 0$ , then (5.6) implies that finding the maximizer h is a convex optimization problem which produces  $h_i$ .

# 5.5. Profile NPMLE Implementation

Successful convergence of the MM algorithm depends on a good starting point  $(\hat{\theta}, \hat{h})$ . We consider the two sample problem:

(5.12) 
$$T_{0,i} \sim F_0(t), \ i = 1, ..., n_0,$$

$$\begin{array}{c} {}^{49}\\ {}^{50}\\ {}^{51}\\ {}^{51} \end{array} (5.13) \qquad T_{1,i} \sim F(t) = \left\{ \begin{array}{c} \frac{1-e^{-\theta F_0(t)}}{1-e^{-\theta}}, \qquad (\theta>0), \ i=1,...,n_1, \\ F_0(t), \qquad (\theta=0), \ i=1,...,n_1, \end{array} \right. \begin{array}{c} {}^{49}\\ {}^{50}\\ {}^{50}\\ {}^{51}\\ {}^{51}\end{array}$$

З

Doksum and Ozeki

where  $F_0(\cdot)$  is an unknown distribution with density  $f_0$ . Note that the density of F(t) is

(5.14) 
$$f(t;\theta) = \begin{cases} \tau(\theta) f_0(t) e^{-\theta F_0(t)}, & (\theta > 0) \\ f_0(t) e^{-\theta F_0(t)}, & (\theta > 0) \end{cases}$$

14) 
$$f(t;\theta) = \begin{cases} f_0(t), & (\theta = 0). \end{cases}$$

We use an algorithm to find  $(\hat{\theta}, \hat{h})$  where  $\theta = \beta$  in this case. For fixed  $F_0$ , (5.14) gives an MOM estimating equation for  $\theta$ . We plug in an estimate  $\hat{F}_0$  for  $F_0$ , and use the following algorithm:

$$step (1) : \hat{F}_0 \to step (2) : \hat{\theta} \to step (3) : (\hat{\theta}_A \leftrightarrow \hat{h}_A)_{until \ \hat{\theta}_A \ converges.}$$

Here  $\hat{\theta}_A$  and  $\hat{h}_A$  are approximated profile NPMLE's from Section 5.2. The details are as follows:

**Step** (1): Compute the empirical distribution  $\hat{F}_0(t)$  based on  $T_{0,i}$  only:

(5.15) 
$$\hat{F}_{0,[0]}(t) \equiv \hat{F}_0(t) = \frac{1}{1+n_0} \sum_{i=1}^{n_0} 1(T_{0,i} \le t).$$

Here  $\hat{F}_0(t) \rightarrow_{a.s.} F_0(t)$  uniformly in t as  $n_0 \rightarrow \infty$ . The subscript [0] indicate iteration zero (starting point) for step (3). Note that the one-to-one relation between  $\hat{F}_0$  and  $\hat{\Lambda}_0$  is used to obtain  $\hat{h}$  by solving for h in the equations:

$$\hat{F}_0(t_i) \equiv \hat{F}_{0,i} = \exp(-\hat{\Lambda}_i), where$$
 25

$$\hat{\Lambda}[t_i] \equiv \quad \hat{\Lambda}_i = \sum_{k=1}^i \hat{h}_k, \ i = 1, ..., n_0.$$
<sup>26</sup>
<sup>27</sup>
<sup>28</sup>
<sup>28</sup>

З

**Step** (2): Solve for  $\hat{\theta}$  based on  $\hat{F}_{0,[0]}(t)$ :

=

$$ar{T}_1 = - au( heta) \int y e^{- heta \hat{F}_{0,[0]}(y)} d\hat{F}_{0,[0]}(y)$$

where  $T_{0,(i)}$  is an order statistics of  $T_{0,1}, \cdots, T_{0,n_0}$  and  $\overline{T}_1 = \sum_{i=1}^{n_1} T_{1,i}/n_1$ . The solution is uniquely determined because the distribution function (5.13) is monotone increasing in  $\theta$  and hence its mean is monotone decreasing in  $\theta$ . If a model is  $\theta \ge 0$ 

and  $\hat{\theta} < 0$ , set  $\hat{\theta} = 0$ . If a model is  $\theta < 0$  and  $\hat{\theta} > 0$ , set  $\hat{\theta} = 0$ . Step (3): Compute  $\theta_A$  and  $h_A$  as follows:

$$\begin{array}{cc} 45 \\ 46 \end{array} (5.18) \qquad \qquad (\hat{\theta}_A \leftrightarrow \hat{h}_A)_{until \ \hat{\theta}_A \ converges.} & 45 \\ 46 \end{array}$$

The first iteration  $\hat{\theta}_{A,[1]}$ , is obtained by maximizing (5.6) with  $\hat{\theta}$  as a starting point, i.e.,  $\hat{\theta}_{A,[0]} = \hat{\theta}$  and with fixed  $\hat{h}_{[0]}$  obtained from (5.16) and  $\hat{F}_{0,[0]}$ , i.e., 

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

Then by the MM algorithm in Section 5.2 with  $\beta_0 = \theta_0 = \hat{\theta}_{A,[1]}$  (see (5.11)), obtain  $\hat{h}_{[1]}$  using the starting point  $\hat{h}_{[0]}$ . Next obtain

(5.20) 
$$\hat{\theta}_{A,[2]} = \arg\max_{\boldsymbol{\rho}} \left\{ l(\boldsymbol{\theta}, \hat{\boldsymbol{h}}_{A,[1]}) : \boldsymbol{\theta} \ge 0 \right\}$$

with starting point  $\hat{\theta}_{A,[1]}$ . Then by the MM algorithm in Section 5.2 with  $\beta_0 = \theta_0 = \hat{\theta}_{A,[2]}$ , obtain  $\hat{h}_{A,[2]}$  with starting point  $\hat{h}_{A,[1]}$ . Repeat the procedure to get  $\hat{\theta}_{A,[j]}$  and  $\hat{h}_{A,[j]}$  until convergence, i.e.,  $|\hat{\theta}_{A,[j]} - \hat{\theta}_{A,[j-1]}| < \epsilon$  for some small  $\epsilon$ .

Numerical optimizations for  $\hat{\theta}$  and  $\hat{\theta}_A$  are carried out by the MATLAB fmincon() function.

**Remark 5.1**: For fixed  $\hat{h}$ ,  $l(\theta, \hat{h})$  is strictly concave and have a unique maximum.

**Remark 5.2**: The estimate of  $\beta$  in the Cox model that we have discussed is asymptotically optimal in the semiparametric sense (Begun, Hall, Huang, and Wellner (1983), Bickel et al. (1993, 1998), van der Vaart (1998)), Murphy and van der Vaart (2000)). These references and others give results that can be used to check the semiparametric asymptotic optimality of the profile NPMLE in the PEHR model.

**Remark 5.3**: Transformation models. We can show that the PEHR is a special case of transformation models as follows: Let  $F_{\lambda}$  be the exponential  $(\lambda)$  distribution function and define

(5.21) 
$$G_0(y|\boldsymbol{x}) = \frac{1 - e^{-\theta F_{\lambda}(y)}}{1 - e^{-\theta}}, \quad y > 0,$$

where  $\theta = g(\boldsymbol{x}, \boldsymbol{\beta})$ . Let  $\psi$  be an increasing function from  $[0, \infty)$  to  $[0, \infty)$  and define the transformation model

(5.22) 
$$G(y|\boldsymbol{x}) = G_0(\psi(y)|\boldsymbol{x})$$

This model is of the form (2.6) with  $F_0 = F_\lambda \psi(t)$ . Klaassen (2007) gives results for general transformation models that can be used to check semiparametric asymptotic efficiency of estimates of  $\beta$  in the model (2.6).

# 5.6. Estimation of the Variance of the Profile NPMLE

Hypothesis tests and confidence intervals require standard errors (estimates of the standard deviation) of  $\hat{\theta}_A$ . An algorithm developed by Tsodikov and Garibotti (2007) combined with the preceding algorithm allows us to compute the profile information matrix which is the observed information matrix derived from the profile likelihood. This provides standard errors  $SE(\hat{\theta}_A)$  of  $\hat{\theta}_A$ .

#### 6. Simulation Results

# 6.1. PEHR Model Estimates



З

| θ                         |       | 2     |       |       | 3     |       |  |
|---------------------------|-------|-------|-------|-------|-------|-------|--|
| n                         | 100   | 200   | 300   | 100   | 200   | 300   |  |
| $E[\hat{\theta}_{MOM}]$   | 1.97  | 2.00  | 2.02  | 3.05  | 3.00  | 3.01  |  |
| $E[\hat{\theta}_A]$       | 1.99  | 2.00  | 2.01  | 3.07  | 3.00  | 3.01  |  |
| $E[\hat{\theta}_{PAR}]$   | 2.02  | 2.00  | 2.01  | 3.03  | 3.01  | 3.00  |  |
| $SD(\hat{\theta}_{MOM})$  | 0.819 | 0.551 | 0.479 | 0.972 | 0.635 | 0.506 |  |
| $SD(\hat{\theta}_A)$      | 0.788 | 0.541 | 0.466 | 0.960 | 0.618 | 0.497 |  |
| $SD(\hat{\theta}_{PAR})$  | 0.541 | 0.387 | 0.317 | 0.605 | 0.431 | 0.335 |  |
| $E[SE(\hat{\theta}_A)]$   | 0.752 | 0.527 | 0.430 | 0.837 | 0.578 | 0.471 |  |
| $MSE[\hat{\theta}_{MOM}]$ | 0.671 | 0.303 | 0.230 | 0.949 | 0.403 | 0.256 |  |
| $MSE[\hat{\theta}_A]$     | 0.621 | 0.292 | 0.217 | 0.926 | 0.383 | 0.248 |  |
| $MSE[\hat{\theta}_{PAR}]$ | 0.289 | 0.150 | 0.101 | 0.367 | 0.186 | 0.112 |  |
| TABLE 1                   |       |       |       |       |       |       |  |

| PEHR    | model    | simulation   | estimates    | (MC = 1000.  | $\theta = 2.3$ | 3) |
|---------|----------|--------------|--------------|--------------|----------------|----|
| 1 11110 | 11000000 | 001100000010 | 000000000000 | 1110 - 10000 | $v = \pm, v$   | 91 |

З

300.  $T_{0,i} \sim EXP(1)$ ,  $i = 1, \dots, n_0$ ,  $T_{1,i} \sim PEHR(\theta = 2, \text{ or } 3)$ ,  $i = 1, \dots, n_1$ ,  $n_0 = n_1 = n/2$ , iid. We compute Monte Carlo estimates of the expected values, standard deviations (SD's), and MSE's of  $\hat{\theta}_{MOM}$ ,  $\hat{\theta}_A$ , and  $\hat{\theta}_{PAR}$  where  $\hat{\theta}_{PAR}$  is the parametric model MLE, obtained by assuming that  $F_0$  is known and equal to the EXP(1) distribution.

We also compute the Monte Carlo estimates of the expected values  $E[SE(\hat{\theta}_A)]$  of the standard errors computed as described in Section 5.5. Table 1 summarizes the result.

Overall  $\hat{\theta}_{MOM}$ ,  $\hat{\theta}_A$ , and  $\hat{\theta}_{PAR}$  have almost no bias in the estimation of  $\theta = 2$  or 3. As expected, the parametric model estimate  $\hat{\theta}_{PAR}$  has the smallest MSE. The approximated profile NPMLE  $\hat{\theta}_A$  has a smaller MSE than  $\hat{\theta}_{MOM}$ , but the difference is small. The approximation to  $SD(\hat{\theta}_A)$  is very good and improves as the sample size increases.

#### 6.2. Model Fit for Misspecified Model

Next consider a model that is neither a Cox PH model nor a PEHR model: i.e.,  $T_{0,i} \sim EXP(1)$  and the true model for  $T_{1,i}$  is : Case 1, Gamma(shape=0.5, scale=0.5), and the : Case 2, Weibull(shape=0.5, scale=0.2). Here the target values  $\theta$  and h are those that minimize the Kullback-Leibler divergence between the true distribution and the model class of distributions (Doksum, Ozeki, Kim and Neto (2007)).

Fig. 2 shows that PEHR gives better fit than the Cox model.

# 7. Estimation in the Normal Copula Model

#### 7.1. The One Covariate Case

46 Assume that the pair (X,Y) has a joint density f(x,y) with respect to Lebesgue 47 measure on  $\mathbb{R}^2$  and a joint distribution function F(x,y). Let  $F_1$  and  $F_2$  be the 48 marginal distribution functions of X and Y, respectively, and let  $\Phi$  denote the 49 standard normal distribution function. Consider the transformations  $X \to Z =$ 50  $\Phi^{-1}(F_1(X)), Y \to W = \Phi^{-1}(F_2(Y))$ . Then the marginal distributions of Z and 51 W are standard normal. The bivariate normal copula model  $\mathcal{F}$  is defined by the

З



FIG 2. Cox and PEHR estimated hazard ratio when  $F_0 \sim EXP(1)$  and  $F \sim Gamma$  or Weibull.

assumption that the joint distribution of (Z,W) is bivariate normal with zero mean, unit variance, and correlation coefficient  $\rho$ . That is,

$$\mathcal{F} = \{ F : (\Phi^{-1}(F_1(X)), \Phi^{-1}(F_2(Y))) \sim N(0, 0, 1, 1, \rho) \},\$$

where  $F_1$  and  $F_2$  are the marginals of F. Let  $\{(X_i, Y_i), i = 1, 2, \dots, n\}$  be independent and identically distributed with distribution function  $F \in \mathcal{F}$ , and set  $Z_i = \Phi^{-1}(F_1(X_i)), W_i = \Phi^{-1}(F_2(Y_i)), i = 1, 2, \dots, n$ . If we (temporarily) assume that  $F_1$  and  $F_2$  are known, then because  $E(ZW) = \rho$ , a method of moments "estimate" of  $\rho$ , is  $r_{MOM} = n^{-1} \sum_{i=1}^{n} Z_i W_i$ . The asymptotic distribution of  $\sqrt{n}(r_{MOM} - \rho)$  is  $N(0, 1 + \rho^2)$  when  $F \in \mathcal{F}$ . Assuming  $F_1$  and  $F_2$  known, the asymptotic distribution of  $\sqrt{n}(r_{MLE} - \rho)$ , where  $r_{MLE}$  is the maximum likelihood "estimate" of  $\rho$  is  $N(0, (1 - \rho^2)^2/(1 + \rho^2))$ . The asymptotic variance  $(1 - \rho^2)^2/(1 + \rho^2)$ of  $r_{MLE}$  is smaller than the asymptotic variance  $(1 - \rho^2)^2$  of the usual Pearson correlation coefficient  $r_P$  and much smaller than the asymptotic variance  $1 + \rho^2$  of  $r_{MOM}$ .

Note that  $\mathcal{F}$  is invariant under coordinate-wise increasing transformations. That is, if  $(X, Y) \sim F \in \mathcal{F}$  and  $U = h_1(X), V = h_2(Y)$  with  $h_1$  and  $h_2$  increasing, then the distribution G of (U,V) is in  $\mathcal{F}$ . If we want methods that are invariant under such transformations, we must use statistics based on the ranks defined in Section 1.

Suppose next that  $F_1$  and  $F_2$  are unknown. It may then make sense to replace the ordered Z's and W's by their expected values. This leads to the Fisher and Yates (1938) or normal scores  $E(Z_{(i)})$ ,  $i = 1, \dots, n$  where  $Z_{(1)}, \dots, Z_{(n)}$  are N(0,1) order statistics. We write  $a(i) = E(Z_{(i)})$ . An accurate approximation to  $E(Z_{(i)})$  is  $\Phi^{-1}[(i - 3/8)/(n + 1/4)]$ , e.g. Cox (2006).

47 Let  $Z'_i = a(R_i)$ ,  $W'_i = a(S_i)$  where  $R_i$  and  $S_i$  are the ranks of  $X_i$  and  $Y_i$  when 48 the X's and Y's are ranked separately. Then we obtain estimates  $\hat{\rho}_{MOM}$ ,  $\hat{\rho}_{MLE}$ , 49 and  $\hat{\rho}_P$  of  $\rho$  when  $F_1$  and  $F_2$  are unknown by replacing  $Z_i$  and  $W_i$  by  $Z'_i$  and  $W'_i$  in 50  $r_{MOM}$ ,  $r_{MLE}$  and  $r_P$ . In this case  $\hat{\rho}_{MOM}$ ,  $\hat{\rho}_P$  are nearly identical and asymptotically 51 equivalent, but they are different from  $\hat{\rho}_{MLE}$ . We will use  $\hat{\rho}_P$  because it is slightly 

З

less biased, and denote it by  $\hat{\rho}_{NS}$  where NS signifies normal scores. Thus

<sup>2</sup>  
<sub>3</sub> (7.1) 
$$\hat{\rho}_{NS} = \sum Z'_i W'_i / \sum a^2(i).$$

It follows from Bhuchongkul (1964) that based on the rank likelihood,  $\hat{\rho}_{NS}$  is, uniformly in  $F_1$  and  $F_2$ , a locally most powerful test statistics in the bivariate normal copula model. Zou and Hall (2002) gave an asymptotic extension of this result. They also computed the rank likelihood estimate of  $\rho$  in the bivariate normal copula model using an improved version of the likelihood sampler in Doksum(1987).

Klaassen and Wellner (1997) found  $\sqrt{n}(\hat{\rho}_{NS}-\rho) \rightarrow_d N(0,(1-\rho^2)^2)$  in the copula model  $\mathcal{F}$  with  $F_1$  and  $F_2$  unknown; the same as for the Pearson correlation in the bivariate normal model. In fact, in a bivariate normal $(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$  model,  $r_P$ is the MLE, and  $r_P$  and  $\hat{\rho}_{NS}$  are asymptotically optimal in the parametric sense.

A fourth possible estimate is the profile NP estimate obtained by fixing  $\rho$  and replacing  $F_1$  and  $F_2$  by step functions with jumps  $\{p_i\}$  and  $\{q_i\}$  at  $(X_i, Y_i)$  in the log likelihood for the normal copula model. That is, ignoring constants ( $\rho$  is fixed), we maximize

$$l(\boldsymbol{p}, \boldsymbol{q}) = \sum \{\log p_i + \log q_i\}$$

$$+\frac{1}{2}((\Phi^{-1}(\sum_{k:X_k\leq X_i}p_k))^2 + \frac{1}{2}(\Phi^{-1}(\sum_{k:Y_k\leq Y_i}q_k))^2$$

(7.2)

$$+ \frac{1}{2}(1-\rho^2)^{-1}[(\Phi^{-1}(\sum_{k:X_k \le X_i} p_k))^2]$$

$$-2\rho\Phi^{-1}(\sum_{k:X_k\leq X_i}p_k)\Phi^{-1}(\sum_{k:Y_k\leq Y_i}q_k)+(\Phi^{-1}(\sum_{k:Y_k\leq Y_i}q_k))^2]\},$$

w.r.t.  $(\boldsymbol{p}, \boldsymbol{q})$  where  $\sum p_i = 1$  and  $\sum q_i = 1$ . Then given  $(\hat{\boldsymbol{p}}, \hat{\boldsymbol{q}})$ , maximize the log likelihood w.r.t.  $\rho$ , which gives  $\hat{\rho}_{PROF}$ , a profile NPMLE.

**Remark 7.1**: An estimate  $\hat{\theta}$  of a parameter  $\theta \in R$  in a semiparametric model is regular if  $\sqrt{n}(\hat{\theta} - \theta) \rightarrow_d N(0, V_{\hat{\theta}}(\theta, \eta))$  for some asymptotic variance  $V_{\hat{\theta}}(\theta, \eta)$  and if  $\hat{\theta}$  satisfies additional regularity conditions given in Bickel et al. (1993, 1998). For  $F \in \mathcal{F}, V_{r_P}(\rho, F)$  depends on F (e.g. Bickel and Doksum (2007), Example 5.3.6), while  $V_{\hat{\rho}_{NS}}(\rho, F)$  does not, as shown by Klaassen and Wellner (1997). Klaassen and Wellner (1997) go on to argue that  $(1 - \rho^2)^2$  is a semiparametric asymptotic variance lower bound for the class S of all regular estimates of  $\rho$ . Thus  $\hat{\rho}_{NS}$  is semiparametrically optimal in the minimax sense:

(7.3) 
$$\sup\{V_{\hat{\rho}_{NS}}(\rho,F):F\in\mathcal{F}\} = \inf_{\hat{\rho}\in\mathcal{S}}\sup\{V_{\hat{\rho}}(\rho,F):F\in\mathcal{F}\}$$

**Remark 7.2**: Recall that  $\hat{\rho}_{NS}$  was obtained by inserting normal scores  $Z'_i$  and  $W'_i$  in the MOM estimate for the model with  $F_1$  and  $F_2$  known, and that the MLE  $r_{MLE}$  for this model has variance  $(1 - \rho^2)^2/(1 + \rho^2)$ . Klaassen and Wellner (1997) have shown that the approximate MLE  $\hat{\rho}_{MLE}$  obtained from  $r_{MLE}$  by replacing  $(Z_i, W_i)$  with  $(Z'_i, W'_i)$  is semiparametrically optimal in the same sense as  $\hat{\rho}_{NS}$ . Because the distribution of the ranks do not depend on  $F_1$  and  $F_2$ , this implies that  $\hat{\rho}_{NS}$  and  $\hat{\rho}_{MLE}$  are asymptotically equivalent for every  $F \in \mathcal{F}$ . We conjecture that  $\hat{\rho}_{PROF}$  is also asymptotically optimal and equivalent to  $\hat{\rho}_{NS}$ . 

#### 7.2. The Multivariate Covariate Case

The normal copula model in the multivariate case is defined as follows: Let  $Y \sim G$ ,  $X_j \sim F_j$ ,  $\mathbf{h}(\mathbf{X}) = (h_1(X_1), \cdots, h_d(X_d))$ , where  $h_j$ ,  $j = 0, \cdots, d$  are increasing functions defined by

(7.4) 
$$h_0(Y) = \Phi^{-1}(G(Y))$$

(7.5) 
$$h_i(X_i) = \Phi^{-1}(F_i(X_i)).$$

The distribution of the untransformed variables  $(\mathbf{X}, Y)$  is a *copula model* if we assume that  $(\mathbf{h}(\mathbf{X}), h_0(Y))$  is multivariate normal with 0 means and unit variances.

#### 8. Transformation and NP Models

Consider a regression experiment with response Y and a random covariate vector  $\mathbf{X} = (X_1, \dots, X_d)^T$ . We will extend the normal scores estimate  $\hat{\rho}_{NS}$  of Section 7 to the d dimensioned case and compare it with estimates appropriate for parametric and nonparametric models. In the *copula regression* model of Section 7.2, we can write

(8.1) 
$$h_0(Y) = \beta^T \boldsymbol{h}(\boldsymbol{X}) + \epsilon, \ \epsilon \sim N(0, \sigma^2),$$

where  $\beta$  is the set of regression coefficients when regressing  $h_0(Y)$  on h(X). The transform both sides Box-Cox model is based on (8.1) with  $h_0(Y) = Y^{(\lambda_{d+1})}$ ,  $h_j(X_j) = X_j^{(\lambda_j)}, j = 1, \dots, d$ , where  $t^{(\lambda)} = (t^{\lambda} - 1)/\lambda$ . Thus for this case, we can write

(8.2) 
$$Y^{(\lambda_{d+1})} = \alpha + \boldsymbol{\beta}^T \boldsymbol{X}^{(\boldsymbol{\lambda})} + \epsilon, \ \epsilon \sim N(0, \sigma^2).$$

We first consider a procedure for estimating the parameters in model (8.2):

# *I* Profile Likelihood for a multivariate model.

Hernandez and Johnson (1980) considered the one sample multivariate Box-Cox transformation model. This was adopted to regression by Doksum, Ozeki, Kim and Neto (2007). We regard  $(Y^{(\lambda_{d+1})}, X^{(\lambda)})$  as a d+1 multivariate normal  $(\mu, \Sigma)$  vector. Regressing  $Y^{(\lambda_{d+1})}$  on  $X^{(\lambda)}$  leads to (8.2). We fix  $\xi \equiv (\lambda, \lambda_{d+1})$  and estimate the parameters in the normal model by maximizing the likelihood thereby obtaining the familiar normal theory estimates  $(\mu(\xi), \Sigma(\xi))$ . We plug these into the likeli-hood and obtain the profile likelihood  $l(\xi)$ , which we maximize to get  $\hat{\xi}$  and the final estimates  $(\hat{\mu}(\hat{\xi}), \hat{\Sigma}(\hat{\xi}))$ . These are the usual linear model estimates with  $Y_i, X_{ij}$ replaced by  $Y_i^{(\hat{\lambda}_{d+1})}, X_{ij}^{(\hat{\lambda}_j)}$ . Similarly, the estimate of  $\beta$  in (8.2) is the usual linear model estimate with  $Y_i$ ,  $X_{ij}$  replaced by  $Y_i^{(\hat{\lambda}_{d+1})}$ ,  $X_{ij}^{(\hat{\lambda}_j)}$ . 

Doksum and Ozeki

Remark 8.1 : We also considered the maximum likelihood estimates of the parameters  $\beta, \sigma^2$ , and  $(\lambda, \lambda_{d+1})$  in model (8.2). This approach has the problem that if we want to test  $H_0$ :  $\beta_i = 0$ , then  $\lambda_i$  is not identifiable under  $H_0$ . Approach I does not have this problem. This is one case where likelihood and profile likelihood are very different. The algorithm for this MLE often failed to converge. When it did converge, it produced results close to those of method I. We omit the details.

**Remark 8.2**: As pointed out by Zou and Hall (2002), when d=1, the MLE of  $\rho$  in the Box-Cox transformation model with unknown transformation parameters and standardized transformations have the same efficiency as the MLE for the model with known transformation parameters because this Box-Cox model is between the bivariate normal model and the bivariate normal copula model and the MLE's in these models have the same asymptotic variance  $(1-\rho^2)^2$ . The result that the efficiency is the same whether or not the  $\lambda$ 's are known in this Box-Cox model was also obtained by Wong (1981). This result is very different from the results of Bickel and Doksum (1981) regarding the estimation of regression coefficients.

Remark 8.3 : Consider the transformation model

(8.3) 
$$h_0(Y) = \boldsymbol{\beta}^T \boldsymbol{X} + \boldsymbol{\epsilon}, \ \boldsymbol{\epsilon} \sim N(0, \sigma)$$

where **X** is a vector of random covariates and  $h(\cdot)$  is increasing. In this case we can consider the rank estimate  $\hat{\beta}_R$  obtained by maximizing the rank likelihood  $l_r(\beta) = P(\mathbf{R} = \mathbf{r})$  defined in Section 3. The results of Bickel and Ritov (1997) imply that in a certain sense  $\hat{\beta}_R$  is semiparametrically optimal for model (8.3). However the normal scores estimate of  $\hat{\boldsymbol{\beta}} = (\boldsymbol{x}^T \boldsymbol{x})^{-1} \boldsymbol{x}^T \boldsymbol{a}$ , where  $\boldsymbol{a} = (a(S_1), \cdots, a(S_n))^T$ and x is a vector of nonrandom covariates, is not asymptotically optimal unless  $|\beta|/\sigma$  tends to zero at a certain rate as  $n \to \infty$  (Doksum (1987)). MC methods for  $\hat{\boldsymbol{\beta}}_R$  is introduced in Bickel and Doksum (2009), Section 10.5.

We next introduce a semiparametric approach for the copula regression model and a nonparametric regression approach.

#### **II** Normal score substitution.

The model (8.1) with  $h_i$ ,  $j = 0, \dots, d$ , satisfying (7.4) and (7.5) is invariant under increasing transformations. As in the d=1 case, this leads to using the ranks  $\{S_i\}$  of the Y's and the ranks  $\{R_{ij}: i = 1, \cdots, n\}$  of  $X_{ij}$  among  $\{X_{ij}: i = 1, \cdots, n\}$  $1, \dots, n, j = 1, \dots, d$ . Because the distribution of the ranks is invariant under increasing transformations, for rank methods, model (8.1) is equivalent to

$$Y' = \boldsymbol{\alpha}^T \boldsymbol{X}' + \boldsymbol{\epsilon}',$$

where  $X'_j \sim N(0,1)$  and  $\epsilon'$  are independent,  $Y' \sim N(0,1)$  and  $\boldsymbol{\alpha}^T = \Sigma^{-1} \boldsymbol{\rho}$  with  $\boldsymbol{\rho} = (Corr(X'_1, Y'), \cdots, Corr(X'_d, Y'))^T$  and  $\Sigma$  the correlation matrix of  $\boldsymbol{X}' =$  $(X'_i)_{d \times 1}$ . Here  $\Sigma$  is assumed to be nonsingular. Based on the distribution of the ranks,  $\alpha_1, \dots, \alpha_d$  are identifiable parameters in model (8.4). These parameters represents the relative importance of the  $X_j$ 's.

The normal scores  $Z'_{ij} = a(R_{ij})$  and  $W'_i = a(S_i)$  have approximately the same distribution as the unobservable  $X'_{ij}$  and  $Y'_i$  in model (8.4). Because  $E(Y'|\mathbf{x}') = T_i I_i$  is the unobservable  $X'_{ij}$  and  $Y'_i$  in model (8.4).  $\boldsymbol{\alpha}^T \boldsymbol{x}'$ , if we replace  $X'_{ij}$  and  $Y'_i$  with  $Z'_{ij}$  and  $W'_i$ , we find that an approximate 



where  $Z_D$  is the no intercept design matrix  $(Z'_{ij})_{n \times d'}$ , d' is the rank of the matrix  $(Z'_{ij})$ , and  $\mathbf{W'} = (W'_1, \dots, W'_n)^T$ . Any subset of variables  $X_j : j \in J$  with the same ranks, say  $R_{1J}, \dots, R_{nJ}$ , is collapsed into one variable denoted as  $X_J$  with ranks  $R_{1J}, \dots, R_{nJ}$  to avoid singularity. Based on Klaassen and Wellner (1997), we conjecture that  $\hat{\boldsymbol{\alpha}}$  is semiparametrically efficient for the multivariate normal copula model.

#### III Nonparametric estimation.

We next introduce a nonparametric approach. We consider the model

(8.6) 
$$Y = m(\boldsymbol{X}) + \epsilon,$$

where m() is unknown and  $\epsilon$  has median zero. To estimate m(), we use a cubic Bspline and the R function smooth.spline(). The number of knots are automatically selected (less than the number of observations n.) The smoothing parameter is chosen by generalized cross validation (GCV.)

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

# 8.1. Simulation Results

We consider the d=1 case and consider the properties of estimates of  $\rho = Corr(h_1(X), h_0(Y))$ . In this case, the method II estimate is  $\hat{\rho}_{NS}$ .

#### 8.1.1. Correctly Specified Model

The true model satisfies

(8.7) 
$$Y^{(\lambda_2)} = \alpha_0 + \alpha_1 X^{(\lambda_1)} + \epsilon,$$

where  $\epsilon \sim N(0, \sigma^2)$  and  $X^{(\lambda_1)} \sim N(\mu_1, \sigma_0^2)$  are independent. This model is a subset of the normal copula model  $\mathcal{F}$  with

(8.8) 
$$F_1(x) = \Phi(\frac{x^{(\lambda_1)} - \mu_1}{\sigma_0}), \ F_2(y) = \Phi(\frac{y^{(\lambda_2)} - \mu_2}{\sigma_2}),$$

where  $\mu_2 = EY^{(\lambda_2)}$ , and  $\sigma_2^2 = VarY^{(\lambda_2)}$ . We use 1000 MC trials and take  $\sigma^2 = 1$ ,  $\sigma_0^2 = 1$ ,  $(\lambda_1, \lambda_2) \in \{(0.5, 0.5), (1, 1)\}$ ,  $\alpha_0 = 6$ ,  $\alpha_1 = \in \{0, 0.1, 0.5, 1, 2\}$ , and  $\mu_1 = 5$ . Fig. 3 shows that methods I and II have similar properties for  $\alpha_1 \leq 0.5$ . For larger  $\alpha_1$ , the normal scores estimate has a downward bias which is negligible for  $n \geq 500$  (not shown here.) Method I converges all the time with the constraint  $-4 \leq \lambda \leq 4$ . Method II does not involve any optimization and hence converges all the time.

We simulate the data from

(8.9) 
$$Y^{(\lambda_2)} = (1 - \gamma)(\alpha_0 + \alpha_1 X^{(\lambda_1)}) + \gamma[L(X)] + \epsilon,$$
<sup>29</sup>
<sub>30</sub>

where L() is a nonlinear function. Thus the model is a Box-Cox model when  $\gamma = 0$ , but when  $\gamma > 0$ , we are checking the performance of the methods when the model generating the methods are misspecified.

For comparisons of methods we need a parameter that makes sense for all three methods. One such parameter is

(8.10) 
$$m(x) = Median(Y|X = x).$$

We consider the 25th, 50th, and 75th population quantiles of X, i.e., our parameters of interest are  $m(x_{0.25})$ ,  $m(x_{0.50})$ , and  $m(x_{0.75})$ .

Methods I and II are based on models of the form

42 (8.11) 
$$h_0(Y) = g(\boldsymbol{X}, \boldsymbol{\beta}) + \boldsymbol{\epsilon},$$

where  $h_0(\cdot)$  is an increasing function. If X and  $\epsilon$  are independent and  $median(\epsilon) = 0$ , then

$$\begin{array}{c} {}^{46}_{47} \\ {}^{47}_{47} \end{array} (8.12) \\ m(x) = h_0^{-1}(g({\bm X},{\bm \beta})). \\ \end{array}$$

<sup>48</sup> For method I, the MLE of m(x) in model (8.7) is,

(8.13) 
$$\hat{m}(x) = (\hat{\lambda}_2 [\hat{\beta}_0 + \hat{\beta}_1 \frac{x^{\hat{\lambda}_1} - 1}{\hat{\lambda}_1}] + 1)^{1/\hat{\lambda}_2}.$$

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

З

 $h_0(Y) = \rho h_1(X) + \epsilon, \ \epsilon \sim N(0, \sigma^2),$ (8.14)where  $\rho \equiv \rho(h_1(X), h_0(Y))$  is the correlation coefficient. Then by (8.12),  $m(x) = h_0^{-1}(\rho h_1(x)),$ (8.15)where  $h_0^{-1}(t) = F_2^{-1}(\Phi(t)).$ (8.16)It follows that  $m(x) = F_2^{-1}(\Phi(\rho \Phi^{-1}(F_1(x)))),$ (8.17)and by replacing  $F_1$  and  $F_2$  by their empiricals, a natural estimate of m(x) is (8.18) $\hat{m}(x) = y_{\left(\left[n\Phi(\hat{q}(x))\right]\right)},$ where  $\hat{g}(x) = \hat{\rho}_{NS} \Phi^{-1}(\hat{F}_1(x))$  and [] is the greatest integer function. For method III, we use the smoothing spline estimate of E(Y|X = x) described earlier. In our models with normal errors, E(Y|X=x) coincide with the conditional median m(x).

In the simulation we use model (8.9) with  $\epsilon \sim N(0, \sigma^2)$  and  $X^{(\lambda_1)} \sim N(\mu_1, \sigma_0^2)$ independent,

(8.19) 
$$L(t) = \alpha_0 + \alpha_1 \mu_1 - 1.25 + 2.5[1 + \exp(-10(t - \mu_1))]^{-1},$$

 $\sigma^2 \in \{0.01, 0.1, 0.5, 1\}, \ \sigma_0^2 = 1, \ (\lambda_1, \ \lambda_2) \in \{(0.5, 0.5), (1, 1)\}, \ \alpha_0 = 6.25, \ \alpha_1 = \in \{0, 0.5, 1, 2\}, \ \mu_1 = 5, \ \text{and} \ \gamma \in \{0, 0.25, 0.5, 0.75, 1\}.$  The sample size is n=512. There are 1000 MC trials.

Fig. 4, 5, 6, and 7 are boxplots of  $\hat{m}(x_{0.25}), \hat{m}(x_{0.50}), \hat{m}(x_{0.75})$  with the setting  $(\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (1, 1, 0.5, 0.5), (\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (1, 1, 0.5, 1), (\lambda_1, \lambda_2, \alpha_1, \sigma^2)$ =(0.5, 0.5, 0.5, 0.5), and  $(\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (0.5, 0.5, 0.5, 1)$  respectively. Fig.8-11 give MSE's for the estimates  $\hat{m}(x_{0.25})$ ,  $\hat{m}(x_{0.50})$ , and  $\hat{m}(x_{0.75})$ .

We see that method I is the best when model (8.2) is correct, that is,  $\gamma = 0$ . However when the model is increasingly misspecified, i.e., as  $\gamma$  increases, its absolute bias increases which leads to low MSE performance.

Method II is overall best in terms of MSE when  $\lambda_1 = \lambda_2 = 0.5$  (Fig 6 and 7). When  $\gamma = 0$ , it is unbiased and its variance is between Method I and Method III (Fig 4, 5, 6, and 7).

Method III is overall best in terms of MSE when  $\lambda_1 = \lambda_2 = 1$  and the model is badly misspecified. It's smaller bias makes up for its large variance in this case. But its MSE suffers at and near model (8.2) (Fig 6, 7,  $\gamma = 0$ ).

In summary, the normal score procedure performs very well at and close to a copula model. For n large, this is to be expected from the results of Klaassen and Wellner (1997). The normal scores estimate is competitive with the Box-Cox estimate in the transform both sides Box-Cox model.

З

For method II, write model (8.1) as


FIG 4. Boxplots of the three estimates of median regression m(x) for model (8.9) with  $(\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (1, 1, 0.5, 0.5)$ . I: profile MLE, II: normal scores, and III: NP, spline. The true value of m(x) is the solid line.



imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009







imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

Doksum and Ozeki



imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009



FIG 7. Boxplots of the three estimates of median regression m(x) for model (8.9) with  $(\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (0.5, 0.5, 0.5, 1)$ . I: profile MLE, II: normal scores, and III: NP, spline. The true value of m(x) is the solid line.



imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

З



FIG 8. MSE of the three estimates of m(x) as a function of the misspecification parameter  $\gamma$  for  $(\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (1, 1, 0.5, 0.5)$ .  $\bigcirc = I$ ,  $\triangle = II$ , + = III.



FIG 9. MSE of the three estimates of m(x) as a function of the misspecification parameter  $\gamma$  for  $(\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (1, 1, 0.5, 1)$ .  $\bigcirc = I, \triangle = II, + = III$ .





imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009



FIG 11. MSE of the three estimates of m(x) as a function of the misspecification parameter  $\gamma$  for  $(\lambda_1, \lambda_2, \alpha_1, \sigma^2) = (0.5, 0.5, 0.5, 1)$ .  $\bigcirc = I, \ \triangle = II, \ + = III$ .

#### Acknowledgements

We thank Peter Bickel, Aad van der Vaart and Jon Wellner for helpful comments.

## References

З

[1] ANDERSEN, P. K., BORGAN, O., GILL, R. D. and KEIDING N. (1996). Statistical Models Based on Counting Processes. Springer-Verlag, New York. [2]BEGUN, J. M., HALL, W. J., HUANG, WEI-MIN and WELLNER, J. A. (1983). Information and Asymptotic Efficiency in Parametric-Nonparametric Models. The Annals of Statistics, 11, 2, 432–452. BELL, C. B. and DOKSUM, K. A. (1966). "Optimal" One-Sample Distribution-Free Tests and Their [3] Two-Sample Extensions. The Annals of Mathematical Statistics, 36, 1, 120-132. BHUCHONGKUL, S. (1964). A Class of Nonparametric Tests for Independence in Bivariate Populations. [4] Ann. Math. Statist., 35,1, 138-149. [5] BICKEL, P. J. and DOKSUM, K. A. (1981). An analysis of transformations revisited. Journal of the American Statistical Association, 76, 296–311. [6] BICKEL, P. J. and DOKSUM, K. A. (2007). Mathematical Statistics: Basic Ideas and Selected Topics, 2nd ed. Vol I, Updated Printing. Pearson Prentice Hall, Upper Saddle River, NJ. BICKEL, P. J. and DOKSUM, K. A. To appear (2009). Mathematical Statistics: Basic Ideas and Selected Topics Vol II. Pearson Prentice Hall, Upper Saddle River, NJ. BICKEL, P. J., KLAASSEN, C. A., RITOV, Y. and WELLNER, J. A. (1993, 1998). Efficient and Adaptive Estimation for Semiparametric Models. Springer-Verlag, New York. [9] BICKEL, P. J. and RITOV, Y. (1997). Local asymptotic normality of ranks and covariates in trans-formation models. In Festschrift for Lucien Le Cam (D. Pollard and G. L. Yang, eds.) Springer, New York. [10] Cox, D. R. (1964). Some applications of exponential ordered scores. J. R. Stat. Soc., B.26, 103–110. Cox, D. R. (1972). Regression models and life tables (with discussion). J. Roy. Statist. Soc., B.34, 187 - 220.[12] Cox, D. R. (1975). Partial likelihood. Biometrika., 62.2, 269-276. [13] Cox, D. R. (2006). Principles of Statistical Inference. Cambridge University Press, Cambridge, [14] DOKSUM, K. A. (1987). An extension of partial likelihood methods from proportional hazard models to general transformation models. Annals of Statistics, 15, 325–345. [15] DOKSUM, K. A., OZEKI, A., KIM, J. and NETO, E. C. (2007). Thinking outside the box: Statistical inference based on Kullback-Leibler Empirical Projections. Statistics and Probability Letters. 77, 1201 - 1213.[16] FISHER, R. A. and YATES, F. (1938). Statistical Tables for Biological, Agricultural and Medical Research, (5th edition, 1957) Oliver and Boyd, Edinburgh. [17] FERGUSON, T. S. (1967). Mathematical Statistics. A Decision Theoretic Approach. New York and London, Academic Press. [18] HAJEK, J. and SIDAK, Z. (1967). Theory of Rank Tests. Academic Press, New York. 

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

#### $Doksum \ and \ Ozeki$

| 1              | [19]  | HODGES, J. L. and LEHMANN, E. L. (1963). Estimates of location based on rank tests. Ann. Math.<br>Stat. 34, 508 611   | 1  |
|----------------|-------|---|----|
| 2              | [20]  | HOEFFDING, W. (1951). 'Optimum' nonparametric tests. Proc. 2nd Berkeley Symposium on Math-  | 2  |
| 3              | L - J | ematical Statistics and Probability, Berkeley, Univ. Calif. Press, 83–92.   | 3  |
| 4              | [21]  | HERNANDEZ, F. and JOHNSON, R. A. (1980). The large-sample behavior of transformations to nor-   | 4  |
| 5              | [22]  | mality. Journal of the American Statistical Association, <b>75</b> , 855–861.   | 5  |
| 6              | [22]  | life model. Biometrika, 60, 267–278.  | 6  |
| 7              | [23]  | KALBFLEICH, J. D. and PRENTICE, R. L. (2002). The Statistical Analysis of failure Time Data, 2nd  | 7  |
| 8              |       | edition. John Wiley and Sons, Inc., Hoboken, New Jersey.  | 8  |
| 9              | [24]  | KLAASSEN, C. A. J. (2007). A Sturm-Liouville Problem in Semiparametric Transformation Models.   | 9  |
| 10             |       | Vijav Nair. World Scientific Pub Co Inc., New Jersev.   | 10 |
| 11             | [25]  | KOSOROK, M. R., LEE B. L. and FINE J. P. (2004). Robust inference for univariate proportional   | 11 |
| 12             | [ ]   | hazards frailty regression models. The Annals of Statistics, <b>32</b> , No. 4, 1448–1491.  | 12 |
| 13             | [26]  | KLAASSEN, C. A. J. and WELLNER, J. A. (1997). Efficient Estimation in the Bivariate Normal Copula<br>Model: Normal Margine Are Least Envourable. <i>Bernoulli</i> <b>3</b> , No. 1, 55–77 | 13 |
| 14             | [27]  | LANGE, K., HUNTER, D. R. and YANG, I. (2000). Optimization Transfer Using Surrogate Objective   | 14 |
| 15             |       | Functions. Journal of Computational and Graphical Statistics, 9, No. 1, 1–20.   | 14 |
| 15             | [28]  | LEHMANN, E. L. (1953). The power of rank tests. Ann. Math. Statist., 24, 23–43.   | 15 |
| 16             | [29]  | MURPHY, S. A. and VAN DER VAART, A. W. (2000). On Profile Likelihood. JASA, 9, No. 450, 449–465.  | 16 |
| 17             | [30]  | natives. Technical Report. University of California. Berkeley.  | 17 |
| 18             | [31]  | OAKES, D. (1992). Bivariate survival models induced by frailties. JASA, 84, 487–493.  | 18 |
| 19             | [32]  | OAKES, D. and JONG-HYEON JEONG (1998). Frailty Models and Rank Tests. Lifetime Data Analysis,   | 19 |
| 20             | [99]  | 4, 3, 209–228.  | 20 |
| 21             | ႞ၜႄၜ႞ | Biometrika, 75. No.2. 237–249.  | 21 |
| 22             | [34]  | OWEN, A. B. (2001). Empirical Likelihood. Chapman and Hall.   | 22 |
| 23             | [35]  | RUYMGAART, F. H. (1974). Asymptotic Normality of Nonparametric Tests for Independence. Ann.   | 23 |
| 24             | [26]  | Statist., 2, 892–910.   | 24 |
| 25             | [30]  | parametric Tests for Independence. Ann. Math. Statist., 43, 1122–1135.  | 25 |
| 26             | [37]  | SAVAGE, I. R. (1956). Contributions to the theory of rank orders statistics: the two-sample case.   | 26 |
| <br>27         |       | Ann. Math. Stat., 27, 590–615.  | 27 |
| 28             | [38]  | SAVAGE, I. R. (1957). Contributions to the Theory of Rank Order Statistics-The "Trend" Case.  | 28 |
| 20             | [39]  | SIBUYA, M. (1968). Generating Doubly Exponential Random Numbers. Annals of the Institute of   | 20 |
| 29             | []    | Statistical Mathematics, Tokyo, Supplement V, 1–7.  | 29 |
| 30             | [40]  | SKLAR, A. (1959). Fonctions de repartition a n dimensions et leurs marges. L'Institut de Statistique  | 30 |
| 31             | [41]  | de L'Universite de Paris. 8, 229–231.<br>TRODUCY A (2002) Sominamentais modeles a generalized celf consistence approach $L$ $R$ . Statist   | 31 |
| 32             | [41]  | Soc., B.65, 3, 759–774.   | 32 |
| 33             | [42]  | TSODIKOV, A. and GARIBOTTI, G. (2007). Profile information matrix for nonlinear transformation  | 33 |
| 34             | [ ]   | models. Lifetime Data Analysis, <b>13</b> , 139–159.  | 34 |
| 35             | [43]  | VAN DER VAART, A. W. (1998). Asymptotic statistics. Cambridge University Press, Cambridge, UK.  | 35 |
| 36             | [44]  | with censored data. J. R. Statist., 69, 4, 1–30.  | 36 |
| 37             | [45]  | ZOU, K. H. and HALL, W. J. (2002). On estimating a transformation correlation coefficient. Journal  | 37 |
| 38             | [10]  | of Applied Statistics, <b>29</b> , 745–760.   | 38 |
| 39             | [46]  | WONG, U. W. (1981). Transformation of Independent Variables in Regression Models. <i>Ph.D. Thesis</i> .   | 39 |
| 40             |       | Department of Statistics, University of Derkeley, CA.   | 40 |
| 41             |       |   | 41 |
| 42             |       |   | 42 |
| 43             |       |   | 43 |
| 10             |       |   | 10 |
| - <del>-</del> |       |   | 44 |
| 40             |       |   | 40 |
| 40             |       |   | 46 |
| 47             |       |   | 47 |
| 48             |       |   | 48 |
| 49             |       |   | 49 |
| 50             |       |   | 50 |
| 51             |       |   | 51 |

imsart-coll ver. 2008/08/29 file: Doksum.tex date: March 25, 2009

| University of Wisconsin–Madison<br>Abstract: The size of the bootstrap test of hypotheses is studied for t<br>normal and exponential one and two-sample problems. It is found that t<br>size depends not only on the problem, but on the choice of test statistic as<br>the nominal level. In some special cases, the bootstrap test is UMP, but<br>other cases, it can be totally useless, such as being completely randomized<br>rejecting the null hypothesis with probability one. More importantly, the si<br>is usually greater than the nominal level, even in the limit as the sample si<br>goes to infinity. | he<br>hd<br>in<br>or<br>ize<br>ize       |                |
|---|--|----------------|
| Abstract: The size of the bootstrap test of hypotheses is studied for t<br>normal and exponential one and two-sample problems. It is found that t<br>size depends not only on the problem, but on the choice of test statistic and<br>the nominal level. In some special cases, the bootstrap test is UMP, but<br>other cases, it can be totally useless, such as being completely randomized<br>rejecting the null hypothesis with probability one. More importantly, the si<br>is usually greater than the nominal level, even in the limit as the sample si<br>goes to infinity.                                   | he<br>he<br>nd<br>in<br>or<br>ize<br>ize |                |
| Contents  |  |                |
|   |  |                |
| 1 Introduction  |  | 93             |
| 2 Testing a Normal Mean   |  | 94             |
| 2.1 Known Variance  |  | 94             |
| 2.1.1 Sample Mean Statistic   |  | 94             |
| 2.1.2 Standard Likelihood Ratio Statistic   |  | 94             |
| 2.1.3 Cox Likelihood Ratio Statistic  |  | 95             |
| 2.2 Unknown Variance  |  | 96             |
| 2.2.1 Standard Likelihood Ratio Statistic   |  | 9'             |
| 2.2.2 Cox Likelihood Ratio Statistic  |  | 99             |
| 3 Testing a Normal Variance, Mean Unknown   |  | 100            |
| 3.1 $H_0: \sigma^2 \le 1$ vs. $H_1: \sigma^2 > 1$   |  | 100            |
| 3.1.1 Standard Likelihood Ratio Statistic   |  | 100            |
| 3.1.2 Cox Likelihood Ratio Statistic $\dots \dots \dots \dots$  |  | 102            |
| 3.2 $H_0: \sigma^2 \ge 1$ vs. $H_1: \sigma^2 < 1$   |  | 103            |
| 3.2.1 Standard Likelinood Ratio Statistic   |  | 100            |
| 3.2.2 Cox Likelihood Ratio Statistic  |  | 104            |
| 4 Testing Difference of Two Normal Means  |  | 100            |
| 4.1.1 Difference of Means Statistic   |  | 10             |
| 4.1.2 Standard Likelihood Ratio Statistic   |  | 106            |
| 4.1.3 Cox Likelihood Ratio Statistic  |  | 106            |
| 4.2 Unknown but Equal Variances   |  | 107            |
| 4.2.1 Difference of Means Statistic   |  | 107            |
| 4.2.2 Standard Likelihood Ratio Statistic   |  | 109            |
| Department of Statistics University of Wissensin Medican 1200 University  | Auonua                                   | 117            |
| 53706, email: loh@stat.wisc.edu   | Avenue                                   | , vv           |
| $^{2}$ Department of Statistics, University of Wisconsin–Madison, 1300 University   | Avenue                                   | , W            |
| 53706, email: zheng@stat.wisc.edu   | E  | ~ <b>4</b> :   |
| ins material is based upon work partially supported by the National Science<br>inder grant DMS-0402470 and the U.S. Army Research Laboratory and the U.S. Arr   | rounda<br>ny Rese                        | atioi<br>earcl |
| Office under grant W911NF-05-1-0047.  |  |                |
| AMS 2000 subject classifications: Primary 62F03, 62F40; secondary 62F05   |  |                |
| distribution, resampling, size, significance level, uniformly most powerful   | expone                                   | snua           |

1

| 4.2.3 Cox Likelihood Ratio Statistic          | ) |
|---|---|
| 5 Testing an Exponential Location Parameter   | ) |
| 5.1 $H_0: \theta \le 0$ vs. $H_1: \theta > 0$ | ) |
| 5.1.1 Standard Likelihood Ratio Statistic     | _ |
| 5.1.2 Cox Likelihood Ratio Statistic          | 2 |
| 5.2 $H_0: \theta \ge 0$ vs. $H_1: \theta < 0$ | 2 |
| 5.2.1 Standard Likelihood Ratio Statistic     | 2 |
| 5.2.2 Cox Likelihood Ratio Statistic          | 2 |
| 6 Conclusion                                  | 3 |
| Appendix                                      | 6 |
| References                                    | ) |

#### 1. Introduction

Owing to its practical convenience and wide applicability, the bootstrap method [7] is used to test statistical hypotheses in many research studies. A sample of recent applications includes evolutionary molecular biology [1], genetic structure [2], gene frequency [11], cancer epidemiology [8], microscopy [3], quality of life [12], economic cycles [5], livestock management [9], and meat demand [6]. Despite its popularity, however, there have been few detailed studies of the theoretical validity of the bootstrap for hypothesis testing. This article addresses this issue for some simple parametric problems where the bootstrap null distributions can be studied analytically. Specifically, we consider one and two-sample problems involving normally and exponentially distributed observations. Our goal is to determine the finite-sample or limiting sizes of the bootstrap tests and compare them with those of the traditional tests.

First, we recall some definitions. Let  $\mathbf{X}_n = (X_1, X_2, \dots, X_n)$  be a vector of n independent observations from  $F_{\mu}$ . In the bootstrap method, we first find an estimate  $\hat{\mu}_0$  of  $\mu$  under  $H_0$  and estimate  $F_{\mu}$  with  $\hat{F} = F_{\hat{\mu}_0}$ . Given a test statistic  $S = S(\mathbf{X}_n)$  for which large values lead to rejection of  $H_0$ , let  $G_{\mu}$  denote the distribution function of S. Let  $\mathbf{X}_n^* = (X_1^*, X_2^*, \dots, X_n^*)$  be a vector of n independent observations from  $\hat{F}$  and define  $S^* = S(\mathbf{X}_n^*)$ . The distribution function  $\hat{G} = G_{\hat{\mu}_0}$  of  $S^*$  is the bootstrap distribution function of S, i.e.,  $\hat{G}$  is the distribution of S under  $\hat{F}$ .

For any nominal level of significance  $\alpha$  ( $0 < \alpha < 1$ ), let  $c_{\alpha}(\hat{\mu}_0)$  be the upper- $\alpha$  quantile of  $\hat{G}$ . Thus  $c_{\alpha}(\hat{\mu}_0)$  is the smallest value such that  $\hat{G}(c_{\alpha}(\hat{\mu}_0)) \ge 1 - \alpha$ . The nominal level- $\alpha$  bootstrap test rejects  $H_0$  with probability 1 if  $S > c_{\alpha}(\hat{\mu}_0)$ , and with probability  $[\alpha - 1 + \hat{G}(c_{\alpha}(\hat{\mu}_0))]/[\hat{G}(c_{\alpha}(\hat{\mu}_0)) - \hat{G}(c_{\alpha}(\hat{\mu}_0)-)]$  if  $S = c_{\alpha}(\hat{\mu}_0)$  and  $\hat{G}(c_{\alpha}(\hat{\mu}_0)) > \hat{G}(c_{\alpha}(\hat{\mu}_0)-)$ .

## 2. Testing a Normal Mean

Let  $X_1, X_2, \ldots, X_n$  be a random sample from  $N(\mu, \sigma^2)$ , a normal distribution with mean  $\mu$  and variance  $\sigma^2$ . Let  $\phi(x)$  and  $\Phi(x)$  denote the density and distribution functions of the N(0, 1) distribution and let  $z_{\alpha}$  be its upper- $\alpha$  critical value, that is,  $1 - \Phi(z_{\alpha}) = \alpha$ . Consider testing

$$H_0: \mu \le 0$$
 vs.  $H_1: \mu > 0.$ 

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

Bootstrap tests

## 2.1. Known Variance

We assume without loss of generality that  $\sigma^2 = 1$ . Let  $\bar{X}_n = n^{-1} \sum_{i=1}^n X_i$ . The unrestricted MLE of  $\mu$  is  $\hat{\mu} = \bar{X}_n$ . Let  $\hat{\mu}_i$  be the MLE of  $\mu$  under  $H_i$  (i = 0, 1). Then  $\hat{\mu}_0 = \bar{X}_n I(\bar{X}_n < 0)$ ,  $\hat{\mu}_1 = \bar{X}_n I(\bar{X}_n > 0)$ , and  $X_1^*, X_2^*, \ldots, X_n^*$  is a bootstrap random sample drawn from  $N(\hat{\mu}_0, 1)$ . Let  $\bar{X}_n^* = n^{-1} \sum_{i=1}^n X_i^*$ .

2.1.1. Sample Mean Statistic

З

**Theorem 2.1.** If  $0 < \alpha \leq 1/2$ , the bootstrap test based on  $\overline{X}_n$  is uniformly most powerful (UMP), but if  $1/2 < \alpha < 1$ , the test rejects  $H_0$  with probability 1.

Proof. Recall that the UMP test rejects  $H_0$  if  $\bar{X}_n \geq z_{\alpha} n^{-1/2}$ . Since  $\bar{X}_n^*$  is normal with mean  $\hat{\mu}_0$  and variance  $n^{-1}$ , its critical value is  $c_{\alpha}(\hat{\mu}_0) = \hat{\mu}_0 + z_{\alpha} n^{-1/2} = \bar{X}_n I(\bar{X}_n < 0) + z_{\alpha} n^{-1/2}$ . Therefore the bootstrap test rejects  $H_0$  if  $\bar{X}_n I(\bar{X}_n > 0) \geq z_{\alpha} n^{-1/2}$ . If  $0 < \alpha \leq 1/2$ , then  $z_{\alpha} \geq 0$  and the bootstrap test is the UMP test. If  $1/2 < \alpha < 1$ , then  $z_{\alpha} < 0$  and the test rejects  $H_0$  w.p.1.

## 2.1.2. Standard Likelihood Ratio Statistic

Given the data and values  $\mu_0$  and  $\mu_1$ , let

$$L(\mu_0, \mu_1, \mathbf{X}_n) = \log \left\{ \prod_{i=1}^n \phi(x_i - \mu_1) \middle/ \prod_{i=1}^n \phi(x_i - \mu_0) \right\}.$$

A general statistic for testing  $H_0$  is the log-likelihood ratio

$$L(\hat{\mu}_{0}, \hat{\mu}, \mathbf{X}_{n}) = \log \left\{ \sup_{\mu} \prod_{i=1}^{n} \phi(x_{i} - \mu) \middle/ \sup_{\mu \in H_{0}} \prod_{i=1}^{n} \phi(x_{i} - \mu) \right\}.$$

Throughout this article, we let Z denote the standard normal variable and  $z_{\alpha}^{+} = \max(z_{\alpha}, 0)$ . We need the following lemma whose proof is given in the Appendix.

**Lemma 2.1.** Let  $\theta \ge 0$ . For fixed  $0 < \alpha < 1$ , the function

(2.2) 
$$P(|Z+\theta| > z_{\alpha}^{+}) - (1-\alpha)E\{\Phi(Z+\theta)^{-1}I(Z+\theta > z_{\alpha}^{+})\}$$

is maximized at  $\theta = 0$  with maximum value

(2.3) 
$$\min(2\alpha, 1) + (1 - \alpha) \log\{1 - \min(\alpha, 1/2)\}\$$

which is greater than  $\alpha$  for all  $0 < \alpha < 1$ .

**Theorem 2.2.** The size of the bootstrap test based on the standard likelihood ratio is  $\min(2\alpha, 1) + (1 - \alpha) \log\{1 - \min(\alpha, 1/2)\}$ .

*Proof.* Since

$$n^{-1}L(\hat{\mu}_0, \hat{\mu}, \mathbf{X}_n) = \bar{X}_n(\hat{\mu} - \hat{\mu}_0) - (\hat{\mu}^2 - \hat{\mu}_0^2)/2$$

$$= X_n(X_n - \hat{\mu}_0) - (X_n^2 - \hat{\mu}_0^2)/2$$

$$= (\bar{X}_n - \hat{\mu}_0)^2/2$$

$$= \bar{X}_n^2 I(\bar{X}_n > 0)/2$$

the test rejects  $H_0$  if  $S = \bar{X}_n I(\bar{X}_n > 0) \ge c_\alpha(\hat{\mu}_0)$ , where the critical value is to be determined. Let  $S^* = \bar{X}_n^* I(\bar{X}_n^* > 0)$  and consider two cases.

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

r f

| 1      | 1. $\bar{X}_n > 0$ . Then $S > 0$ , $\hat{\mu}_0 = 0$ , and $\bar{X}_n^*$ has a $N(0, n^{-1})$ distribution. For any   | 1      |
|--------|--|--------|
| 2      | $x \ge 0, P(S \le x) = P(X_n \le x) = \Phi(xn + 1).$ Therefore if $0 < \alpha < 1/2$ ,   | 2      |
| 3      | $c_{\alpha}(\mu_0) = z_{\alpha}n^{-1}$ . Otherwise, if $\alpha \geq 1/2$ , then $c_{\alpha}(\mu_0) = 0$ and the bootstrap  | 3<br>1 |
| т<br>5 | $2  \overline{X} \leq 0$ Then $S = 0$ $\hat{\mu}_{0} = \overline{X} \leq 0$ and $S^{*}$ has a $N(\overline{X} = n^{-1})$ distribution  | 5      |
| 6      | 2. $X_n \leq 0$ . Then $S = 0$ , $\mu_0 = X_n \leq 0$ , and $S$ has a $N(X_n, n)$ distribution<br>left truncated at 0 with $P(S^* = 0) = P(\bar{X}^* \leq 0) = \Phi(-n^{1/2}\bar{X})$ . Thus             | 6      |
| 7      | Here truncated at 0 with $\Gamma(D = 0) = \Gamma(X_n \le 0) = \Psi(-n + X_n)$ . Thus   | 7      |
| 0      | $(\bar{\mathbf{v}} + m^{-1}/2r)$ ; $f \bar{\mathbf{v}} + m^{-1}/2r > 0$  | ,<br>0 |
| 0      | $c_{\alpha}(\hat{\mu}_{0}) = \begin{cases} \Lambda_{n} + n^{-\gamma-z_{\alpha}}, & \Pi_{n} \Lambda_{n} + n^{-\gamma-z_{\alpha}} > 0 \\ 0 & \vdots & \tilde{\nabla}_{n} + \frac{-1/2}{2} < 0 \end{cases}$ | 0      |
| 10     | $(0, \qquad \text{If } X_n + n^{-1/2} z_\alpha \le 0.$   | 10     |
| 10     | $= (\bar{X}_n + n^{-1/2} z_\alpha)^+.$   | 10     |
| 12     |  | 12     |
| 13     | Since $S = 0$ , the bootstrap test never rejects $H_0$ if $\bar{X}_n + n^{-1/2} z_\alpha > 0$ .  | 13     |
| 14     | Otherwise, the test is randomized and rejects $H_0$ with probability $\{\alpha - 1 +$  | 14     |
| 15     | $\Phi(-n^{1/2}ar{X}_n)\}/\Phi(-n^{1/2}ar{X}_n).$   | 15     |
| 16     |  | 16     |
| 17     | Thus for $0 < \alpha < 1$ ,  | 17     |
| 18     |  | 18     |
| 19     | $P\{\text{Reject } H_0\} = P\{\text{Reject } H_0, X_n > 0\} + P\{\text{Reject } H_0, X_n < 0\}$  | 19     |
| 20     | $= P(S > z_{\alpha}^+, \bar{X}_n > 0)$   | 20     |
| 21     | $+ P\{\text{Reject } H_0, \bar{X}_n + n^{-1/2} z_\alpha \leq 0, \bar{X}_n < 0\}$   | 21     |
| 22     | $= P(\bar{X}) > z^{+}n^{-1/2}$   | 22     |
| 23     | $= T \left( \frac{1}{2} \overline{\mathbf{x}} \right) $  | 23     |
| 24     | $+ E[\{\alpha - 1 + \Psi(-n^{-\gamma} X_n)\}/\Psi(-n^{-\gamma} X_n)]I(-n^{-\gamma} X_n \ge z_{\alpha})$  | 24     |
| 25     | $= P( W  > z_{\alpha}^{+}) - (1 - \alpha)E\{\Phi(W)^{-1}I(W > z_{\alpha}^{+})\}$   | 25     |
| 26     | 1/0  | 26     |
| 27     | where W is normally distributed with mean $-n^{1/2}\mu$ and variance 1. By Lemma 2.1,  | 27     |
| 28     | the supremum of the rejection probability under $H_0$ is attained when $\mu = 0$ and is  | 28     |
| 29     | given by $(2.3)$ . Figure 1 shows a plot of this function.   | 29     |
| 30     |  | 30     |
| 31     |  | 31     |
| 32     | 2.1.3. Cox Likelihood Ratio Statistic  | 32     |
| 33     |  | 33     |
| 34     | Cox (1961) proposed the following alternative likelihood ratio statistic for testing   | 34     |
| 35     | separate families of hypotheses:   | 35     |
| 36     |  | 36     |
| 37     | $\begin{pmatrix} n \\ \mathbf{H} \end{pmatrix}$  | 37     |
| 38     | $L(\hat{\mu}_0, \hat{\mu}_1, \mathbf{X}_n) = \log \left\{ \sup_{H} \prod \phi(x_i - \mu) / \sup_{H} \prod \phi(x_i - \mu) \right\}.$   | 38     |
| 39     | $\begin{pmatrix} II_1 & i=1 & I & II_0 & i=1 \end{pmatrix}$  | 39     |
| 40     | For the current problem  | 40     |
| 41     | For the current problem,   | 41     |
| 42     | $I(\hat{\mu}, \hat{\mu}, \mathbf{X}) = m[\bar{\mathbf{X}}(\hat{\mu}, \hat{\mu}) - (\hat{\mu}^2, \hat{\mu}^2)/2]$   | 42     |
| 43     | $L(\mu_0, \mu_1, \mathbf{A}_n) = n\{\mathbf{A}_n(\mu_1 - \mu_0) - (\mu_1 - \mu_0)/2\}$   | 43     |
| 44     | $= n(X_n X_n  - X_n X_n /2)$   | 44     |
| 45     | $= n\bar{X}_n^2 \operatorname{sgn}(\bar{X}_n)/2.$  | 45     |
| 46     |  | 46     |
| 47     | Therefore rejecting $H_0$ for large values of $L(\hat{\mu}_0, \hat{\mu}_1, \mathbf{X}_n)$ is equivalent to rejecting   | 47     |
| 48     | for large values of $\bar{X}_n$ , and the next theorem follows directly from Theorem 2.1.  | 48     |
| 49     |  | 49     |
| 50     | <b>1 neorem 2.3.</b> If $0 < \alpha \le 1/2$ , the bootstrap test based on the Cox likelihood ratio  | 50     |
| 51     | has size $\alpha$ and is UMP. If $1/2 < \alpha < 1$ , it rejects $H_0$ with probability 1.   | 51     |
|        |  |        |

W. Loh and W. Zheng

96

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009



FIG 1. Size (2.3) of bootstrap test for the normal mean based on the standard likelihood ratio, for known  $\sigma$ . The dashed line is the identity function.

## 2.2. Unknown Variance

Now suppose we test the hypotheses (2.1) without assuming that  $\sigma$  is known. The log-likelihood function is

$$l(\mu, \sigma) = -n \log \sigma - \sum_{i} (X_i - \mu)^2 / (2\sigma^2) - (n/2) \log(2\pi)$$

and its derivatives are  $\partial l/\partial \mu = -\sigma^{-2} \sum (X_i - \mu)$  and  $\partial l/\partial \sigma = -n\sigma^{-1} + \sigma^{-3} \sum (X_i - \mu)^2$ . Hence the unrestricted and restricted (under  $H_0$  and  $H_1$ ) maximum likelihood estimates (MLEs) of  $\mu$  and  $\sigma^2$  are, respectively,

$$\hat{\mu} = \bar{X}_n$$
  
 $\hat{\sigma}^2 = n^{-1} \sum (X_i - \bar{X}_n)^2$ 

$$\hat{\mu}_0 = X_n I(X_n < 0)$$
  $\hat{\sigma}_0^2 = n^{-1} \sum (X_i - \hat{\mu}_0)^2$ 

$$\hat{\mu}_1 = \bar{X}_n I(\bar{X}_n > 0)$$
  $\hat{\sigma}_1^2 = n^{-1} \sum (X_i - \hat{\mu}_1)^2$ 

giving the log-likelihood ratio statistics:

47 **Standard:** 
$$n \log(\hat{\sigma}_0/\hat{\sigma}) = (n/2) \log\{\sum (X_i - \hat{\mu}_0)^2 / \sum (X_i - \bar{X}_n)^2\}$$
  
48 **Cox:**  $n \log(\hat{\sigma}_0/\hat{\sigma}_1) = (n/2) \log\{\sum (X_i - \hat{\mu}_0)^2 / \sum (X_i - \hat{\mu}_1)^2\}.$ 

The corresponding bootstrap tests reject 
$$H_0$$
 for large values of  $\sum (X_i - \hat{\mu}_0)^2 / \sum (X_i - 50 \bar{X}_n)^2$  and  $\sum (X_i - \hat{\mu}_0)^2 / \sum (X_i - \hat{\mu}_1)^2$ , respectively.

З

2.2.1. Standard Likelihood Ratio Statistic

Let

(2.4) 
$$T_n = n^{1/2} \bar{X}_n \left\{ \sum (X_i - \bar{X}_n)^2 / (n-1) \right\}^{-1/2}.$$

The standard log-likelihood ratio statistic is

$$\frac{\sum (X_i - \hat{\mu}_0)^2}{\sum (X_i - \bar{X}_n)^2} = \frac{\sum \{X_i - \bar{X}_n I(\bar{X}_n < 0)\}^2}{\sum (X_i - \bar{X}_n)^2}$$

$$= \begin{cases} 1, & \text{if } \bar{X}_n < 0\\ \sum X_i^2 / \sum (X_i - \bar{X}_n)^2, & \text{if } \bar{X}_n \ge 0 \end{cases}$$

$$= \begin{cases} 1, & \text{if } \bar{X}_n < 0\\ 1 + n \bar{X}^2 / \sum (X_1 - \bar{X}_1)^2 & \text{if } \bar{X}_n > 0 \end{cases}$$

$$\begin{pmatrix} 1 + nA_n / \sum (A_i - A_n) , & \text{if } A_n \ge 0 \\ 1, & \text{if } \bar{X}_n < 0 \end{cases}$$

$$= \begin{cases} 1, & \text{if } X_n < 0 \\ 1 + (n-1)^{-1} T_n^2, & \text{if } \bar{X}_n \ge 0. \end{cases}$$
<sup>16</sup>

Thus  $H_0$  is rejected for large values of  $S = T_n I(T_n > 0)$ . Let  $t_{\nu,\delta}$  denote the noncentral *t*-distribution with  $\nu$  degrees of freedom and noncentrality parameter  $\delta$  and let  $t_{\nu,\delta,\alpha}$  denote its upper- $\alpha$  critical point.

## **Lemma 2.2.** For any $\nu$ and $\alpha$ , $t_{\nu,\delta,\alpha}$ is an increasing function of $\delta$ .

*Proof.* Let Z denote a standard normal variable independent of  $\chi^2_{\nu}$ . Since

$$P(t_{\nu,\delta} \le x) = P\left(\frac{Z+\delta}{\sqrt{\chi_{\nu}^2/\nu}} \le x\right)$$
<sup>25</sup>
<sup>26</sup>
<sup>27</sup>
<sup>26</sup>
<sup>27</sup>

$$= P\left(\frac{Z}{\sqrt{\chi_{\nu}^{2}/\nu}} \le x - \frac{\delta}{\sqrt{\chi_{\nu}^{2}/\nu}}\right)$$

$$= P\left(\frac{Z}{\sqrt{\chi_{\nu}^{2}/\nu}} \le x - \frac{\delta}{\sqrt{\chi_{\nu}^{2}/\nu}}\right)$$

$$= \frac{28}{30}$$

we see that  $P(t_{\nu,\delta} \leq x)$  is a decreasing function of  $\delta$ . Therefore  $t_{\nu,\delta,\alpha}$  is an increasing function of  $\delta$ .

**Theorem 2.4.** If  $\sigma$  is unknown, the size of the nominal level- $\alpha$  test of  $H_0: \mu \leq 0$ vs.  $H_1: \mu > 0$  based on the standard likelihood ratio has lower bound

$$\min(\alpha, 1/2) + E\left\{\frac{\alpha - 1 + \Phi\left(-t_{n-1}\sqrt{n/(n-1)}\right)}{\Phi\left(-t_{n-1}\sqrt{n/(n-1)}\right)} I\left(t_{n-1}\sqrt{\frac{n}{n-1}} < -z_{\alpha}^{+}\right)\right\}$$

where  $t_{n-1}$  has a (central) t-distribution with n-1 degrees of freedom. As  $n \to \infty$ , the bound tends to (2.3), the size for the case where  $\sigma$  is known and n is finite.

# *Proof.* Again, consider two cases.

1.  $\bar{X}_n > 0$ . Then S > 0 and  $\hat{\mu}_0 = 0$ . The bootstrap distribution of  $T_n^*$  is a central  $t_{n-1}$ -distribution and that of  $S^*$  is a central  $t_{n-1}$ -distribution left-truncated at 0. If  $0 < \alpha < 1/2$ , the test rejects  $H_0$  whenever  $T_n > t_{n-1,0,\alpha}$ . Otherwise, if  $\alpha \ge 1/2$ , the test rejects  $H_0$  with probability 1.

2. 
$$\bar{X}_n < 0$$
. Then  $S = 0$ ,  $\hat{\mu}_0 < 0$ , and  $S^*$  has a left-truncated noncentral  $t_{n-1,\delta}$ -distribution with  $n-1$  degrees of freedom and noncentrality parameter

(2.5) 
$$\delta = n^{1/2} \hat{\mu}_0 / \hat{\sigma}_0 = n \bar{X}_n / \sqrt{\sum (X_i - \bar{X}_n)^2} = T_n \sqrt{n/(n-1)}$$

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

З

| 1        | and probability $P(\bar{X}_n^* \le 0) = P\{n^{1/2}(\bar{X}_n^* - \hat{\mu}_0)/\hat{\sigma}_0 \le -n^{1/2}\hat{\mu}_0/\hat{\sigma}_0\} = \Phi(-\delta)$  | 1  |
|----------|---|----|
| 2        | at 0.   | 2  |
| 3        | If $t_{n-1,\delta,\alpha} > 0$ , the bootstrap test does not reject $H_0$ because $S = 0$ . Other-  | 3  |
| 4        | wise, if $t_{n-1,\delta,\alpha} \leq 0$ , the test is randomized and rejects $H_0$ with probability   | 4  |
| 5        | $\{\alpha - 1 + \Phi(-\delta)\}/\Phi(-\delta)$ . Note that the event $t_{n-1,\delta,\alpha} \leq 0$ occurs if and only  | 5  |
| 6        | if $\alpha \geq P(T_n^* > 0   \mu_0, \sigma_0)$ . But   | 6  |
| 7        |   | 7  |
| 8        | $P(T_n^* > 0   \hat{\mu}_0, \hat{\sigma}_0) = P(X_n^* > 0   \hat{\mu}_0, \hat{\sigma}_0) = 1 - \Phi(-\delta).$  | 8  |
| 9        |   | 9  |
| 10       | Therefore $t_{n-1,\delta,\alpha} \leq 0$ if and only if $\delta \leq -z_{\alpha}$ .   | 10 |
| 12       | Let P denote probabilities when $\mu = n$ and $\sigma = \tau$ . The size of the test for  | 12 |
| 13       | Let $\eta_{\eta,\tau}$ denote probabilities when $\mu = \eta$ and $v = \tau$ . The size of the test for $0 < \alpha < 1$ is   | 13 |
| 14       | $0 < \alpha < 1$ 15   | 14 |
| 15       | $\sup P \setminus \{\text{Roject } H_n\}$   | 15 |
| 16       | $H_0$   | 16 |
| 17       | $= \sup[P_{\mu\sigma}\{\text{Reject } H_0, \bar{X}_n > 0\} + P_{\mu\sigma}\{\text{Reject } H_0, \bar{X}_n < 0\}]$   | 17 |
| 18       | $H_0$   | 18 |
| 19       | $= \sup P_{\mu,\sigma}[\{T_n I(T_n > 0) > t_{n-1,0,\alpha}, \bar{X}_n > 0\}$  | 19 |
| 20       | $H_0$ –   | 20 |
| 21       | $+ P_{\mu,\sigma} \{ \text{Reject } H_0,  t_{n-1,\delta,\alpha} \le 0,  X_n < 0 \} ]$   | 21 |
| 22       | $= \sup \{P_{\mu,\sigma}\{T_n > \max(t_{n-1,0,\alpha}, 0)\}\}$  | 22 |
| 23       | $H_0$   | 23 |
| 24       | $+ E_{\mu,\sigma}[\{\alpha - 1 + \Phi(-\delta)\}/\Phi(-\delta)]I\{\delta < \min(-z_{\alpha}, 0)\}]$   | 24 |
| 25       | $\geq P_{0,1}\{T_n > \max(t_{n-1,0,\alpha}, 0)\} + E_{0,1}[\{\alpha - 1 + \Phi(-\delta)\}/\Phi(-\delta)I\{\delta < -z_{\alpha}^+\}]$  | 25 |
| 26       | $\left(\alpha - 1 + \Phi\left(-t - \sqrt{n}\right)\right)$  | 26 |
| 27       | $= \min(\alpha \ 1/2) + E \left\{ \frac{\alpha - 1 + \Psi \left( -t_{n-1} \sqrt{n-1} \right)}{2} I \left( t_{n-1} \sqrt{n} < -z^{+} \right) \right\}$   | 27 |
| 28       | $\Phi\left(-t_{n-1}\sqrt{\frac{n}{n-1}}\right) = \Phi\left(-t_{n-1}\sqrt{\frac{n}{n-1}}\right)$   | 28 |
| 29       | $\begin{pmatrix} & & & & \\ & & & & & \\ & & & & & \end{pmatrix}$   | 29 |
| 30       | by equation (9.5) Cines $t = \sqrt{7}$ in distribution as $r = 1$ as where $7$ is a standard  | 30 |
| 31       | by equation (2.5). Since $\iota_{n-1} \to \mathbb{Z}$ in distribution as $n \to \infty$ , where $\mathbb{Z}$ is a standard normal variable.   | 31 |
| 32       | normai variable,  | 32 |
| 30       |   | 34 |
| 35       | $\alpha - 1 + \Phi\left(-t_{n-1}\sqrt{\frac{n}{n-1}}\right) + \left(-t_{n-1}\sqrt{\frac{n}{n-1}}\right) + \left(-t_{n-1}$ | 35 |
| 36       | $\lim_{n \to \infty} E\left\{ \underbrace{-\frac{1}{2}}_{\Phi\left(-t, -\sqrt{n}\right)} I\left(t_{n-1}\sqrt{n-1} < -z_{\alpha}\right) \right\}$  | 36 |
| 37       | $\left( \begin{array}{c} \Psi \left( -\iota_{n-1} \sqrt{n-1} \right) \end{array} \right) $  | 37 |
| 38       | $\sum_{\mathbf{F}} \left\{ \alpha - 1 + \Phi(-Z) \right\}_{\mathbf{F}(Z) < -\infty^+} $   | 38 |
| 39       | $\rightarrow E \left\{ \frac{-1}{\Phi(-Z)} I(Z < -z_{\alpha}) \right\}$   | 39 |
| 40       | $t^{-z_{\alpha}^+}$   | 40 |
| 41       | $= (\alpha - 1) \int \phi(z)/\Phi(-z) dz + \Phi(-z_{\alpha}^{+})$   | 41 |
| 42       | $J_{-\infty}$   | 42 |
| 43       | $= (\alpha - 1) \int_{-\infty}^{\infty} \phi(z)/\Phi(z) dz + \min(\alpha, 1/2)$   | 43 |
| 44       | $J_{z_{\alpha}^{+}}$  | 44 |
| 45       | $= (1-\alpha)\log\Phi(z_{\alpha}^{+}) + \min(\alpha, 1/2)$  | 45 |
| 46       | $= (1 - \alpha) \log \{ \max(1 - \alpha, 1/2) \} + \min(\alpha, 1/2) \}$  | 46 |
| 47       | $= (1 - \alpha) \log\{1 - \min(\alpha \ 1/2)\} + \min(\alpha \ 1/2)\}$  | 47 |
| 48       | $= (1 \alpha) \log(1 \min(\alpha, 1/2)) + \min(\alpha, 1/2)$  | 48 |
| 49       |   | 49 |
| 50<br>54 | Thus the limiting size is $2\min\{\alpha, 1/2\} + (1 - \alpha)\log\{1 - \min\{\alpha, 1/2\}\}$  | 50 |
| 10       | Thus the minimum size is $2 \min(\alpha, 1/2) + (1 - \alpha) \log\{1 - \min(\alpha, 1/2)\} > \alpha$ .  | 51 |
|          |   |    |
|          | <pre>imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009</pre>   |    |

2.2.2. Cox Likelihood Ratio Statistic З **Theorem 2.5.** If  $\sigma^2$  is unknown, the size of the bootstrap test of (2.1) based on the Cox likelihood ratio has lower bound  $\min(\alpha, 1/2) + P(t_{n-1, t_{n-1}\sqrt{n/(n-1)}, \alpha} < t_{n-1} < 0) \ge \alpha.$ Proof. The Cox log-likelihood ratio statistic is  $\frac{\sum (X_i - \hat{\mu}_0)^2}{\sum (X_i - \hat{\mu}_1)^2} = \frac{\sum \{X_i - \bar{X}_n I(\bar{X}_n < 0)\}^2}{\sum \{X_i - \bar{X}_n I(\bar{X}_n > 0)\}^2} \\ = \begin{cases} \{1 + (n-1)T_n^2\}^{-1}, & \text{if } \bar{X}_n < 0\\ 1, & \text{if } \bar{X}_n = 0\\ 1 + (n-1)T_n^2, & \text{if } \bar{X}_n > 0 \end{cases}$ where  $T_n$  is defined in (2.4). Thus rejecting for large values of the statistic is equiv-alent to rejecting for large values of  $S = T_n$ . 1.  $\bar{X}_n > 0$ . Then  $\hat{\mu}_0 = 0, T_n > 0$ , and  $T_n^*$  has a central t-distribution with n-1degrees of freedom. Thus the test rejects  $H_0$  if  $T_n > t_{n-1,0,\alpha}$ . If  $1/2 \le \alpha < 1$ , then  $t_{n-1,0,\alpha} \leq 0$  and the test rejects w.p.1. 2.  $\bar{X}_n < 0$ . Then  $\hat{\mu}_0 < 0$ ,  $T_n < 0$ , and  $T_n^*$  has a noncentral t-distribution with n-1 degrees of freedom and noncentrality parameter  $\delta$  given in (2.5). Hence  $H_0$  is rejected if  $T_n > t_{n-1,\delta,\alpha}$ . Since  $T_n < 0$ , rejection occurs only if  $t_{n-1,\delta,\alpha} < 0$ If  $0 < \alpha < 1/2$ ,  $\sup_{H_0} P_{\mu,\sigma}(\text{Reject } H_0) = \sup_{H_0} [P_{\mu,\sigma}\{\text{Reject } H_0, \bar{X}_n > 0\}$  $+ P_{\mu,\sigma} \{ \text{Reject } H_0, \bar{X}_n < 0 \} ]$  $= \sup_{H_0} [P_{\mu,\sigma}(T_n > t_{n-1,0,\alpha}) + P_{\mu,\sigma}(t_{n-1,\delta,\alpha} < T_n < 0)]$  $\geq P_{0,1}(T_n > t_{n-1,0,\alpha}) + P_{0,1}(t_{n-1,\delta,\alpha} < T_n < 0)$  $= \alpha + P(t_{n-1,\delta,\alpha} < t_{n-1} < 0)$  $> \alpha$ . If  $1/2 \le \alpha < 1$ ,  $\sup_{H_0} P_{\mu,\sigma}(\text{Reject } H_0) = \sup_{H_0} [P_{\mu,\sigma}\{\text{Reject } H_0, \bar{X}_n > 0\}$  $+ P_{\mu,\sigma} \{ \text{Reject } H_0, \bar{X}_n < 0 \} ]$  $= \sup_{\mu} [P_{\mu,\sigma}(\bar{X}_n > 0) + P_{\mu,\sigma}(t_{n-1,\delta,\alpha} < T_n < 0)]$  $\geq P_{0,1}(\bar{X}_n > 0) + P_{0,1}(t_{n-1,\delta,\alpha} < T_n < 0)$  $= 1/2 + P(t_{n-1,\delta,\alpha} < t_{n-1} < 0)$  $> 1/2 + P(t_{n-1,0,\alpha} < t_{n-1} < 0)$  $= \alpha$ by Lemma 2.2. 

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

3. Testing a Normal Variance, Mean Unknown

Let  $\chi^2_{\nu}$  denote a chi-squared random variable with  $\nu$  degrees of freedom,  $\chi^2_{\nu,\alpha}$  its upper- $\alpha$  point, and  $\Psi_{\nu}(.)$  its cumulative distribution function.

**Lemma 3.1.** 
$$\Psi_{n-1}(n^2/\chi^2_{n-1,\alpha}) \to \alpha \text{ and } \Psi_{n-1}(n) \to 1/2 \text{ as } n \to \infty.$$

*Proof.* Let  $Z_1, Z_2, \ldots$  be independent N(0, 1) variables. Then

$$\Psi_{n-1}(n^2/\chi_{n-1,\alpha}^2) = P\left(\sum_{i=1}^{n-1} Z_i^2 \le n^2/\chi_{n-1,\alpha}^2\right)$$

$$= P\left(\frac{\sum_{i=1}^{n-1} (Z_i^2 - 1)}{\sqrt{2(n-1)}} \le \sqrt{\frac{n-1}{2}} \left\{ \frac{n^2}{(n-1)^2 \chi_{n-1,\alpha}^2} - 1 \right\} \right)$$

$$\Phi\left(\sqrt{(n-1)/2}\left\{\frac{n^2}{(n-1)^2\chi^2_{n-1,\alpha}}-1\right\}\right) \text{ as } n \to \infty.$$

By the Wilson-Hilferty (1931) approximation,  $\nu/\chi^2_{\nu,\alpha} = 1 - z_{\alpha}(2/\nu)^{1/2} + o(\nu^{-1})$ . Therefore

$$\sqrt{(n-1)/2} \left\{ \frac{n^2}{(n-1)^2 \chi^2_{n-1,\alpha}} - 1 \right\} \to -z_{\alpha}$$

which yields the first result. The second result is similarly proved.

Let  $\mathbf{X}_n = (X_1, X_2, \dots, X_n)$  be a vector of n independent observations from  $N(\mu, \sigma^2)$ , with  $\mu$  and  $\sigma$  unknown, and let  $\hat{\sigma}^2 = n^{-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$  denote the unrestricted MLE of  $\sigma^2$ .

# 3.1. $H_0: \sigma^2 < 1$ vs. $H_1: \sigma^2 > 1$

 $\approx$ 

Let  $\hat{\sigma}_i^2$  be the MLE of  $\sigma^2$  under  $H_i$  (i = 0, 1). Then  $\hat{\sigma}_0^2 = \min(\hat{\sigma}^2, 1)$  and  $\hat{\sigma}_1^2 = \max(\hat{\sigma}^2, 1)$ . Define the log-likelihood ratio

$$M(\mu_0, \mu_1, \sigma_0, \sigma_1, \mathbf{X}_n) = \log\left(\frac{\prod_i \sigma_1^{-1} \phi\{\sigma_1^{-1}(x_i - \mu_1)\}}{\prod_i \sigma_0^{-1} \phi\{\sigma_0^{-1}(x_i - \mu_0)\}}\right).$$

3.1.1. Standard Likelihood Ratio Statistic

**Theorem 3.1.** The size of the bootstrap test based on the standard likelihood ratio for testing  $H_0: \sigma^2 \leq 1$  vs.  $H_1: \sigma^2 > 1$ , with  $\mu$  unknown, is bounded below by (3.1)

$$\min\{\alpha, 1-\Psi_{n-1}(n)\} + E\left[\frac{\alpha - 1 + \Psi_{n-1}(n^2/\chi_{n-1}^2)}{\Psi_{n-1}(n^2/\chi_{n-1}^2)} I\left\{\chi_{n-1}^2 \le \min\left(n, \frac{n^2}{\chi_{n-1,\alpha}^2}\right)\right\}\right].$$

*Proof.* The standard log-likelihood ratio statistic is  $M(\hat{\mu}, \hat{\mu}, \hat{\sigma}_0, \hat{\sigma}, \mathbf{X}_n)$  and

$$2n^{-1}M(\hat{\mu},\hat{\mu},\hat{\sigma}_0,\hat{\sigma},\mathbf{X}_n) = \log(\hat{\sigma}_0^2\hat{\sigma}^{-2}) + \hat{\sigma}^2(\hat{\sigma}_0^{-2} - \hat{\sigma}^{-2})$$

$$= \hat{\sigma}^2 \hat{\sigma}_0^{-2} - \log(\hat{\sigma}^2 \hat{\sigma}_0^{-2}) - 1 \tag{49}$$

50  
51
$$= \begin{cases} 0, & \text{if } \hat{\sigma}^2 \leq 1 \\ \hat{\sigma}^2 - \log(\hat{\sigma}^2) - 1, & \text{otherwise.} \end{cases}$$
51

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

З

W. Loh and W. Zheng

Since the function  $x - \log(x) - 1$  increases monotonically from 0 for x > 1, rejecting for large values of the statistic is equivalent to rejecting for large values of З  $S = n \max(\hat{\sigma}^2, 1) = \max\left\{\sum (X_i - \bar{X}_n)^2, n\right\}.$ Let S<sup>\*</sup> denote the bootstrap version of S under resampling from  $N(\bar{X}_n, \hat{\sigma}_0^2)$ . To find the critical point of the distribution of  $S^*$ , consider two cases. 1.  $\hat{\sigma}^2 > 1$ . Then S > n,  $\hat{\sigma}_0^2 = 1$ , and the distribution of  $S^*$  is  $\chi^2_{n-1}$  left-truncated at n, i.e., it has probability mass  $\Psi_{n-1}(n)$  at n. If  $0 < \alpha < 1 - \Psi_{n-1}(n)$ , the critical point of the bootstrap distribution is  $\chi^2_{n-1,\alpha}$ . Otherwise, the critical point is n and the test rejects  $H_0$  w.p.1. 2.  $\hat{\sigma}^2 \leq 1$ . Then S = n,  $\hat{\sigma}_0^2 = \hat{\sigma}^2$ , and the distribution of  $S^*$  is  $\hat{\sigma}^2 \chi^2_{n-1}$  left truncated at n. Thus the test does not reject  $H_0$  if  $\hat{\sigma}^2 \chi^2_{n-1,\alpha} > n$ . On the other hand, if  $\hat{\sigma}^2 \chi^2_{n-1,\alpha} \leq n$ , then the critical point is n and the test rejects  $H_0$  randomly with probability  $\{\alpha - 1 + \Psi_{n-1}(n\hat{\sigma}_0^{-2})\}/\Psi_{n-1}(n\hat{\sigma}_0^{-2}).$ Since  $\alpha < 1 - \Psi_{n-1}(n)$  if and only if  $n < \chi^2_{n-1,\alpha}$ , we have  $P_{\mu,\sigma}(\text{Reject } H_0, \, \hat{\sigma}^2 > 1) = \begin{cases} P_{\mu,\sigma}(S > \chi^2_{n-1,\alpha}, \, n\hat{\sigma}^2 > n), & \text{if } n < \chi^2_{n-1,\alpha} \\ P_{\mu,\sigma}(n\hat{\sigma}^2 > n), & \text{otherwise} \end{cases}$  $= \begin{cases} P_{\mu,\sigma}(n\hat{\sigma}^{2} > \chi^{2}_{n-1,\alpha}), & \text{if } n < \chi^{2}_{n-1,\alpha} \\ P_{\mu,\sigma}(n\hat{\sigma}^{2} > n), & \text{otherwise} \end{cases} \\ = \begin{cases} 1 - \Psi_{n-1}(\sigma^{-2}\chi^{2}_{n-1,\alpha}), & \text{if } n < \chi^{2}_{n-1,\alpha} \\ 1 - \Psi_{n-1}(n\sigma^{-2}), & \text{otherwise} \end{cases}$  $= 1 - \Psi_{n-1}(\sigma^{-2}\max\{n,\chi^2_{n-1}\})$ and  $P_{\mu,\sigma}(\text{Reject } H_0, \hat{\sigma}^2 \leq 1) = P_{\mu,\sigma}(\text{Reject } H_0, \hat{\sigma}^2 \chi^2_{n-1,\sigma} \leq n, n\hat{\sigma}^2 \leq n)$  $= E_{\mu,\sigma} \left[ \frac{\alpha - 1 + \Psi_{n-1}(n\hat{\sigma}_0^{-2})}{\Psi_{n-1}(n\hat{\sigma}_0^{-2})} I(\hat{\sigma}_0^2 \chi_{n-1,\alpha}^2 \le n, n\hat{\sigma}^2 \le n) \right]$  $= E\left[\frac{\alpha - 1 + \Psi_{n-1}(n^2 \sigma^{-2}/\chi_{n-1}^2)}{\Psi_{n-1}(n^2 \sigma^{-2}/\chi_{n-1}^2)}I(\chi_{n-1}^2 \le \sigma^{-2}\min\{n, n^2/\chi_{n-1,\alpha}^2\})\right].$ The choice  $\sigma^2 = 1$  yields the lower bound  $\sup_{H_0} P_{\mu,\sigma}(\text{Reject } H_0)$  $\geq P_{\mu,1}(\text{Reject } H_0, \hat{\sigma}^2 > 1) + P_{\mu,1}(\text{Reject } H_0, \hat{\sigma}^2 \leq 1)$  $= 1 - \Psi_{n-1}(\max\{n, \chi^2_{n-1,\alpha}\})$  $+ E\left[\frac{\alpha - 1 + \Psi_{n-1}(n^2/\chi_{n-1}^2)}{\Psi_{n-1}(n^2/\chi_{n-1}^2)} I\left\{\chi_{n-1}^2 \le \min\left(n, \frac{n^2}{\chi_{n-1}^2\alpha}\right)\right\}\right]$  $\min\{\alpha, 1 - \Psi_{n-1}(n)\}$ =  $+ E\left[\frac{\alpha - 1 + \Psi_{n-1}(n^2/\chi_{n-1}^2)}{\Psi_{n-1}(n^2/\chi_{n-1}^2)} I\left\{\chi_{n-1}^2 \le \min\left(n, \frac{n^2}{\chi_{n-1}^2}\right)\right\}\right].$ Figure 2 shows graphs of the lower bound (3.1) for n = 5, 10, 100, and 500. 



FIG 2. Lower bounds (3.1) on the size of the bootstrap test of  $H_0: \sigma^2 \leq 1$  vs.  $H_1: \sigma^2 > 1$  based on the standard likelihood ratio, for n = 5, 10, 100, and 500. The 45-degree line is the identity function.

#### 3.1.2. Cox Likelihood Ratio Statistic

**Theorem 3.2.** If  $\mu$  is unknown, the bootstrap test of  $H_0: \sigma^2 \leq 1$  vs.  $H_1: \sigma^2 > 1$ based on the Cox likelihood ratio has size  $\alpha$  and is UMP for  $\chi^2_{n-1,\alpha} > n$ . It rejects  $H_0$  w.p.1 for other values of  $\alpha$ .

*Proof.* The Cox log-likelihood ratio statistic is  $M(\hat{\mu}, \hat{\mu}, \hat{\sigma}_0, \hat{\sigma}_1, \mathbf{X}_n)$ . Since

$$\hat{\sigma}_{0}^{2} \, \hat{\sigma}_{1}^{-2} = \begin{cases} \hat{\sigma}^{2}, & \text{if } \hat{\sigma}^{2} \leq 1 \\ \hat{\sigma}^{-2}, & \text{if } \hat{\sigma}^{2} > 1 \end{cases}$$

and

$$\hat{\sigma}_0^{-2} - \hat{\sigma}_1^{-2} = \begin{cases} \hat{\sigma}^{-2} - 1, & \text{if } \hat{\sigma}^2 \le 1\\ 1 - \hat{\sigma}^{-2}, & \text{if } \hat{\sigma}^2 > 1 \end{cases}$$
<sup>37</sup>
<sup>38</sup>

we have

$$2n^{-1}M(\hat{\mu},\hat{\mu},\hat{\sigma}_0,\hat{\sigma}_1,\mathbf{X}_n) = \log(\hat{\sigma}_0^2\hat{\sigma}_1^{-2}) + \hat{\sigma}^2(\hat{\sigma}_0^{-2} - \hat{\sigma}_1^{-2})$$

$$= \begin{cases} \log(\hat{\sigma}^2) - \hat{\sigma}^2 + 1, & \text{if } \hat{\sigma}^2 \le 1 \\ \log(\hat{\sigma}^2) + \hat{\sigma}^2 - 1 & \text{if } \hat{\sigma}^2 \le 1 \end{cases}$$

$$\left\{ -\log(\hat{\sigma}^2) + \hat{\sigma}^2 - 1, \text{ if } \hat{\sigma}^2 > 1 \right\}$$

which is strictly increasing in  $\hat{\sigma}^2$ . Therefore rejecting for large values of the statistic is equivalent to rejecting for large values of  $\hat{\sigma}^2$ . Since the bootstrap null distribution of  $n\hat{\sigma}^2$  is  $\hat{\sigma}_0^2\chi_{n-1}^2$ , the bootstrap critical point of  $\hat{\sigma}^2$  is  $n^{-1}\hat{\sigma}_0^2\chi_{n-1,\alpha}^2$ . Thus the bootstrap test rejects  $H_0$  if  $\hat{\sigma}^2\hat{\sigma}_0^{-2} > n^{-1}\chi_{n-1,\alpha}^2$ , or equivalently,  $\max(\hat{\sigma}^2, 1) >$  $n^{-1}\chi_{n-1,\alpha}^2$ . If  $\alpha$  is so large that  $n^{-1}\chi_{n-1,\alpha}^2 \leq 1$ , the bootstrap test rejects  $H_0$ regardless of the data. On the other hand, if  $n^{-1}\chi_{n-1,\alpha}^2 > 1$ , the test rejects  $H_0$  if  $\hat{\sigma}^2 > n^{-1}\chi_{n-1,\alpha}^2$ , which coincides with the UMP test [10, p. 88].

3.2.  $H_0: \sigma^2 > 1$  vs.  $H_1: \sigma^2 < 1$ Next suppose we reverse the hypotheses and test  $H_0: \sigma^2 \ge 1$  versus  $H_1: \sigma^2 < 1$ . Then  $\hat{\sigma}_0^2 = \max(\hat{\sigma}^2, 1)$  and  $\hat{\sigma}_1^2 = \min(\hat{\sigma}^2, 1)$ . 3.2.1. Standard Likelihood Ratio Statistic **Theorem 3.3.** For any  $0 < \alpha < 1$  and  $\mu$  unknown, the size of the bootstrap test for  $H_0: \sigma^2 \ge 1$  vs.  $H_1: \sigma^2 < 1$ , based on the standard likelihood ratio, is bounded below by (3.2) $\min\{\alpha, \Psi_{n-1}(n)\} + E\left[\frac{\alpha - \Psi_{n-1}(n^2/\chi_{n-1}^2)}{1 - \Psi_{n-1}(n^2/\chi_{n-1}^2)}I\left(\chi_{n-1}^2 \ge \max\{n, n^2/\chi_{n-1, 1-\alpha}^2\}\right)\right].$ *Proof.* Direct computation yields  $2n^{-1}M(\hat{\mu},\hat{\mu},\hat{\sigma}_0,\hat{\sigma},\mathbf{X}_n) = \begin{cases} \hat{\sigma}^2 - \log(\hat{\sigma}^2) - 1, & \text{if } \hat{\sigma}^2 \le 1\\ 0, & \text{otherwise.} \end{cases}$ Since the function  $x - \log(x) - 1$  decreases monotonically for  $0 < x \leq 1$ , the test rejects  $H_0$  for small values of  $S = n \min(\hat{\sigma}^2, 1)$ . Let  $S^*$  denote the bootstrap version of S under resampling from  $N(\bar{X}_n, \hat{\sigma}_0^2)$ . 1.  $\hat{\sigma}^2 < 1$ . Then  $\hat{\sigma}_0 = 1$  and the distribution of  $S^*$  is  $\chi^2_{n-1}$  right-truncated at n, with probability mass  $1 - \Psi_{n-1}(n)$  there. If  $0 < \alpha < \Psi_{n-1}(n)$ , the critical point of the bootstrap distribution is  $\chi^2_{n-1,1-\alpha}$ . Otherwise, the critical point is n and the test rejects w.p.1, because S < n. 2.  $\hat{\sigma}^2 \geq 1$ . Then  $\hat{\sigma}_0^2 = \hat{\sigma}^2$ , S = n and the distribution of  $S^*$  is  $\hat{\sigma}^2 \chi^2_{n-1}$  right-truncated at *n*. The test does not reject  $H_0$  if  $\hat{\sigma}^2 \chi^2_{n-1,1-\alpha} < n$ . Otherwise, if  $\hat{\sigma}^2 \chi^2_{n-1,1-\alpha} \ge n$ , the test rejects  $H_0$  with probability  $\{\alpha - \Psi_{n-1}(n\hat{\sigma}^{-2})\}/\{1 - \Psi_{n-1}(n\hat{\sigma}^{-2})\}$ . Since  $\alpha < \Psi_{n-1}(n)$  if and only if  $\chi^2_{n-1,1-\alpha} < n$ ,  $P_{\mu\sigma}$  (Reject  $H_0, \hat{\sigma}^2 < 1$ )  $= \begin{cases} P_{\mu,\sigma}(S < \chi^2_{n-1,1-\alpha}, n\hat{\sigma}^2 < n), & \text{if } 0 < \alpha < \Psi_{n-1}(n) \\ P_{\mu,\sigma}(n\hat{\sigma}^2 < n), & \text{otherwise} \end{cases}$  $= \begin{cases} P_{\mu,\sigma}(n\hat{\sigma}^2 < \chi^2_{n-1,1-\alpha}, n\hat{\sigma}^2 < n), & \text{if } 0 < \alpha < \Psi_{n-1}(n) \\ P_{\mu,\sigma}(n\hat{\sigma}^2 < n), & \text{otherwise} \end{cases}$  $= \begin{cases} P_{\mu,\sigma}(n\hat{\sigma}^2 < \chi^2_{n-1,1-\alpha}), & \text{if } 0 < \alpha < \Psi_{n-1}(n) \\ P_{\mu,\sigma}(n\hat{\sigma}^2 < n), & \text{otherwise} \end{cases}$  $= P_{\mu,\sigma}(n\hat{\sigma}^2 < \min\{\chi^2_{n-1,1-\alpha}, n\})$  $= \Psi_{n-1}(\sigma^{-2}\min\{\chi_{n-1,1-\alpha}^2,n\})$ and  $P_{\mu,\sigma}(\text{Reject } H_0, \hat{\sigma}^2 > 1)$ =  $P_{\mu,\sigma}$  (Reject  $H_0, \hat{\sigma}^2 \chi^2_{n-1,1-\sigma} \ge n, \hat{\sigma}^2 \ge 1$ )  $= E_{\mu,\sigma} \left[ \frac{\alpha - \Psi_{n-1}(n\hat{\sigma}^{-2})}{1 - \Psi_{n-1}(n\hat{\sigma}^{-2})} I(\hat{\sigma}^2 \chi_{n-1,1-\alpha}^2 \ge n, n\hat{\sigma}^2 \ge n) \right]$  $= E\left[\frac{\alpha - \Psi_{n-1}(n^2\sigma^{-2}/\chi^2_{n-1})}{1 - \Psi_{n-1}(n^2\sigma^{-2}/\chi^2_{n-1})}I(\sigma^2\chi^2_{n-1} \ge \max\{n, n^2/\chi^2_{n-1, 1-\alpha}\})\right].$ 



FIG 3. Lower bounds (3.2) on the size of the bootstrap test of  $H_0: \sigma^2 \ge 1$  vs.  $H_1: \sigma^2 < 1$  based on the standard likelihood ratio  $M_n^{(1)}$ . The 45-degree line is the identity function.

## 3.2.2. Cox Likelihood Ratio Statistic

**Theorem 3.4.** The bootstrap test of  $H_0: \sigma^2 \ge 1$  vs.  $H_1: \sigma^2 < 1$  based on the Cox likelihood ratio has size  $\alpha$  and is UMP if  $\chi^2_{n-1,\alpha} > n$ . Otherwise, it rejects  $H_0$  w.p.1.

*Proof.* Since

$$\hat{\sigma}_0^2 \hat{\sigma}_1^{-2} = \begin{cases} \hat{\sigma}^{-2}, & \text{if } \hat{\sigma}^2 < 1\\ \hat{\sigma}^2, & \text{if } \hat{\sigma}^2 \ge 1 \end{cases}$$

$$\hat{\sigma}_0^{-2} - \hat{\sigma}_1^{-2} = \begin{cases} 1 - \hat{\sigma}^{-2}, & \text{if } \hat{\sigma}^2 < 1 \\ \hat{\sigma}^{-2} - 1, & \text{if } \hat{\sigma}^2 \ge 1 \end{cases}$$

З

we have

$$2n^{-1}M(\hat{\mu},\hat{\mu},\hat{\sigma}_0,\hat{\sigma}_1,\mathbf{X}_n) = \log(\hat{\sigma}_0^2\hat{\sigma}_1^{-2}) + \hat{\sigma}^2(\hat{\sigma}_0^{-2} - \hat{\sigma}_1^{-2})$$

$$\int -\log(\hat{\sigma}^2) + \hat{\sigma}^2 - 1$$
, if  $\hat{\sigma}^2 < 1$ 

$$= \begin{cases} \log(\hat{\sigma}^2) + \hat{\sigma}^2 + 1, & \text{if } \hat{\sigma}^2 \ge 1 \\ \log(\hat{\sigma}^2) - \hat{\sigma}^2 + 1, & \text{if } \hat{\sigma}^2 \ge 1 \end{cases}$$

a strictly decreasing function of  $\hat{\sigma}^2$ . Thus the test statistic rejects  $H_0$  for small values of  $\hat{\sigma}^2$ . The bootstrap null distribution of  $\hat{\sigma}^2$  is  $n^{-1}\hat{\sigma}_0^2\chi_{n-1}^2$ , with critical value  $n^{-1}\hat{\sigma}_0^2\chi_{n-1,1-\alpha}^2$ . Hence the bootstrap test rejects  $H_0$  if  $\hat{\sigma}^2 \hat{\sigma}_0^{-2} < n^{-1}\chi_{n-1,1-\alpha}^2$ . But the left side of the inequality is never greater than 1, because

$$\hat{\sigma}^2 \, \hat{\sigma}_0^{-2} = \begin{cases} \hat{\sigma}^2, & \text{if } \hat{\sigma}^2 < 1 \\ 1, & \text{otherwise.} \end{cases}$$

Therefore, if  $\alpha$  is so large that  $n^{-1}\chi^2_{n-1,1-\alpha} \geq 1$ , the bootstrap test rejects  $H_0$  w.p.1. Otherwise, if  $n^{-1}\chi^2_{n-1,1-\alpha} < 1$ , the test rejects  $H_0$  if  $\hat{\sigma}^2 < n^{-1}\chi^2_{n-1,1-\alpha}$ , which coincides with the classical UMP unbiased test [10, pp. 154].

## 4. Testing Difference of Two Normal Means

Let  $X_1, \ldots, X_m$  and  $Y_1, \ldots, Y_n$  be independent random samples from  $N(\mu, \sigma^2)$  and  $N(\eta, \tau^2)$ , respectively, and N = m + n > 2. We want to test

(4.1) 
$$H_0: \eta \le \mu \quad \text{vs.} \quad H_1: \eta > \mu.$$

The likelihood function for this case is

$$L(\mu, \tau)$$

$$= (2\pi)^{-(m+n)/2} \sigma^{-m} \tau^{-n} \exp\left\{-(2\sigma^2)^{-1} \sum (X_i - \mu)^2 - (2\tau^2)^{-1} \sum (Y_j - \eta)^2\right\}$$

$$= (2\pi)^{-(m+n)/2} \sigma^{-m} \tau^{-n} \exp\{-(2\sigma^2)^{-1} \sum (X_i - X_m)^2 - (2\tau^2)^{-1} \sum (Y_j - Y_n)^2 m(2\sigma^2)^{-1} (\mu - \bar{\mathbf{X}})^2 + m(2\sigma^2)^{-1} (\mu - \bar{\mathbf{X}})^2 \}$$

$$-m(2\sigma^{-})^{-1}(\mu - X_m)^{-} - n(2\tau^{-})^{-1}(\eta - Y_n)^{-}\}$$

and the unrestricted MLE of  $(\mu, \eta)$  is  $(\hat{\mu}, \hat{\eta}) = (\bar{X}_m, \bar{Y}_n)$ .

## 4.1. Known Variances

Let 
$$V = (m\tau^2 \bar{X}_m + n\sigma^2 \bar{Y}_n)/(m\tau^2 + n\sigma^2)$$
. The MLE of  $(\mu, \eta)$  under  $H_0$  is

$$(\hat{\mu}_0, \hat{\eta}_0) = \begin{cases} (\bar{X}_m, \bar{Y}_n), & \bar{Y}_n \leq \bar{X}_m \\ (V, V), & \bar{Y}_n > \bar{X}_m \end{cases}$$

and that under  $H_1$  is

$$(\hat{\mu}_{1}, \hat{\eta}_{1}) = \begin{cases} (V, V), & \bar{Y}_{n} > \bar{X}_{m} \\ (\bar{X}_{m}, \bar{Y}_{n}), & \bar{Y}_{n} \le \bar{X}_{m}. \end{cases}$$

4.1.1. Difference of Means Statistic

**Theorem 4.1.** The size of the bootstrap test of (4.1) based on  $\bar{Y}_n - \bar{X}_m$  is  $\alpha$  if  $\alpha < 1/2$  and is 1 if  $\alpha \ge 1/2$ . 

З

#### Bootstrap tests

*Proof.* Let  $S = \bar{Y}_n - \bar{X}_m$ . The bootstrap test statistic  $S^* = \bar{Y}_n^* - \bar{X}_m^*$  has a normal distribution with mean  $\hat{\eta}_0 - \hat{\mu}_0 = S I(S < 0)$  and variance  $m^{-1}\sigma^2 + n^{-1}\tau^2$ . Thus its nominal level- $\alpha$  bootstrap critical value is  $SI(S < 0) + z_{\alpha} \{m^{-1}\sigma^2 + n^{-1}\tau^2\}^{1/2}$ and the rejection region is  $\max(S, 0) > z_{\alpha} \{m^{-1}\sigma^2 + n^{-1}\tau^2\}^{1/2}$ . Clearly, the size of the test is attained at the boundary  $\mu = \eta$ . If  $\alpha < 1/2$ , the probability of rejecting  $H_0$  when  $\mu = \eta$  is exactly  $\alpha$ . On the other hand, if  $\alpha \ge 1/2$ , then  $z_{\alpha} \le 0$  and the test rejects  $H_0$  w.p.1. 4.1.2. Standard Likelihood Ratio Statistic **Theorem 4.2.** The size of the bootstrap test of (4.1) based on the standard likeli-hood ratio is  $\min(2\alpha, 1) + (1 - \alpha) \log\{1 - \min(\alpha, 1/2)\} > \alpha$ . *Proof.* The log-likelihood ratio statistic is  $\log\{L(\hat{\mu}, \hat{\tau})/L(\hat{\mu}_0, \hat{\tau}_0)\}$  $= \{m(2\sigma^2)^{-1}(V-\bar{X}_m)^2 + n(2\tau^2)^{-1}(V-\bar{Y}_n)^2\}I(\bar{Y}_n > \bar{X}_m)$  $= mn(m\tau^2 + n\sigma^2)^{-2}(\bar{Y}_n - \bar{X}_m)^2 I(\bar{Y}_n > \bar{X}_m).$ Thus the test statistic is equivalent to  $S = (\bar{Y}_n - \bar{X}_m) I(\bar{Y}_n > \bar{X}_m)$ . The bootstrap distribution of  $S^*$  is normal with mean  $\hat{\eta}_0 - \hat{\mu}_0$  and variance  $n^{-1}\tau^2 + m^{-1}\sigma^2$ , left-truncated at 0 with  $P(S^* = 0) = \Phi\{(\hat{\mu}_0 - \hat{\eta}_0)(n^{-1}\tau^2 + m^{-1}\sigma^2)^{-1/2}\}$ . Let  $\delta = (\mu - \eta)(n^{-1}\tau^2 + m^{-1}\sigma^2)^{-1/2}$  and  $W = (\bar{X}_m - \bar{Y}_n)(n^{-1}\tau^2 + m^{-1}\sigma^2)^{-1/2} \sim N(\delta, 1)$ . We consider two cases. 1.  $\bar{Y}_n \leq \bar{X}_m$ . Then S = 0,  $\hat{\eta}_0 - \hat{\mu}_0 = \bar{Y}_n - \bar{X}_m$ , and  $\Phi(W) \geq 1/2$ . If  $1 - \Phi(W) < \alpha$ , the test is randomized and rejects  $H_0$  with probability  $\{\alpha - 1 + \Phi(W)\}/\Phi(W)$ . Otherwise, if  $1 - \Phi(W) \ge \alpha$ , the test does not reject  $H_0$ . 2.  $\bar{Y}_n > \bar{X}_m$ . Then  $S = \bar{Y}_n - \bar{X}_m > 0$ ,  $\hat{\eta}_0 - \hat{\mu}_0 = 0$ , and  $P(S^* = 0) = 1/2$ . If  $\alpha < 1/2$ , then the test rejects  $H_0$  if  $\bar{Y}_n - \bar{X}_m > z_\alpha (n^{-1}\tau^2 + m^{-1}\sigma^2)^{-1/2}$ , i.e.,  $W < -z_{\alpha}$ . Otherwise, if  $\alpha \geq 1/2$ , then the critical value is 0 and the test rejects w.p.1. Therefore  $P(\text{Reject } H_0)$  $= P(\text{Reject } H_0, \bar{Y}_n \leq \bar{X}_m, 1 - \Phi(W) < \alpha)$ + P(Reject  $H_0, \bar{Y}_n > \bar{X}_m) I(\alpha < 1/2) + P(\text{Reject } H_0, \bar{Y}_n > \bar{X}_m) I(\alpha \ge 1/2)$  $= E[\Phi(W)^{-1}\{\alpha - 1 + \Phi(W)\} I(W > z_{\alpha}^{+})]$  $+ P(W < -z_{\alpha}) I(\alpha < 1/2) + P(W < 0) I(\alpha > 1/2)$  $= E[\Phi(W)^{-1}\{\alpha - 1 + \Phi(W)\} I(W > z_{\alpha}^{+})] + P(W < -z_{\alpha}^{+})$ and the result follows from Lemma 2.1. 4.1.3. Cox Likelihood Ratio Statistic **Theorem 4.3.** The bootstrap test of (4.1) based on the Cox likelihood ratio statistic is the same as that based on the difference of sample means; its size is  $\alpha$  if  $\alpha < 1/2$ and is 1 if  $\alpha \geq 1/2$ . 

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

*Proof.* The Cox log-likelihood ratio statistic is

$$\log\left\{\frac{L(\hat{\mu}_1, \hat{\tau}_1)}{L(\hat{\mu}_0, \hat{\tau}_0)}\right\} = \frac{mn(\bar{Y}_n - \bar{X}_m)^2}{2(m\tau^2 + n\sigma^2)} \{I(\bar{Y}_n > \bar{X}_m) - I(\bar{Y}_n \le \bar{X}_m)\}.$$

З

Thus the test statistic is equivalent to  $S = \overline{Y}_n - \overline{X}_m$  and the result follows from Theorem 4.1.

## 4.2. Unknown but Equal Variances

Suppose that  $\tau^2 = \sigma^2$  but their value is unknown. Then the likelihood function is

$$L(\mu,\tau,\sigma) = (2\pi\sigma^2)^{-N/2} \exp\left[-(2\sigma^2)^{-1} \left\{\sum_i (X_i - \mu)^2 + \sum_j (Y_j - \eta)^2\right\}\right]$$

giving the unrestricted MLE

$$(\hat{\mu}, \hat{\eta}, \hat{\sigma}^2) = \left( \bar{X}_m, \bar{Y}_n, N^{-1} \left\{ \sum_i (X_i - \bar{X}_m)^2 + \sum_j (Y_j - \bar{Y}_n)^2 \right\} \right).$$

Let  $V = N^{-1}(m\bar{X}_m + n\bar{Y}_n)$  and

$$\widetilde{\sigma}^{2} = N^{-1} \left\{ \sum_{i} (X_{i} - V)^{2} + \sum_{i} (Y_{j} - V)^{2} \right\}$$

$$22$$
23
24
23
24

$$= \hat{\sigma}^2 + mnN^{-2}(\bar{Y}_n - \bar{X}_m)^2.$$
<sup>25</sup>
<sup>26</sup>

The corresponding MLEs under  $H_0$  and  $H_1$  are, respectively,

$$(\hat{\mu}_0, \hat{\eta}_0, \hat{\sigma}_0^2) = \begin{cases} (\bar{X}_m, \bar{Y}_n, \hat{\sigma}^2), & \text{if } \bar{Y}_n < \bar{X}_m \\ (V, V, \tilde{\sigma}^2), & \text{if } \bar{Y}_n \ge \bar{X}_m \end{cases}$$

$$(\hat{\mu}_{1}, \hat{\eta}_{1}, \hat{\sigma}_{1}^{2}) = \begin{cases} (V, V, \tilde{\sigma}^{2}), & \text{if } \bar{Y}_{n} < \bar{X}_{m} \\ (\bar{X} - \bar{V} - \hat{\sigma}^{2}), & \text{if } \bar{Y}_{n} < \bar{X}_{m} \end{cases}$$
<sup>31</sup>
<sup>32</sup>

$$(X_m, Y_n, \hat{\sigma}^2), \quad \text{if } Y_n \ge X_m.$$

4.2.1. Difference of Means Statistic

Suppose  $S = \bar{Y}_n - \bar{X}_m$ . Then  $S^* = \bar{Y}_n^* - \bar{X}_m^*$  has a normal distribution with mean  $\hat{\eta}_0 - \hat{\mu}_0$  and variance  $N(mn)^{-1}\hat{\sigma}_0^2$ . Let  $\Upsilon_{\nu}$  denote the *t* distribution function with  $\nu$  degrees of freedom and let  $s^2 = \hat{\sigma}^2 N(N-2)^{-1}$  be the usual unbiased estimate of  $\sigma^2$ .

**Theorem 4.4.** The size of the bootstrap test of (4.1) based on  $\bar{Y}_n - \bar{X}_m$  is

$$\sup P(Reject H_0)$$

$$H_0$$
 ( 0,

$$\begin{array}{ll} 45\\ 46\\ 47 \end{array} & (4.2) \\ 47 \end{array} = \begin{cases} 0, & if \ \alpha \le 1 - \Psi(\sqrt{N}) \\ 1 - \Upsilon_{N-2}\left(z_{\alpha}\sqrt{\frac{N-2}{N-z_{\alpha}^{2}}}\right), & if \ 1 - \Phi(\sqrt{N}) < \alpha < 1/2\\ 1, & if \ \alpha \ge 1/2 \end{cases}$$

50 
$$(1, if \alpha \ge 1/2)$$
 50  
51  $as N \to \infty.$  51

if  $\alpha < 1 - \Phi(\sqrt{N})$ 

 $Bootstrap \ tests$ 

*Proof.* The hypothesis  $H_0$  is rejected if

$$\bar{Y}_n - \bar{X}_m > \hat{\eta}_0 - \hat{\mu}_0 + z_\alpha \hat{\sigma}_0 \sqrt{N/(mn)}$$

$$= \begin{cases} \bar{Y}_n - \bar{X}_m + z_\alpha \hat{\sigma} \sqrt{N/(mn)}, & \text{if } \bar{Y}_n < \bar{X}_m \\ z_\alpha \tilde{\sigma} \sqrt{N/(mn)}, & \text{if } \bar{Y}_n \ge \bar{X}_m. \end{cases}$$

1.  $\alpha < 1/2$ . If  $\bar{Y}_n < \bar{X}_m$ , the test does not reject  $H_0$ . Otherwise, if  $\bar{Y}_n \ge \bar{X}_m$ , the test rejects  $H_0$  if

$$(\bar{Y}_n - \bar{X}_m)^2 > z_\alpha^2 \tilde{\sigma}^2 N/(mn)$$

 $1 - \Upsilon_{N-2} \left( z_{\alpha} \sqrt{(N-2)/(N-z_{\alpha}^2)} \right).$ 2.  $\alpha \ge 1/2$ . The test rejects  $H_0$  w.p.1 because  $z_{\alpha} < 0$ .

$$\iff (\bar{Y}_n - \bar{X}_m)^2 (1 - N^{-1} z_\alpha^2) > N z_\alpha^2 \hat{\sigma}^2 / (mn)$$

Therefore if  $\alpha \leq 1 - \Phi(\sqrt{N})$ , the test does not reject  $H_0$ . Otherwise, if  $1 - \Phi(\sqrt{N})$ 

 $\Phi(\sqrt{N}) < \alpha < 1/2$ , the rejection probability is maximized when  $\eta = \mu$  at

Hence the result (4.2). The limit is due to  $\Upsilon_{\nu}(x) \to \Phi(x)$  as  $\nu \to \infty$ , for every x.

$$\iff \quad \alpha > 1 - \Phi(\sqrt{N}) \quad \text{and} \quad \sqrt{\frac{mn}{N}} \, \frac{\bar{Y}_n - \bar{X}_m}{s} > \sqrt{\frac{N-2}{N-z_\alpha^2}}.$$



FIG 4. Size of bootstrap test for a difference of normal means based on the difference of sample means, for equal but unknown variances (4.2). The 45-degree line is the identity function.

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

З

| 1  | 4.2.2. Standard Likelihood Ratio Statistic   | 1  |
|--|--|--|
| 2<br>3   | The log-likelihood ratio statistic is  | 2<br>3   |
| 4  | $(N/2)\log(\tilde{\sigma}^2/\hat{\sigma}^2)I(\bar{Y}_n > \bar{X}_m) = (N/2)\log\{1 + mnN^{-2}(\bar{Y}_n - \bar{X}_m)^2\hat{\sigma}^{-2}\}I(\bar{Y}_n > \bar{X}_m)$   | 4  |
| 6  | which is equivalent to the positive part of the $t$ -statistic:  | 5<br>6   |
| 7  | $S = \sqrt{mn/N}s^{-1}(\bar{Y}_n - \bar{X}_m)^+.$  | 7  |
| 9<br>10  | <b>Theorem 4.5.</b> The size of the bootstrap test of $(4.1)$ based on the standard likelihood ratio is bounded below by   | 9<br>10  |
| 11<br>12<br>13<br>14<br>15                         | $\min(\alpha, 1/2) + E\left[\frac{\alpha - 1 + \Phi\left(\sqrt{N/(N-2)}t_{N-2}\right)}{\Phi\left(\sqrt{N/(N-2)}t_{N-2}\right)} I\left(\sqrt{N/(N-2)}t_{N-2} \ge z_{\alpha}^{+}\right)\right] \\ \to \min(2\alpha, 1) + (1-\alpha)\log\{1 - \min(\alpha, 1/2)\},  N \to \infty.$  | 11<br>12<br>13<br>14<br>15                         |
| 16<br>17   | <i>Proof.</i> We again consider two situations.  | 16<br>17   |
| 18<br>19<br>20<br>21<br>22<br>23<br>24<br>25<br>26 | <ol> <li>\$\bar{Y}_n &gt; \bar{X}_m\$. The bootstrap distribution of \$S^*\$ is a \$t_{N-2}\$ distribution left-truncated at 0 with probability 1/2. If \$\alpha &lt; 1/2\$, then the test rejects \$H_0\$ if \$S &gt; t_{N-2,\alpha\$}\$. Otherwise, if \$\alpha \ge 1/2\$, the test rejects \$H_0\$ w.p.1\$.</li> <li>\$\bar{Y}_n \le \bar{X}_m\$. The bootstrap distribution of \$S^*\$ consists of the positive part of a noncentral \$t\$ with \$N-2\$ degrees of freedom and noncentrality parameter \$\delta = \sqrt{mn/N} \bar{\alpha^{-1}}(\bar{Y}_n - \bar{X}_m) = \$S\sqrt{N/(N-2)}\$ and probability at 0 equal to \$\Phi(-\delta)\$. Since \$S = 0\$, the test does not reject \$H_0\$ if \$\alpha &lt; 1 - Φ(-\delta)\$ \$\epsilon\$ = \$\delta_\$\$, then the test is randomized, rejecting \$H_0\$ with probability \$\{\alpha - 1 + Φ(-\delta)\}/Φ(-\delta)\$.</li> </ol>   | 18<br>19<br>20<br>21<br>22<br>23<br>24<br>25<br>26 |
| 27<br>28   | Thus   | 27<br>28   |
| 29<br>30<br>31<br>32                               | $P(\text{Reject } H_0)$ $= P(S > t_{N-2,\alpha}, \bar{Y}_n > \bar{X}_m) I(\alpha < 1/2) + P(\bar{Y}_n > \bar{X}_m) I(\alpha \ge 1/2)$ $+ P(\text{Reject } H_0, -\delta \ge z_\alpha, \bar{Y}_n \le \bar{X}_m)$   | 29<br>30<br>31<br>32                               |
| 33<br>34<br>35                                     | $= P(S > t_{N-2,\alpha}) I(\alpha < 1/2) + (1/2)I(\alpha \ge 1/2) + E\left\{\frac{\alpha - 1 + \Phi(-\delta)}{\Phi(-\delta)} I(-\delta \ge z_{\alpha}, \bar{Y}_{n} \le \bar{X}_{m})\right\}$   | 33<br>34<br>35                                     |
| 36<br>37<br>38<br>39                               | $= \min(\alpha, 1/2) + E\left\{\frac{\alpha - 1 + \Phi(-\delta)}{\Phi(-\delta)}I(-\delta \ge z_{\alpha}, \bar{Y}_{n} \le \bar{X}_{m})\right\}$ $= \min(\alpha, 1/2) + E\left\{\frac{\alpha - 1 + \Phi(-\delta)}{\Phi(-\delta)}I(-\delta \ge z_{\alpha}^{+})\right\}.$  | 36<br>37<br>38<br>39                               |
| 40<br>41   | Evolution this probability at $\mu$ , projection   | 40<br>41   |
| 42<br>43   | Evaluating this probability at $\mu = \eta$ yields<br>sup $P(\text{Reject } H_0)$  | 42<br>43   |
| 44   | $H_0$ $\geq \min(\alpha, 1/2)$   | 44   |
| 45<br>46<br>47                                     | $= E\left[\frac{\alpha - 1 + \Phi\left(\sqrt{N/(N-2)}t_{N-2}\right)}{E\left(\sqrt{N/(N-2)}t_{N-2}\right)}I\left(\sqrt{N/(N-2)}t_{N-2} \ge z_{\alpha}^{+}\right)\right]$  | 45<br>46<br>47                                     |
| 48<br>49   | $\left[ \begin{array}{c} \Phi\left(\sqrt{N/(N-2)t_{N-2}}\right) \\ \vdots \\ \left(2 - 1\right) + \left(1 - 1\right) \\ 1 - \left(1 - 1\right) \\$ | 48<br>49   |
| 50   | $\rightarrow \min(2\alpha, 1) + (1 - \alpha)\log\{1 - \min(\alpha, 1/2)\}$   | 50   |
| 51   | as $N \to \infty$ by Lemma 2.1.  | 51   |

4.2.3. Cox Likelihood Ratio Statistic

**Theorem 4.6.** The size of the bootstrap test of (4.1) based on the Cox likelihood ratio or the ordinary t-statistic is

$$P(t_{N-2,t_{N-2}\sqrt{N/(N-2)},\alpha} < t_{N-2} < 0) + P(t_{N-2} > t_{N-2,\alpha}^+) \ge \alpha.$$

*Proof.* The Cox log-likelihood ratio simplifies to

$$(N/2)\log\{1+mn(\bar{Y}_n-\bar{X}_m)^2N^{-2}\hat{\sigma}^{-2}\}\{(I(\bar{Y}_n\geq\bar{X}_m)-I(\bar{Y}_n<\bar{X}_m)\}$$

which is an increasing function of the Student t statistic  $S = \sqrt{mn/N}s^{-1}(\bar{Y}_n - \bar{X}_m)$ . The bootstrap distribution of S is a noncentral  $t_{N-2,\delta}$  with N-2 degrees of freedom and noncentrality parameter

$$\delta = \sqrt{mn/N}\hat{\sigma}_{0}^{-1}(\hat{\eta}_{0} - \hat{\mu}_{0})$$

$$\int \sqrt{mn/N}\hat{\sigma}_{0}^{-1}(\bar{Y}_{n} - \bar{X}_{m}) \quad \text{if } \bar{Y}_{n} \leq \bar{X}_{m}$$
18

$$\begin{cases} \sqrt{mn/N} \sigma^{-1}(Y_n - X_m), & \text{if } Y_n < X_m \\ 0, & \text{if } \bar{Y}_n > \bar{X}_m \end{cases}$$

$$\begin{cases} \sqrt{N/(N-2)}S, & \text{if } \bar{Y}_n < \bar{X}_m \\ \sqrt{N/(N-2)}S, & \text{if } \bar{Y}_n < \bar{X}_m \end{cases}$$

$$= \begin{cases} \sqrt{1} - \sqrt{1$$

Therefore

$$P(\text{Reject } H_0) = P(S > t_{N-2,\delta,\alpha}, \bar{Y}_n < \bar{X}_m) + P(S > t_{N-2,\alpha}, \bar{Y}_n \ge \bar{X}_m)$$
  
=  $P(t_{N-2,\delta,\alpha} < S < 0) + P(S > t_{N-2,\alpha}^+).$ 

Evaluating the probabilities at  $\mu = \eta$  yields

=

$$\sup_{H_0} P(\text{Reject } H_0) \geq P(t_{N-2, t_{N-2}}\sqrt{N/(N-2)}, \alpha < t_{N-2} < 0) + P(t_{N-2} > t_{N-2, \alpha}^+)$$
  
$$\geq P(t_{N-2, 0, \alpha} < t_{N-2} < 0) + \min(\alpha, 1/2)$$
  
$$= (\alpha - 1/2) I(\alpha > 1/2) + \min(\alpha, 1/2)$$
  
$$= \alpha.$$

## 5. Testing an Exponential Location Parameter

Let  $\operatorname{Exp}(\theta, \tau)$  denote the distribution with density  $\tau^{-1} \exp\{-\tau^{-1}(x-\theta)\}, x \geq \theta$ . We consider testing hypotheses about  $\theta$  with  $\tau = 1$ . The likelihood for a sample  $X_1, \ldots, X_n$  from an  $\text{Exp}(\theta, 1)$  distribution is  $\prod \exp\{-(x_i - \theta)\}I(x_{(1)} \ge \theta)$ , where  $x_{(1)}$  is the smallest order statistic. The unconstrained MLE is  $\hat{\theta} = X_{(1)}$ .

# 5.1. $H_0: \theta \leq 0$ vs. $H_1: \theta > 0$

The MLE of  $\theta$  is  $\hat{\theta}_0 = \min(X_{(1)}, 0)$  and  $\hat{\theta}_1 = \max(X_{(1)}, 0)$  under  $H_0$  and  $H_1$ , respectively. Given  $X_{(1)}$ , the bootstrap data are independent observations from an  $\operatorname{Exp}(\hat{\theta}_0, 1)$  distribution. 

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

5.1.1. Standard Likelihood Ratio Statistic The standard log-likelihood ratio statistic is  $S = \sum_{i=1}^{n} \log \left| \frac{\exp\{-(X_i - \hat{\theta})\} I(X_{(1)} \ge \hat{\theta})}{\exp\{-(X_i - \hat{\theta}_0)\} I(X_{(1)} \ge \hat{\theta}_0)} \right|$  $= n(\hat{\theta} - \hat{\theta}_0) \\ = \begin{cases} 0, & X_{(1)} \le 0 \\ nX_{(1)}, & X_{(1)} \ge 0. \end{cases}$ Given  $\hat{\theta}_0$ , the bootstrap distribution of S is  $\text{Exp}(n\hat{\theta}_0, 1)$ , left-truncated at 0 with probability mass  $1 - \exp(n\hat{\theta}_0)$  there. 1.  $X_{(1)} \ge 0$ . Then  $S = nX_{(1)}$ ,  $\hat{\theta}_0 = 0$ , the distribution of  $S^*$  is Exp(0,1) with upper- $\alpha$  critical point  $\log(1/\alpha)$ , and the test rejects  $H_0$  if  $nX_{(1)} > -\log \alpha$ . 2.  $X_{(1)} \leq 0$ . Then S = 0,  $\hat{\theta}_0 = X_{(1)}$ , and the distribution of  $S^*$  is  $\text{Exp}(nX_{(1)}, 1)$ , left-truncated at 0 with probability  $1 - \exp(nX_{(1)})$  there. If  $\alpha < \exp(nX_{(1)})$ , the test never rejects  $H_0$ . Otherwise, the test rejects  $H_0$  with probability  $\{\alpha - \exp(nX_{(1)})\}/\{1 - \exp(nX_{(1)})\}.$ Since  $nX_{(1)}$  has an  $\text{Exp}(n\theta, 1)$  distribution,  $P_{\theta}$  {Reject  $H_0$  }  $= P_{\theta} \{ \text{Reject } H_0, X_{(1)} \geq 0 \} + P_{\theta} \{ \text{Reject } H_0, X_{(1)} < 0 \}$  $= P_{\theta}(nX_{(1)} > -\log \alpha, X_{(1)} \ge 0)$  $+ P_{\theta} \{ \text{Reject } H_0, X_{(1)} < 0, \exp(n\theta) < \exp(nX_{(1)}) \le \alpha \} \}$  $= P_{\theta}(nX_{(1)} > -\log \alpha) + E_{\theta} \left[ \frac{\alpha - \exp(nX_{(1)})}{1 - \exp(nX_{(1)})} I(n\theta < nX_{(1)} \le \log \alpha) \right]$  $= \begin{cases} \alpha \exp(n\theta), & \alpha = -c \\ \alpha \exp(n\theta) + \int_{n\theta}^{\log \alpha} \{\alpha - \exp(y)\} \exp(n\theta - y) / \{1 - \exp(y)\} \, dy, & n\theta \le \log \alpha. \end{cases}$ Now for  $n\theta < \log \alpha$ ,  $\int_{n\theta}^{\log \alpha} \frac{\alpha - \exp(y)}{1 - \exp(y)} \exp\{-(y - n\theta)\} \, dy$  $= \exp(n\theta) \int_{\exp(n\theta)}^{\alpha} \frac{\alpha - z}{z^2(1 - z)} dz$  $= \exp(n\theta) \int_{\alpha z = 0}^{\alpha} [\alpha z^{-2} - (1-\alpha) \{z^{-1} + (1-z)^{-1}\}] dz$  $= \exp(n\theta) \left[ -\alpha z^{-1} + (1-\alpha) \{ \log(1-z) - \log z \} \right]_{\exp(n\theta)}^{\alpha}$  $\exp(n\theta)[(1-\alpha)\log(\alpha^{-1}-1) - 1 + \alpha\exp(-n\theta) - (1-\alpha)\log\{\exp(-n\theta) - 1\}].$ = Therefore  $P_{\theta} \{ \text{Reject } H_0 \} = \begin{cases} \alpha \exp(n\theta), & n\theta \ge \log \alpha \\ g_{\alpha}(\exp(n\theta)), & n\theta \le \log \alpha \end{cases}$ where  $g_{\alpha}(z) = \alpha + z(1-\alpha)[\log(\alpha^{-1}-1) - \log(z^{-1}-1) - 1], \quad 0 < z < \alpha.$ Since  $\lim_{z\to 0} g_{\alpha}(z) = \alpha$ ,  $\lim_{z\to \alpha} g_{\alpha}(z) = \alpha^2$ , and  $g''_{\alpha}(z) > 0$  for  $0 < z < \alpha$ , we conclude that  $\sup_{H_0} P_{\theta} \{ \text{Reject } H_0 \} = \lim_{\theta \to -\infty} g_{\alpha}(\exp(n\theta)) = \alpha.$  5.1.2. Cox Likelihood Ratio Statistic

 The Cox log-likelihood ratio statistic is

$$S = \sum_{i=1}^{n} \log \left[ \frac{\exp\{-(X_{i} - \hat{\theta}_{1})\}I(X_{(1)} \ge \hat{\theta}_{1})}{\hat{x}} \right]$$

$$\sum_{i=1}^{2} \exp\left\{-(X_i - \hat{\theta}_0)\right\} I(X_{(1)} \ge \hat{\theta}_0)\right]$$

 $= \int_{-\infty}^{\infty} -\infty, \quad X_{(1)} < 0$ 

$$nX_{(1)}, X_{(1)} \ge 0.$$

It follows that the bootstrap test behaves the same as that based on the standard likelihood ratio. We therefore have the following theorem.

**Theorem 5.1.** For testing  $H_0: \theta \leq 0$  vs.  $H_1: \theta > 0$  for a sample from an  $Exp(\theta, \theta)$ 1) distribution, the bootstrap tests based on the standard and Cox likelihood ratios have size  $\alpha$ .

5.2.  $H_0: \theta \ge 0$  vs.  $H_1: \theta < 0$ 

The MLEs under  $H_0$  and  $H_1$  are  $\hat{\theta}_0 = \max(X_{(1)}, 0)$  and  $\hat{\theta}_1 = \min(X_{(1)}, 0)$ , respectively.

5.2.1. Standard Likelihood Ratio Statistic

**Theorem 5.2.** The bootstrap test of  $H_0$ :  $\theta \ge 0$  vs.  $H_1$ :  $\theta < 0$  based on the standard likelihood ratio test is completely randomized.

*Proof.* The standard log-likelihood ratio statistic is

$$S = \sum_{i=1}^{n} \log \left[ \frac{\exp\{-(X_i - \hat{\theta})\}I(X_{(1)} \ge \hat{\theta})}{((X_i - \hat{\theta}))I(X_i \ge \hat{\theta})} \right]$$

$$\overline{\sum_{i=1}^{i=1}} \left[ \exp\{-(X_i - \theta_0)\} I(X_{(1)} \ge \theta_0) \right]$$
$$\int \infty, \quad X_{(1)} < 0$$

$$= \begin{cases} \infty, & X_{(1)} < 0 \\ 0, & X_{(1)} \ge 0. \end{cases}$$

Since  $\hat{\theta}_0 \geq 0$ , the distribution of  $S^*$  is degenerate at 0. On the other hand, S = 0w.p.1 under  $H_0$ . Therefore the bootstrap test based on S is completely randomized.

5.2.2. Cox Likelihood Ratio Statistic

**Theorem 5.3.** The bootstrap test of  $H_0: \theta \ge 0$  vs.  $H_1: \theta < 0$  based on the Cox likelihood ratio test rejects  $H_0$  w.p.1 for any  $0 < \alpha < 1$ .

*Proof.* The Cox log-likelihood ratio statistic is

$$S = \sum_{i=1}^{n} \log \left[ \frac{\exp\{-(X_i - \hat{\theta}_1)\}I(X_{(1)} \ge \hat{\theta}_1)}{\exp\{-(X_i - \hat{\theta}_0)\}I(X_{(1)} \ge \hat{\theta}_0)} \right]$$

- $= \begin{cases} \infty, & X_{(1)} < 0\\ n(\hat{\theta}_1 \hat{\theta}_0), & X_{(1)} \ge 0 \end{cases}$

$$= \begin{cases} \infty, & X_{(1)} < 0 \\ -nX_{(1)}, & X_{(1)} \ge 0. \end{cases}$$
 50

З 

| <ol> <li>X<sub>(1)</sub> &lt; 0. Then θ̂<sub>0</sub> = 0, the bootstrap data have an Exp(0,1) distribution, and the distribution of S* is the negative of an Exp(0,1) distribution. Since S = ∞, the test rejects H<sub>0</sub> w.p.1 for any 0 &lt; α &lt; 1.</li> <li>X<sub>(1)</sub> ≥ 0. Then θ̂<sub>0</sub> &gt; 0, and the bootstrap data have an Exp(X<sub>(1)</sub>, 1) distribution. The distribution of S* is the negative of an Exp(nX<sub>(1)</sub>, 1) distribution, with support (-∞, -nX<sub>(1)</sub>). Since S = -nX<sub>(1)</sub>, the test rejects H<sub>0</sub> w.p.1 for any 0 &lt; α &lt; 1.</li> </ol>   | 1<br>2<br>3<br>4<br>5<br>6<br>7<br>8                     |
|---|--|
| 6. Conclusion   | 9<br>10  |
| The results show that the size of the bootstrap test of hypotheses is unpredictable.<br>It depends on the problem as well as the choice of test statistic. For example, in the case of testing a normal mean with known variance, the test based on the sample mean or the Cox likelihood ratio is UMP for $0 < \alpha \leq 1/2$ , but it is sub-optimal when it is based on the standard likelihood ratio. On the other hand, if $\alpha > 1/2$ , the test often rejects $H_0$ w.p.1. The overall conclusion is that the size of the test is typically larger than its nominal level. This may explain the high power that the test is found to possess in simulation experiments.   | 11<br>12<br>13<br>14<br>15<br>16<br>17<br>18<br>19<br>20 |
| Appendix  | 21   |
| Proof of Lemma 2.1.   | 23<br>24   |
| First note that   | 25   |
| (6.1) $\phi(x) - \phi(x - \theta) \begin{cases} > 0, & \text{if } x < \theta/2 \\ < 0, & \text{if } x > \theta/2. \end{cases}$  | 20<br>27<br>28   |
| Let $f(\theta)$ denote the function (2.2). We consider two cases.   | 30   |
| 1. $\alpha \ge 1/2$ . Since $z_{\alpha}^+ = 0$ , we have $f(\theta) = 1 - (1 - \alpha) \int_0^\infty \phi(x - \theta) / \Phi(x)  dx$ and  | 31<br>32   |
| $f(0) - f(\theta) \qquad \int_{-\infty}^{\infty} \phi(x - \theta) - \phi(x)$  | 33   |
| $\frac{1}{1-\alpha} = \int_0^{\infty} \frac{1}{\Phi(x)} dx$   | 34<br>35   |
| $\int^{\theta/2} \phi(x-\theta) - \phi(x) dx + \int^{\infty} \phi(x-\theta) - \phi(x) dx$   | 36   |
| $= \int_0 -\frac{\Phi(x)}{\Phi(x)} dx + \int_{\theta/2} -\frac{\Phi(x)}{\Phi(x)} dx$  | 37   |
| $\sum_{n=0}^{\infty} \int_{-\infty}^{\theta/2} \left( \int_{-\infty}^{\infty} \left( $ | 38<br>39   |
| $\geq 2\int_{0} \{\phi(x-\theta) - \phi(x)\} dx + \int_{\theta/2} \{\phi(x-\theta) - \phi(x)\} dx$  | 40   |
| $= 2\{\Phi(-\theta/2) - \Phi(-\theta) - \Phi(\theta/2) + 1/2\} - \Phi(-\theta/2) + \Phi(\theta/2)$  | 41   |
| $= 2\{\Phi(\theta) - \Phi(\theta/2)\}$  | 42   |
| $\geq 0$  | 43<br>44   |
|   | 45   |
| where we use $(0.1)$ in the first inequality. Hence   | 46   |
| $f(	heta) \ \leq \ f(0)$  | 47   |
| $-1-(1-\alpha)\int_{-\infty}^{\infty}\phi(x)/\Phi(x)dx$   | 48<br>49   |
| $= 1  (1  \alpha) \int_0^{\infty} \psi(x) / x(x)  \alpha x$   | 50   |
| $= 1 + (1 - \alpha) \log(1/2).$   | 51   |

W. Loh and W. Zheng

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

Bootstrap tests

2.  $\alpha < 1/2$ . Write  $f(\theta) = J_1(\theta) + J_2(\theta)$ , where  $J_1(\theta) = \alpha (1-\alpha)^{-1} \Phi(\theta - z_\alpha) + \Phi(-z_\alpha - \theta)$ З  $J_2(\theta) = E\left\{\frac{(1-2\alpha)\Phi(Z+\theta) - (1-\alpha)^2}{(1-\alpha)\Phi(Z+\theta)}I(Z+\theta \ge z_\alpha)\right\}.$ Since  $\partial J_1(\theta)/\partial \theta = \alpha (1-\alpha)^{-1} \phi(\theta-z_\alpha) - \phi(\theta+z_\alpha)$  $= \phi(\theta + z_{\alpha}) \{ \alpha (1 - \alpha)^{-1} \exp(2\theta z_{\alpha}) - 1 \},\$  $J_1(\theta)$  is decreasing-increasing, with minimum at  $\theta_0 = (2z_\alpha)^{-1} \log\{(1-\alpha)/\alpha\} >$ 0. Further,  $J_1(0) = \lim_{\theta \to \infty} J_1(\theta) = \alpha (1-\alpha)^{-1}$ . Therefore  $J_1(\theta) < J_1(0)$  for  $\theta > 0.$ To obtain a similar result for  $J_2$ , let  $q(x) = \{(1 - 2\alpha)\Phi(x) - (1 - \alpha)^2\} / \{(1 - \alpha)\Phi(x)\}$ which is increasing in x with  $g(z_{\alpha}) = -\alpha/(1-\alpha)$  and  $g(x) \to -\alpha^2/(1-\alpha)$ as  $x \to \infty$ . Hence g(x) < 0 for  $x \ge z_{\alpha}$ (a)  $0 < \theta \leq 2z_{\alpha}$ . Since  $\phi(x) \leq \phi(x-\theta)$  for  $x \geq z_{\alpha}$ ,  $J_2(0) - J_2(\theta) = \int_0^\infty g(x)\phi(x) \, dx - \int_0^\infty g(x)\phi(x-\theta) \, dx \ge 0.$ (b)  $\theta > 2z_{\alpha}$ . From (6.1),  $J_2(0) - J_2(\theta)$  $= \int_{x}^{\infty} g(x)\phi(x)\,dx - \int_{x}^{\infty} g(x)\phi(x-\theta)\,dx$  $= \int_{-\infty}^{\theta/2} g(x) [\phi(x) - \phi(x-\theta)] dx + \int_{\theta/2}^{\infty} g(x) [\phi(x) - \phi(x-\theta)] dx$  $> -\frac{\alpha}{1-\alpha} \int_{-\infty}^{\theta/2} [\phi(x) - \phi(x-\theta)] \, dx - \frac{\alpha^2}{1-\alpha} \int_{\theta/2}^{\infty} [\phi(x) - \phi(x-\theta)] \, dx$  $= \alpha (1-\alpha)^{-1} [-\{\Phi(\theta/2) - (1-\alpha) - \Phi(-\theta/2) + \Phi(z_{\alpha} - \theta)\}]$  $+ \alpha \{ \Phi(\theta/2) - \Phi(-\theta/2) \} ]$  $= \alpha(1-\alpha)^{-1} \{K_1(\theta) + K_2(\theta)\}$ where  $K_1(\theta) = \Phi(-\theta/2) - \Phi(z_{\alpha} - \theta)$  $K_2(\theta) = 1 - \alpha - \Phi(\theta/2) + \alpha \Phi(\theta/2) - \alpha \Phi(-\theta/2).$ Now  $K_1(\theta) > 0$  for  $\theta > 2z_{\alpha}$ ,  $K'_2(\theta) = (\alpha - 1/2)\phi(\theta/2) < 0$ , and  $K_2(\theta) \rightarrow 0$ 0 as  $\theta \to \infty$ . Thus  $J_2(0) - J_2(\bar{\theta}) \ge 0$ . Therefore  $f(\theta) \leq f(0) = 2\Phi(-z_{\alpha}) - (1-\alpha) \int_{z_{\alpha}}^{\infty} \phi(x) / \Phi(x) \, dx = 2\alpha + (1-\alpha) \int_{z_{\alpha}}^{\infty} \phi(x) / \Phi(x) \, dx$  $\alpha$ ) log $(1 - \alpha)$ . It remains to show that  $f(0) > \alpha$  for all  $0 < \alpha < 1$ . Let  $h(\alpha) = f(0) - \alpha$ . Then  $h(\alpha) = \left\{ \begin{array}{ll} \alpha + (1-\alpha)\log(1-\alpha), & \text{if } 0 < \alpha \leq 1/2 \\ 1 - \alpha - (1-\alpha)\log 2, & \text{if } 1/2 \leq \alpha < 1 \end{array} \right.$ 

imsart-coll ver. 2008/08/29 file: Loh.tex date: April 10, 2009

W. Loh and W. Zheng

and h is continuous with h(0) = h(1) = 0,  $h(1/2) = (1 - \log 2)/2 > 0$ , and  $h'(\alpha) = \begin{cases} -\log(1-\alpha) > 0, & \text{if } 0 < \alpha \le 1/2 \\ -1 + \log 2 < 0, & \text{if } 1/2 \le \alpha < 1. \end{cases}$ Therefore  $h(\alpha) > 0$  for  $0 < \alpha < 1$ , concluding the proof. References [1] ANÉ, C., BURLEIGH, J. G., MCMAHON, M. M. and SANDERSON, M. J. (2005). Covarion structure in plastid genome evolution: A new statistical test. Molecular Biology and Evolution, 22, 914-924. CARSTENS B. C., BANKHEAD III, A. AND JOYCE P. and SULLIVAN J. (2005). Testing population genetic structure using parametric bootstrapping and MIGRATE-N. Genetica, 124, 71-75. [3] CHUNG P. J. and MOURA J. (2004). A GLRT and bootstrap approach to detection in magnetic resonance force microscopy. In IEEE International Conference on Acoustics, Speech, and Signal Processing. Cox, D. R. (1961). Tests of separate families of hypotheses. In Proceedings of the Fourth Berkeley [4]Symposium on Mathematical Statistics and Probability, 1. Berkeley, CA, University of California Press. 105-123. DALLA V. and HIDALGO J. (2004). A parametric bootstrap test for cycles. Journal of Econometrics, 129, 219-261. [6] DAMEUSA A., RICHTER, F. G. C., BRORSEN, B. W. and SUKHDIAL K. P. (2002). AIDS versus the Rotterdam Demand System: a Cox Test with Parametric Bootstrap. Journal of Agricultural and Resource Economics, 27, 335–347. EFRON, B. (1979). Bootstrap methods: another look at the jackknife. Ann. Statist., 7, 1–26. HUNSBERGER S., ALBERT, P. S., FOLLMANN, D. A. and SUH E. (2002). Parametric and semiparametric approaches to testing for seasonal trend in serial count data. Biostatistics, 3, 289-298. KAITIBIE, S., NGANJE, W. E., WADE BRORSEN, B. and EPPLIN, F. M. (2007). A Cox Parametric Boot-strap Test of the von Liebig Hypotheses. Canadian Journal of Agricultural Economics, 55, 15–25. [10] LEHMANN E. L. and ROMANO, J. P. (2005). Testing Statistical Hypotheses, 3rd ed. Springer, New York. [11] SOLOW A. R., COSTELLO, C. J. and WARD, M. (2003). Testing the Power Law Model for Discrete Size Data. The American Naturalist, 162, 685-689. [12] WALTERS S. J. and CAMPBELL, M. J. (2004). The use of bootstrap methods for analysing health-related quality of life outcomes (particularly the SF-36). Health Qual Life Outcomes, 2, 70. [13] WILSON, E. B. and HILFERTY, M. M. (1931). The distribution of chi-square. Proceedings of the National Academy of Sciences, 17, 684–688. 

|               | Nonparametric Estimation for Lévy  |
|---------------|--|
|               | Models Based on Discrete-Sampling  |
|               | José E. Figueroa-López <sup>1,*</sup>  |
|               | Purdue University  |
|               | Abstract: A Lévy model combines a Brownian motion with drift and a pure-<br>jump homogeneous process such as a compound Poisson process. The estima-<br>tion of the Lévy density, the infinite-dimensional parameter controlling the<br>jump dynamics of the process, is studied under a discrete-sampling scheme.<br>In that case, the jumps are latent variables whose statistical properties can<br>in principle be assessed when the frequency of observations increase to in-<br>finity. We propose nonparametric estimators for the Lévy density following<br><i>Grenander's method of sieves</i> . The associated problem of selecting a suitable<br>approximating sieve is subsequently investigated using regular piece-wise poly-<br>nomials as sieves and assuming standard smoothness conditions on the Lévy<br>density. By sampling the process at a high enough frequency relative to the<br>time horizon $T$ , we show that it is feasible to choose the dimension of the<br>sieve so that the rate of convergence of the <i>risk of estimation off the origin</i><br>is the best possible from a minimax point of view, and even if the estimation<br>were based on the whole sample path of the process. The sampling frequency<br>necessary to attain the optimal minimax rate is explicitly identified. The pro-<br>posed method is illustrated by simulation experiments in the case of variance<br>Gamma processes. |
| $\mathbf{C}$  | ontents  |
| 1<br>2<br>3   | Introduction111.1Motivation and Some Background111.2The Statistical Problems and Methodology111.3An Overview of the Estimators and the Results121.4Outline12The Estimators and Central Limit Theorems12The Model Selection Problem123.1Analysis of the Variance Term123.2The Approximation Error for Besov Type Smooth Functions123.3Rate of Convergence for Smooth Functions Via Splines133.4About the Critical Mesh13  |
| $4 \\ 5 \\ 6$ | Minimax Risk of Estimation for Smooth Lévy Densities13A Data-Driven Selection Method and Adaptability13An Example: Estimation of Variance Gamma Processes13  |
| st            | 6.1 The Model  |

imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009

| 6.3<br>Append<br>Append | The Simulation Procedure       135         The Numerical Results       135         ix A: Technical Proofs       137         ix B: Figures       142 | 1<br>2<br>3<br>4 |
|-------------------------|---|------------------|
| Referen                 | oduction  | (                |
| 1.1. N                  | otivation and Some Background   | 2                |
|                         |   | 1                |

$$(1.1) S_t := S_0 e^{X_t},$$

where  $X := \{X_t\}_{t \ge 0}$  is a Lévy process. Even this simple extension of the classical Black-Scholes model, in which X is simply a Brownian motion with drift, is able to account for some fundamental empirical features commonly observed in time series of asset returns such as heavy tails, high-kurtosis, and asymmetry. More re-cently, other Lévy based models have been proposed to account for more stylized features of stock prices. These models include exponential time-changed Lévy processes (cf. [7]-[9]), and stochastic differential equations driven by multivariate Lévy processes (cf. [1], [31]). Lévy processes, as models capturing some of the most impor-tant features of returns and as "first-order approximations" to other more accurate models, should be considered first in developing and testing a successful statistical methodology. However, even in such parsimonious models, there are several issues in performing statistical inference by standard likelihood-based methods. 

A Lévy process is the "discontinuous sibling" of a Brownian motion. Concretely,  $X = \{X_t\}_{t\geq 0}$  is a Lévy process if X has independent and stationary increments, its paths are right-continuous with left limits, and it has no fixed jump times. The later condition means that, for any t > 0,

$$\mathbb{P}\left[\Delta X_t \neq 0\right] = 0,$$
<sup>37</sup>

where  $\Delta X_t := X(t) - \lim_{s \geq t} X_s$  is the magnitude of the "jump" of X at time t. It can be proved that the only Lévy process with continuous paths is essentially the Brownian motion  $W := \{W_t\}_{t\geq 0}$  up to a drift term bt (hence, the well-known Gaussian distribution of the increments of W is a byproduct of the stationarity and independence of its increments). The only deterministic Lévy process is of the form  $X_t := bt$ , for a constant b. Another distinguished type of Lévy process is a compound Poisson process defined as

where N is a homogeneous Poisson process and the random variables  $\xi_i$ ,  $i \ge 1$ , are mutually independent from one another, independent from N, and with common

distribution  $\rho$ . The process N dictates the jump times, which can occur "homoge-neously" across time with an (average) intensity of  $\lambda$  jumps per unit time, while the sequence  $\{\xi_i\}_{i>1}$  determines the sizes of the jumps.

It turns out that the most general Lévy process is the superposition of a Brownian motion with drift,  $\sigma W_t + bt$ , a compound Poisson process, and the limit process resulting from making the jump intensity of a compensated compound Poisson process,  $Y_t - \mathbb{E} Y_t$ , to go to infinity while simultaneously allowing jumps of smaller sizes. The latter limiting process is governed by a measure  $\nu$  such that the intensity of jumps is  $\lambda_{\varepsilon} := \nu(\varepsilon \leq |x| < 1)$ , the common distribution of the jump sizes is  $\rho_{\varepsilon}(dx) := \mathbf{1}_{\{|x| \ge \varepsilon\}} \nu(dx) / \lambda_{\varepsilon}$ , and the limit is when  $\varepsilon \searrow 0$ . For such a limit to converge to a "steady" process it must hold that

$$\int_{\{|x|<1\}} x^2 \nu(dx) < \infty.$$

The previous fundamental decomposition of a Lévy process is called the Lévy-Itô decomposition (see Section 19 in [29] for the details).

J

In summary, Lévy processes are determined by three "parameters": a nonnegative real  $\sigma^2$ , a real b, and a measure  $\nu$  on  $\mathbb{R}\setminus\{0\}$  such that  $\int (x^2 \wedge 1)\nu(dx) < \infty$ . The measure  $\nu$  controls the jump dynamics of the process X in that for any  $A \in \mathcal{B}(\mathbb{R})$  whose indicator  $\chi_A$  vanishes in a neighborhood of the origin,

$$\nu(A) = \frac{1}{t} \mathbb{E}\left[\sum_{s \le t} \chi_A\left(\Delta X(s)\right)\right],$$

for any t > 0 (see Section 19 of [29]). Thus,  $\nu(A)$  gives the average number of jumps (per unit time) whose magnitudes fall in the set A. A common assumption in Lévybased financial models is that  $\nu$  is determined by a function  $s : \mathbb{R} \setminus \{0\} \to [0, \infty)$ , called the *Lévy density*, as follows

 $\nu(A) = \int_A s(x) dx, \ \forall A \in \mathcal{B}(\mathbb{R} \setminus \{0\}).$ 

Intuitively, the value of s at  $x_0$  provides information on the frequency of jumps with sizes "close" to  $x_0$ . In the case of the compound Poisson process (1.2), the Lévy measure is  $\nu(dx) = \lambda \rho(dx)$ . 

By allowing a general Lévy process X in (1.1), instead of just a Brownian motion with drift as in the Black-Scholes model, one can incorporate two very appealing features: sudden changes in the price dynamics and some freedom in the distribution for the log return  $\log\{S_t/S_s\} = X_t - X_s$ . The possible distributions belong to the socalled class of infinitely-divisible distributions, a very rich class which include most known parametric families of distributions. We recall that an infinitely divisible distribution  $\mu$  is characterized by the so-called *Lévy-Khinchin representation* of its characteristic function.

There are two key properties of a Lévy process that are exploited in this work. The first property relates  $\nu$  with the short-term moments of  $X_t$ . Concretely, if  $\varphi$  is  $\nu$ -continuous, bounded, and vanishing in a neighborhood of the origin, then

(1.3) 
$$\lim_{\Delta \to 0} \frac{1}{\Delta} \mathbb{E} \varphi(X_{\Delta}) = \int \varphi(x) \nu(dx) = \int \varphi(x) s(x) dx.$$

The limiting relation (1.3) is straightforward when X is a compound Poisson process. A proof of (1.3) for a general Lévy process can be found in [29] (see his З

Corollary 8.9). Let us remark that (1.3) is also valid for certain unbounded functions  $\varphi$ , which does not necessarily vanish in a neighborhood of the origin, but rather converge to 0 at a proper rate (see [15] for more details). The second key property is related to the decomposition of X into two *independent* processes: one accounting for the "small" jumps and a compound Poisson process collecting the "big" jumps. Concretely, let

$$\widetilde{X}_t^{\varepsilon} := \sum_{s < t} \Delta X_s \mathbf{1}_{\{|\Delta X_s| > \varepsilon\}},$$

be the piece-wise constant process associated with those jumps of X with sizes larger than  $\varepsilon$ . Then,  $\tilde{X}^{\varepsilon}$  is a compound Poisson process independent of  $X - \tilde{X}^{\varepsilon}$ .

## 1.2. The Statistical Problems and Methodology

We are interested in estimating the Lévy density s on a window of estimation  $D := [a, b] \subset \mathbb{R} \setminus \{0\}$ , based on discrete observations of the process on a finite interval [0,T]. We remark that the domain D is "separated" from the origin; that is to say, the estimation window D lies outside of a neighborhood of the origin. If the whole path of the process were available (and hence, the jumps of the process would be available), the problem would be identical to the estimation of the intensity of a non-homogeneous Poisson process on a fixed time interval, say [0, 1], based on T independent copies of the process. However, under discrete-sampling, the times and sizes of jumps are latent (unobservable) variables, whose statistical properties can be assessed when the frequency of observations increase to infinity at a certain speed relative to the time horizon. Hence, we will aim at determining the performance of our estimation method as both frequency and time horizon increase. 

We adopt the so-called *method of sieves* originally proposed by [18] and im-plemented by Birgé, Massart, and others (see e.g. [3] & [5]) in several classical nonparametric problems such as density estimation and regression. This approach consists of the following general steps. First, choose a family of finite-dimensional *linear models* of functions, called *sieves*, with good approximation properties. Com-mon sieves are splines, trigonometric polynomials, or wavelets. Second, specify a distance between functions relative to which the best approximation to s, in a given linear model, is going to be defined and characterized. Finally, devise an estimator, called the *projection estimator*, for the best approximation of s in the given linear model. It is important to point out that in principle there is no guarantee that the projection estimator will be nonnegative. In practice, one barely faces this prob-lem when working with a large sample size, which is exactly the situation when nonparametric methods are recommended. 

A linear model has the generic form

(1.4) 
$$\mathcal{S} := \{\beta_1 \varphi_1 + \dots + \beta_d \varphi_d : \beta_1, \dots, \beta_d \in \mathbb{R}\},\$$

where  $\varphi_1, \ldots, \varphi_d$  are given functions, typically taken to be orthonormal with respect to the inner product  $\langle p, q \rangle := \int_D p(x)q(x)dx$ . In the sequel,  $\|\cdot\|$  stands for the associated norm  $\langle \cdot, \cdot \rangle^{1/2}$  on  $\mathbb{L}^2(D, dx)$ . Relative to the distance induced by  $\|\cdot\|$ , the element of  $\mathcal{S}$  closest to s, i.e. the *orthogonal projection* of s on  $\mathcal{S}$ , is

$$\frac{d}{2}$$

50 (1.5) 
$$s^{\perp}(x) := \sum_{j=1}^{50} \nu(\varphi_j) \varphi_j(x),$$
 51 51

imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009

where  $\nu(\varphi_j) := \langle \varphi_j, s \rangle = \int \varphi_j(x) s(x) dx$ . Then, the method of sieves boils down to estimate the orthogonal projection (1.5) on an "adequate" sieve S. The core problem in this paper is to determine what a good sieve is. A very large linear model  $\mathcal{S}$  will allow to attain a close approximation to s, but will entail necessarily a high estimation variance as the result of the large number of coefficients  $\beta_i$  to be estimated. Therefore, an essential task, called *model selection*, consists of selecting a linear model  $\mathcal{S}$  accomplishing a good tradeoff between the error of approximation (or mis-specification error) and the standard error of the estimation. Concretely, one wishes to minimize the risk of the estimator  $\hat{s}$ , which in turn can be decomposed into two antagonist terms as follows: 

(1.6) 
$$\mathbb{E} \|s - \hat{s}\|^2 = \|s - s^{\perp}\|^2 + \mathbb{E} \|s^{\perp} - \hat{s}\|^2.$$

The first term, called the *bias term*, accounts for the error of the approximation, while the second, called the *variance term*, accounts for the standard error of the estimation.

## 1.3. An Overview of the Estimators and the Results

We assume that the Lévy process  $\{X_t\}_{t\geq 0}$  is being sampled over a time horizon [0,T] at discrete times  $0 < t_1^n < \cdots < t_n^n = T$ . In the sequel,  $t_0^n := 0, \pi^n := \{t_k^n\}_{k=0}^n$ , and  $\bar{\pi}^n := \max_k \{t_k^n - t_{k-1}^n\}$ , the so-called mesh of the partition. We shall sometimes drop the superscript n in  $\pi^n$  and  $t_i^n$ . The following statistics are the main building blocks of our estimators:

(1.7) 
$$\hat{\beta}^{\pi^{n}}(\varphi) := \frac{1}{t_{n}} \sum_{k=1}^{n} \varphi \left( X_{t_{k}^{n}} - X_{t_{k-1}^{n}} \right).$$

In the case of a quadratic function  $\varphi(x) = x^2$ ,  $\sum_{k=1}^n \varphi\left(X_{t_k^n} - X_{t_{k-1}^n}\right)$  is called the realized quadratic variation (or variance) of the process. Thus, the statistics (1.7) can be interpreted as the average realized  $\varphi$ -variation of the process per unit time based on the observations  $X_{t_1^n}, \ldots, X_{t_n^n}$ .

To explain the motivation behind the estimator in (1.7), let us assume for now that sampling observations are equally-spaced in time so that  $\Delta_n := t_i^n - t_{i-1}^n = T/n$  for all *i*, and hence,

$$\mathbb{E}\left\{\hat{\beta}^{\pi^{n}}(\varphi)\right\} = \frac{1}{\Delta_{n}} \mathbb{E}\varphi\left(X_{\Delta_{n}}\right),$$
<sup>36</sup>
<sub>37</sub>
<sub>37</sub>

$$\operatorname{Var}\left\{\hat{\beta}^{\pi^{n}}(\varphi)\right\} = \frac{1}{T} \left(\frac{1}{\Delta_{n}} \mathbb{E} \varphi^{2}\left(X_{\Delta_{n}}\right)\right) - \frac{1}{n} \left(\frac{1}{\Delta_{n}} \mathbb{E} \varphi\left(X_{\Delta_{n}}\right)\right)^{2}.$$

In view of (1.3), it is now evident that

(1.8) 
$$\lim_{n \to \infty} \mathbb{E}\left\{\hat{\beta}^{\pi^{n}}(\varphi)\right\} = \int \varphi(x)s(x)dx, \text{ and } \lim_{T \to \infty} \sup_{n} \operatorname{Var}\left(\hat{\beta}^{\pi^{n}}(\varphi)\right) = 0,$$

if  $\varphi$  is  $\nu$ -continuous, bounded, and vanishing in a neighborhood of the origin. In statistical terms, (1.8) means that the statistic  $\hat{\beta}^{\pi^n}(\varphi)$  is an asymptotically unbiased estimator of  $\int \varphi(x) s(x) dx$  with associated risk vanishing uniformly when time horizon T increases. The previous argument leads us to propose

49  
50 (1.0) 
$$\hat{s}^{\pi^{n}}(x) := \sum_{n=0}^{d} \hat{\beta}^{\pi^{n}}(x) \cdot (x)$$
50

50 (1.9) 
$$\hat{s}^{\pi^n}(x) := \sum_{j=1} \hat{\beta}^{\pi^n}(\varphi_j)\varphi_j(x),$$
 51 51

imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009
as a natural estimator for the orthogonal projection  $s^{\perp}$  defined in (1.5). In view of (1.8),  $\hat{s}^{\pi^n}$  is a "consistent" estimator for  $s^{\perp}$ , in the integrated mean-square sense, as both the time horizon  $T = t_n^n$  and the sampling frequency  $n/t_n^n$  go to  $\infty$ . The general sampling case will be considered in Section 2 as well as other statistical properties. It is worth pointing out that  $\hat{s}^{\pi^n}$  is independent of the specific orthonormal basis of  $\mathcal{S}$  as it can be proved that  $\hat{s}^{\pi^n}$  is the unique solution of the minimization problem

$$\min_{f\in\mathcal{S}}\gamma_{D}^{\pi}(f),$$

where  $\gamma_{D}^{\pi^{n}}: L^{2}\left(D, dx\right) \to \mathbb{R}$  is given by

(1.10) 
$$\gamma_{D}^{\pi^{n}}(f) \equiv -\frac{2}{t_{n}^{n}} \sum_{k=1}^{n} f(X_{t_{k}^{n}} - X_{t_{k-1}^{n}}) + \int f^{2}(x) dx.$$

In the literature of model selection (see e.g. [4] and [25]),  $\gamma_D^{\pi^n}$  is called the *contrast* function.

Finding the best sieve S to estimate s, even if we stick with using the class of projection estimators in (1.9), is impossible because s is unknown. However, it is possible to select a reasonably good model under certain qualitative assumptions on the parameter s, typically expressed by requiring s to be a member of a certain class  $\Theta$  of smooth functions. Concretely, suppose we are interested in selecting a good model out of a family of linear models  $\{S_m\}_{m \in \mathcal{M}}$  (here,  $\mathcal{M}$  is a suitable set of labels). Let  $m^* := m^*(\pi)$  be the optimal minimax element of  $\{\hat{s}_m\}_{m \in \mathcal{M}}$  on  $\Theta$ , defined as

$$m^* := \operatorname{arginf}_{m \in \mathcal{M}} \{ \sup_{s \in \mathcal{O}} \mathbb{E} \, \| s - \hat{s}_m \|^2 \}.$$

By requiring certain conditions on  $\Theta$  and by choosing a suitable family of sieves  $\{S_m\}_{m \in \mathcal{M}}$ , we can ensure that

(1.11) 
$$\mathbb{E} \| s - \hat{s}_{m^*(\pi)} \|^2 \to 0$$

as the mesh of the partition  $\pi = \{t_k\}_{k\geq 1}$  vanishes and the time horizon  $T := t_n$  goes to infinity. Our goal will be to select a linear model  $\hat{m}(\pi) \in \mathcal{M}$  so that the projection estimator on this model,  $\hat{s}_{\hat{m}(\pi)}$ , "attains" the minimax rate of convergence in (1.11), in the sense that

(1.12) 
$$\limsup \frac{\mathbb{E} \|s - \hat{s}_{\hat{m}(\pi)}\|^2}{\mathbb{E} \|s - \hat{s}_{m^*(\pi)}\|^2} < \infty,$$

where the limit is taken as  $\bar{\pi} \to 0$  and  $T \to \infty$ . In order to be able to determine in a "simple" way the rate of convergence of  $\hat{s}_{\hat{m}(\pi)}$ , we shall control the sampling frequency, measured by  $\bar{\pi}$ , in function of the time horizon T. In the case of a finitejump activity process with a jump intensity of  $\lambda$  jumps per unit time, we can expect that it suffices to sample at a faster rate than 1/T (that is,  $\bar{\pi}T \to 0$ ). It is intuitive that in general the sampling frequency will depends on how close the window of estimation D is to the origin (see Section 3.4). The limit result (1.12) and the rate of convergence of projection estimators for a certain class of smooth Lévy densities are addressed in Section 3.

In this paper, we will show that the rate of convergence that can be attained
using projection estimation on sieves is actually the best possible among all feasible
estimators, given the information available on s (namely, that s belongs to a certain

З

class  $\Theta$  of smooth functions), and even if the estimators were based on continuoustime sampling of the process. Concretely, define  $\hat{s}_{\tau}^{*}$  be the minimax estimator,

$$\hat{s}_{_{T}}^{*} := \arg \inf_{\hat{s}} \sup_{s \in \Theta} \mathbb{E} \|s - \hat{s}\|^{2} < \infty,$$

where the infimum is over all the estimators  $\hat{s}$  of s based on  $\{X(t)\}_{0 \le t \le T}$ . Then, by sampling at a high enough frequency (relative to T), we can accomplish that

$$\limsup_{T \to \infty} \frac{\mathbb{E} \|s - \hat{s}_{\hat{m}}\|^2}{\mathbb{E} \|s - \hat{s}_T^*\|^2} < \infty.$$

The rate of convergence of the minimax estimator will be provided in Section 4.

Let us finish by pointing out that the model selection problem was already analyzed in Figueroa-López & Houdré (2006) using the statistics

which intrinsically required continuous-time sampling of the process to determine the jumps  $\Delta X_t$ . In the cited paper, the statistics (1.7) were proposed as good proxies of (1.13). Indeed, convergence in distribution is not hard to check, but moreover, recently [20] prove that (1.7) converges in probability to (1.13) when  $n \to \infty$  (for fixed T). To the best of our knowledge, an analysis of the model selection problem for Lévy densities, under discrete sampling schemes, has not been considered before the present work.

## 

### 1.4. Outline

The paper is structured as follows. In Section 2, we introduce the estimators proposed in this paper and study some basic statistical properties. In particular, we prove a CLT for the estimator  $\hat{\beta}^{\pi}(\varphi)$  of (1.7) centered at the inner product  $\beta(\varphi) = \int \varphi(x) s(x) dx$ . In Section 3, we describe how to control the risk of the projection estimators by imposing three conditions. First, the time horizon T should be large enough (compared to the complexity of the sieves). Second, the time span between consecutive observations should be small enough compared to the time horizon. Finally, the sieves should have good approximating properties in general classes of smooth functions. We show that by ensuring the three previous conditions and by suitably choosing the dimension of the sieve (in terms of the presumed smoothness of the function s), the rate of convergence of the risk is of order  $O(T^{-2\alpha/(2\alpha+1)})$  as  $T \to \infty$  provided that the parameter s has "degree of smoothness"  $\alpha$ .

In Section 4, the minimax risk of estimation, defined by

$$\inf_{\hat{s}_T} \sup_{s \in \Theta} \mathbb{E}_s \|s - \hat{s}_T\|^2,$$

is studied. Here, the infimum is over all estimators  $\hat{s}_{\tau}$  which can be computed from the whole sample paths of X on the interval [0, T] and the supremum is over all Lévy densities in a class  $\Theta$  of functions that are smooth in D = [a, b]. We found that the minimax risk converges at an order of  $O(T^{-2\alpha/(2\alpha+1)})$ , where  $\alpha$  is a parameter that measures the smoothness of the functions on  $\Theta$ . For instance, if s has d continuous derivatives in D, then  $\alpha > d$ . The rate of convergence of the estimation is faster when  $\alpha$  increases. Sections 3 and 4 justify the claim of the abstract: "...we show that 

З

it is feasible to choose the dimension of the sieve so that the rate of convergence of the *risk of estimation off the origin* is the best possible from a minimax point of view, and even if the estimation were based on the whole sample path of the process".

In Section 5, we propose a data-driven selection method for the sieve. Instead of deciding the dimension of the sieve from a presumed degree of smoothness of s(as it was suggested in Section 3), we propose to choose the sieve that minimizes an unbiased estimator of the risk of the projection estimator corresponding to that sieve. Since the proposed estimator of the risk will require the knowledge of all jumps of X up to time T, we replace it by a natural discrete-based proxy, where the jumps  $\Delta X_t$  are replaced by the increments  $X_{t_k} - X_{t_{k-1}}$ . Section 6 illustrates the statistical methods using simulation experiments in the case of a variance gamma Lévy model. We finish with an Appendix where some technical proofs are given.

# 2. The Estimators and Central Limit Theorems

In this section, our goal is to survey some statistical properties of the estimators (1.7) and (1.9). We already mentioned a few of these in the case of *regular sampling*<sup>1</sup> and of bounded  $\nu$ -continuous *test functions*  $\varphi$  vanishing in a neighborhood of the origin. In the framework of this paper, this kind of test functions indeed suffices to recover and estimate the Lévy density off the origin. Our first result is a simple application of the Central Limit Theorem (CLT) for independent random variables (cf. [20] [Theorem 3.2] for the case of regular sampling). In the following results, Z stands for a standard Normal random variable.

**Proposition 2.1.** Let  $\varphi$  be  $\nu$ -continuous, bounded, and such that  $\varphi(x) = o(|x|)$ , as  $x \to 0$ . Then,

(2.1) 
$$\sqrt{t_n} \left( \hat{\beta}^{\pi}(\varphi) - \mathbb{E} \, \hat{\beta}^{\pi}(\varphi) \right) \xrightarrow{\mathfrak{D}} \nu(\varphi^2)^{\frac{1}{2}} Z,$$

as  $t_n \to \infty$  and  $\bar{\pi} \to 0$ .

*Proof.* Let  $\Gamma_t(\varphi) := \mathbb{E} \varphi^2(X_t) - \{\mathbb{E} \varphi(X_t)\}^2$  and  $\Delta_k := t_k - t_{k-1}$ . We can write

$$\sqrt{t_n} \left( \hat{\beta}^{\pi}(\varphi) - \mathbb{E} \, \hat{\beta}^{\pi}(\varphi) \right) = \sum_{k=1}^n \xi_k^{\pi}, \qquad 34$$

where  $\xi_k^{\pi} = \frac{1}{\sqrt{t_n}} \left\{ \varphi \left( X_{t_k} - X_{t_{k-1}} \right) - \mathbb{E} \varphi \left( X_{t_k - t_{k-1}} \right) \right\}$ . Under the assumption of this Proposition, it turns out that  $\lim_{t\to 0} \frac{1}{t} \Gamma_t(\varphi) = \nu(\varphi^2)$  (see Lemma 5.5 in Jacod (2007)), and thus,

$$t_{n,\pi} := \operatorname{var}_{k=1} \zeta_k = t_n \sum_{k=1}^{r} \Delta_k(\varphi) \to \nu(\varphi),$$

as the mesh  $\bar{\pi} := \max_k \{t_k - t_{k-1}\} \to 0$ . Due to the boundedness of  $\varphi$ , we have that, for  $\bar{\pi}$  small enough,

$$\frac{\xi_k^{\pi}|}{\bar{\tau}_{n,\pi}} \le C \frac{1}{\sqrt{t_n}} \to 0,$$

as  $t_n \to \infty$ . Then, (2.1) follows from the Central Limit Theorem for independent random variables (see e.g. the Corollary following Theorem 7.1.2 in [10]).  $\bar{\sigma}^2$ 

<sup>&</sup>lt;sup>1</sup>Sampling equally spaced in time.

In order to provide an explicit centering in (2.1), we need to estimate the rate of convergence of the bias  $\mathbb{E}\hat{\beta}^{\pi}(\varphi) - \nu(\varphi)$ . Since

$$\mathbb{E}\,\hat{\beta}^{\pi}(\varphi) - \nu(\varphi) = \frac{1}{t_n} \sum_{k=1}^n \Delta_k \left\{ \frac{1}{\Delta_k} \mathbb{E}\,\varphi(X_{\Delta_k}) - \nu(\varphi) \right\},\,$$

the problem is equivalent to analyzing the rate of convergence in (1.3). To achieve this goal, we need to impose some regularity on either the Lévy process or the moment functions  $\varphi$ . Following the second approach, [15] shed light on this problem for functions  $\varphi \in C_b^2(\mathbb{R})$ ; namely, twice-continuously differentiable functions  $\varphi$  such that  $\limsup_{|x|\to\infty} |\varphi^{(i)}(x)| < \infty$ , for i = 0, 1, 2. Below, \* denotes the convolution operator  $\nu_1 * \nu_2(\varphi) := \iint \varphi(x_1 + x_2)\nu(dx_1)\nu(dx_2)$ , and L denotes the infinitesimal generator of the process X (see e.g. Sato (1999)), which is known to be given by

(2.2) 
$$(L\varphi)(x) := \frac{\sigma^2}{2} \varphi''(x) + b\varphi'(x) + \int \left(\varphi(y+x) - \varphi(x) - y\varphi'(x)\mathbf{1}_{\{|y| \le 1\}}\right) \nu(dy);$$

see Theorem 31.5 in [29] and also in Proposition 2.3 in [15]. The following result can be found in [15] (see Proposition 3.1), where a proof is provided for a certain class of unbounded functions  $\varphi$ :

**Lemma 2.2.** If  $\varphi \in C_b^2(\mathbb{R})$  vanishes in a neighborhood of the origin, then

(2.3) 
$$\lim_{t \to 0} \frac{1}{t} \left\{ \frac{1}{t} \ \mathbb{E} \varphi(X_t) - \nu(\varphi) \right\} = \nu_{\varepsilon}(L\varphi) - \frac{1}{2} \nu_{\varepsilon} * \nu_{\varepsilon}(\varphi),$$

where  $\nu_{\varepsilon}(dx) := \mathbf{1}_{\{|x| > \varepsilon\}} \nu(dx).$ 

The following is an easy consequence of the previous two results.

**Theorem 2.3.** Under the assumptions of Lemma 2.2,

(2.4) 
$$\sqrt{t_n} \left( \hat{\beta}^{\pi}(\varphi) - \nu(\varphi) \right) \xrightarrow{\mathfrak{D}} \nu(\varphi^2)^{\frac{1}{2}} Z$$

as  $t_n \to \infty$  and  $\bar{\pi} \to 0$  so that  $\bar{\pi}\sqrt{t_n} \to 0$ .

*Proof.* It suffices to prove that

$$D_n := \sqrt{t_n} \left( \mathbb{E} \ \hat{\beta}^{\pi}(\varphi) - \int \varphi(x) \nu(dx) \right) \longrightarrow 0.$$
<sup>37</sup>
<sub>38</sub>

Writing  $\Delta_k = t_k - t_{k-1}$  and using (2.3), for  $\bar{\pi} = \max_k \Delta_k$  small enough, there exists a constant C such that

$$|D_n| \le \frac{1}{\sqrt{t_n}} \sum_{k=1}^n \Delta_k \left| \frac{1}{\Delta_k} \mathbb{E} \varphi(X_{\Delta_k}) - \nu(\varphi) \right| \le C \frac{1}{\sqrt{t_n}} \sum_{k=1}^n \Delta_k^2 \le C \bar{\pi} \sqrt{t_n} \to 0,$$

by assumption.

**Remark 2.4.** As a direct consequence, it follows that  $\hat{\beta}^{\pi}(\varphi)$  is a consistent esti-48 mator for  $\nu(\varphi)$  as  $t_n \to \infty$  and  $\bar{\pi} \to 0$  so that  $\bar{\pi}\sqrt{t_n} \to 0$ . As a matter of fact, it 49 suffices that  $t_n \to \infty$  and  $\bar{\pi} \to 0$ , provided that e.g. f is  $\nu$ -continuous, bounded, and 50  $f(x) = o(|x|^2)$ , as  $x \to 0$ . A proof of this statement is outlined in [32] for regular 51 sampling observations, while the general case is considered in [14]. З

In view of the linearity of  $\hat{\beta}^{\pi}(\cdot)$  and  $\nu(\cdot)$ , we conclude that

**Corollary 2.5.** Let  $\Xi$  be the class of functions  $\varphi \in C_0^2(\mathbb{R})$  that vanish in a neighborhood of the origin. Suppose that the linear model S in (1.4) is such that  $\{\varphi_j\}_{j=1}^d \subseteq \Xi$ . Then, the projection estimator  $\hat{s}^{\pi}(x)$  in (1.9) satisfies the limiting relation

(2.5) 
$$\sqrt{t_n} \left( \hat{s}^{\pi}(x) - s^{\perp}(x) \right) \xrightarrow{\mathfrak{D}} V^{1/2}(x) Z$$

as  $t_n \to \infty$  and  $\bar{\pi} \to 0$  so that  $\bar{\pi}\sqrt{t_n} \to 0$ , where  $V(x) := \int f^2(y)\nu(dy)$  with  $f(y) := \sum_{j=1}^d \varphi_j(x)\varphi_j(y)$ .

**Remark 2.6.** Notice that we have the following bound for the variance

$$V(x) \le \|s\|_{\infty,D} \sum_{j=1}^d \varphi_j^2(x),$$

where  $||s||_{\infty,D} := \sup_{y \in D} s(y).$ 

We can relax the regularity conditions on the moment functions  $\varphi$  by using a simple integration by parts formula (see Remark 3.3 below). A different approach could be to impose additional regularity conditions on the Lévy process itself. In this direction, [28] studies series expansions for the transition density  $p_t(x)$  of  $X_t$  as powers of t. For instance, one of their results states that if  $p_t$  is monotonically decreasing for x > b and x < -c, for some b, c > 0, then for any  $\eta > 0$ , there exists  $\varepsilon' > 0$  and  $t_0 > 0$ , such that

(2.6) 
$$\frac{1}{t}p_t(x) = e^{-\int_{\{|y|>\varepsilon\}} s(y)dy} s(x) + O_{\varepsilon,\eta}(t),$$
27
28

for  $|x| > \eta$ . Such a result will allow us to estimate the rate of convergence in (1.3) if  $\varphi$  vanishes around the origin, since

$$\frac{1}{\Delta} \mathbb{E} \varphi(X_{\Delta}) - \nu(\varphi) = \int \varphi(x) \left\{ \frac{1}{\Delta} p_{\Delta}(x) - s(x) \right\} dx.$$

However, we should warn that the derivation of (2.6) in [28] is not completely formal<sup>2</sup>, and hence, we avoid to use such an approach in the sequel. See [17] for more insight on the small-time polynomial expansions of the transition distributions of the Lévy process.

#### 3. The Model Selection Problem

In this part we describe how to control the risk (1.6) of the projection estimators by imposing two conditions. First, the time horizon T should be large enough (compared to the complexity of the sieves), while the sampling frequency is kept small compared to the time horizon. These conditions will ensure that the variance term of (1.6) is of order  $O(T^{-1})$ . Second, the sieves should have good approximating properties in general classes of smooth functions so that when the Lévy density is

З

 $2^{2}$  The main problem arises from the application of Lemma 1 in [28]. The value of  $t_{0}$  actually 50 depends on  $\delta$ . Later on in their proof,  $\delta$  is taken arbitrarily small, which is likely to result in  $t_{0} \rightarrow 0$ 51 (unless otherwise proved).

presumed to have "degree of smoothness  $\alpha$ ", the bias term of (1.6) is of order  $O(m^{-\alpha})$ , where *m* is the dimension of the sieve (see Section 3.2 for the details). We prove that under the above conditions, we can tune up the dimension of the sieve to the presumed smoothness of *s* so that the rate of convergence of the risk is of order  $O(T^{-2\alpha/(2\alpha+1)})$ .

# 3.1. Analysis of the Variance Term

Consider the setting and notation of the introduction. For simplicity, we focus on estimation windows D in the positive reals (that is, D := [a, b], for some  $0 < a < b \le \infty$ ). By making the sampling frequency per unit time high enough relative to the sampling horizon T, we can estimate the rate at which the variance term of the risk (1.6) decreases in the time horizon T. In the subsequent sections, we will see that this estimate actually leads to a rate of convergence for the risk which is optimal, even if our estimation were based on the whole sample path  $\{X_t\}_{t \le T}$ . We shall need the following technical lemma, which we prove in the appendix for the sake of completeness.

**Lemma 3.1.** For any T > 0, there exist  $\delta_T > 0$  and k > 0 (independent of T) such that

$$\sup_{y \in D} \left| \frac{1}{\Delta} \mathbb{P} \left[ X_{\Delta} \ge y \right] - \nu([y, \infty)) \right| < k \frac{1}{T}$$

for all  $\Delta < \delta_{T}$ .

The mesh size  $\delta_T$  will play a very important role below as the asymptotic results in the sequel will hold true as far as the sampling frequency, measured by  $\bar{\pi} := \max\{t_k - t_{k-1}\}$ , is such that  $\bar{\pi} < \delta_T$ . Thus, from a practical point of view, estimating  $\delta_T$  is crucial. We will discuss this point in more detail in the Section 3.4.

The following easy estimate will be useful in the sequel.

**Lemma 3.2.** Suppose that  $\varphi$  has support  $[c, d] \subset \mathbb{R}_+ \setminus \{0\}$ , where  $\varphi$  is continuous with continuous derivative. Then,

$$\left|\frac{\mathbb{E}\,\varphi\left(X_{\Delta}\right)}{\Delta}-\nu(\varphi)\right|\leq\left(\left|\varphi(c)\right|+\int_{c}^{d}\left|\varphi'(u)\right|\,du\right)M_{\Delta}([c,d]),$$

where  $M_{\Delta}([c,d]) := \sup_{y \in [c,d]} \left| \frac{1}{\Delta} \mathbb{P} \left[ X_{\Delta} \ge y \right] - \nu([y,\infty)) \right|.$ 

*Proof.* The result follows from the following identities

$$\mathbb{E}\,\varphi(X_{\Delta}) = \varphi(c)\,\mathbb{P}\,\left[X_{\Delta} \ge c\right] + \int_{c}^{\infty} \varphi'(u)\,\mathbb{P}\,\left[X_{\Delta} \ge u\right] du,$$

$$\int \varphi(x)\nu(dx) = \varphi(c)\nu\left([c,\infty)\right) + \int_c^\infty \varphi'(u)\nu\left([u,\infty)\right) du.$$

These are standard consequences of Fubini's Theorem.

**Remark 3.3.** We can apply the previous two lemmas to obtain CLTs for  $\hat{\beta}^{\pi}$  and  $\hat{s}^{\pi}$ . 48 Indeed, if  $\varphi$  is as in Lemma 3.2 and, for each T, the partition  $\pi_T$  has mesh smaller 49 than  $\delta_T$ , the critical value in Lemma 3.1, then (2.3) holds true. The projection 50 estimator  $\hat{s}^{\pi}$  will satisfy (2.5) provided that the basis functions  $\varphi$  are as in Lemma 51 3.2. We are now ready to estimate the variance term. We shall impose conditions on the approximating linear models so that the estimates of the above Lemmas are applicable.

**Standing assumption 1.** The linear model S of (1.4) is generated by an orthonormal basis  $\mathcal{G} := \{\varphi_j\}_{j=1}^d$  such that each  $\varphi_j$  is bounded with continuous derivative on the interior of its support, which is assumed to be of the form  $[x_{j-1}, x_j] \subset D$ .

In the sequel, we will need the following notation:

(3.1) 
$$D_1(\mathcal{S}) := \inf_{\mathcal{G}} \max\left\{ \|\varphi\|_{\infty}^2 : \varphi \in \mathcal{G} \right\},$$

<sup>11</sup> (3.2) 
$$D_2(\mathcal{S}) := \inf_{\mathcal{G}} \max\left\{ \|\varphi'\|_1^2 : \varphi \in \mathcal{G} \right\},$$

where the infimums are over all orthonormal bases  $\mathcal{G}$  of  $\mathcal{S}$ .

**Proposition 3.4.** There exists a constant K > 0 such that

<sup>16</sup>  
<sub>17</sub> (3.3) 
$$\mathbb{E} \|s^{\perp} - \hat{s}^{\pi}\|^2 \le K \frac{\dim(\mathcal{S})}{T},$$
 <sup>16</sup>

for any linear model S satisfying the Standing Assumption 1, and for any partition  $\pi : 0 = t_0 < \cdots < t_n = T$  such that  $T > \max\{D_1(\mathcal{S}), D_2(\mathcal{S})\}$  and  $\bar{\pi} < \delta_T$ , where  $\delta_{\tau}$  is the "critical" mesh size introduced in Lemma 3.1.

*Proof.* Fix an orthonormal basis  $\mathcal{G} := \{\varphi_j\}_{j=1}^d$  of  $\mathcal{S}$ . Let  $D_t(\varphi) := \frac{1}{t} \mathbb{E} \varphi(X_t) - \nu(\varphi)$ . For any  $\varphi_j \in \mathcal{G}$ , we have

$$\mathbb{E}\left\{\hat{\beta}^{\pi}(\varphi_j) - \nu(\varphi_j)\right\}^2 \leq \frac{1}{t_n} \int \varphi_j^2(x)\nu(dx) + 25$$

$$+\frac{1}{t_n^2}\sum_{k=1}^n \left|D_{\Delta_k}(\varphi_j^2)\right|\Delta_k + \left\{\frac{1}{t_n}\sum_{k=1}^n \left|D_{\Delta_k}(\varphi_j)\right|\Delta_k\right\}^2,$$

where  $\Delta_k := t_k - t_{k-1}$ . Then, from the previous two lemmas, when  $\bar{\pi} < \delta_T$ ,

$$\mathbb{E}\left\{\hat{\beta}^{\pi}(\varphi_j) - \nu(\varphi_j)\right\}^2 \leq \frac{1}{T} \int \varphi_j^2(x)\nu(dx)$$

$$+\frac{k}{T^2}\left(\left|\varphi_j^2(x_{j-1})\right| + \int_{x_{j-1}}^{x_j} \left|2\varphi_j(u)\varphi_j'(u)\right| du\right)$$

$$+\frac{k^2}{T^2}\left(|\varphi_j(x_{j-1})|+\int_{x_{j-1}}^{x_j}|\varphi_j'(u)|\,du\right)^2,$$

which can be simplified further as follows

$$\mathbb{E}\left\{\hat{\beta}^{\pi}(\varphi_j) - \nu(\varphi_j)\right\}^2 \leq \frac{1}{T} \int \varphi_j^2(x)\nu(dx) + \frac{2k^2}{T^2} \left(\|\varphi_j\|_{\infty} + \|\varphi_j'\|_1\right)^2$$

$$\leq \frac{\|s \cdot \chi_{\scriptscriptstyle D}\|_{\infty}}{T} + 8k^2 \frac{\max_{j'} \|\varphi_{j'}\|_{\infty}^2 + \|\varphi_{j'}'\|_1^2}{T^2}.$$

Then,

$$\mathbb{E} \|s^{\perp} - \hat{s}^{\pi}\|^2 \le \frac{\dim(\mathcal{S})}{T} \left\{ \|s \cdot \chi_{_D}\|_{\infty} + 8k^2 \frac{\max_{j'} \|\varphi_{j'}\|_{\infty}^2 + \|\varphi_{j'}'\|_1^2}{T} \right\}.$$

$$\mathbb{E} \|s - s\| \leq \frac{T}{T} \left\{ \|s \cdot \chi_D\|_{\infty} + \delta \kappa - \frac{T}{T} \right\}.$$

Now, it is evident that (3.3) holds whenever  $T > \max\{D_1(\mathcal{S}), D_2(\mathcal{S})\}$ . 

З

### 3.2. The Approximation Error for Besov Type Smooth Functions

As it is customary, the bias term in (1.6) will be estimated by imposing certain degree of smoothness on the function s. Concretely, the restriction of the Lévy density s to D := [a, b] is assumed to belong to the Besov space  $\mathcal{B}^{\alpha}_{\infty}(L^p([a, b]))$ for some  $p \in [2, \infty]$  and  $\alpha > 0$  (see for instance [12] and references therein for background on these spaces). The space  $\mathcal{B}^{\alpha}_{\infty}(L^p([a, b]))$  consists of those functions f belonging to  $L^p([a, b])$  if  $0 (or being uniformly continuous if <math>p = \infty$ ) such that

$$|f|_{\mathcal{B}^{\alpha}_{\infty}(L^{p})} \equiv \sup_{\delta>0} \frac{1}{\delta^{\alpha}} \sup_{0 < h \le \delta} \|\Delta_{h}^{r}(f, \cdot)\|_{p} < \infty,$$

with  $r := [\alpha] + 1$ . Here,  $\Delta_h(f, x) \equiv f(x+h) - f(x)$  and  $\Delta_h^r(f, x)$  is the  $r^{th}$ -order difference of f defined recursively by

$$\Delta_h^r(f,x) \equiv \Delta_h(\Delta_h^{r-1}(f,\cdot),x),$$

for x's such that  $x + rh \in D$  and  $r \in \mathbb{N}$ .

The Besov class is closely related to the so-called class of Lipschitz functions. For constants  $k \in \mathbb{N}$  and  $\beta \in (0, 1]$ , f is said to belong to  $\operatorname{Lip}(k + \beta, L^p([a, b]))$  if  $f, \ldots, f^{(k-1)}$  are absolutely continuous (on [a, b]) and  $f^{(k)}$  belongs to  $L^p((a, b))$  and satisfies

(3.4) 
$$\sup_{h>0} \frac{1}{h^{\beta}} \|\Delta_h(f^{(k)}, \cdot)\|_p < \infty.$$

It is known that if  $\beta < 1$  and  $1 \le p \le \infty$ , then  $f \in \operatorname{Lip}(k + \beta, L^p([a, b]))$  if and only if f is a.e. equal to a function in  $\mathcal{B}^{\alpha}_{\infty}(L^p([a, b]))$  with  $\alpha := k + \beta$ . In general,  $\operatorname{Lip}(k + \beta, L^p([a, b])) \subset \mathcal{B}^{k+\beta}_{\infty}(L^p([a, b]))$ , for any 0 (see e.g. [12]). Notice $that when <math>p = \infty$ , the condition (3.4) takes the form:

(3.5) 
$$|f^{(k)}(x) - f^{(k)}(y)| \le L|x - y|^{\beta},$$

for all  $x, y \in (a, b)$  and some  $L < \infty$ .

An important reason for working with the Besov-type smooth functions is the availability of estimates of the approximation error by splines, trigonometric polynomials, and wavelets (see [12] and [3] for more details). For instance, if  $S_{k,m}$  denotes the space of piecewise polynomials of degree at most k, based on a regular partition of [a, b] with m classes, and  $s \in \mathcal{B}^{\alpha}_{\infty}(L^p([a, b]))$  with  $\alpha < k + 1$ , then there exists a constant  $c_{[\alpha]} < \infty$  such that

(3.6) 
$$\inf_{f \in S_{k,m}} \|s - f\|_p \le c_{[\alpha]} |s|_{\mathcal{B}^{\alpha}_{\infty}(L^p)} (b - a)^{\alpha} m^{-\alpha}.$$

Thus, when  $p \ge 2$ , the orthogonal projection of s on  $\mathcal{S}_{k,m}$ , denoted by  $s_m^{\perp}$ , is such that

(3.7) 
$$\|s - s_m^{\perp}\| \le c_{[\alpha]} (b - a)^{\frac{1}{2} - \frac{1}{p} + \alpha} |s|_{\mathcal{B}_{\infty}^{\alpha}(L^p)} m^{-\alpha}.$$

<sup>47</sup> Notice that the elements of  $S_{k,m}$  are not necessarily smooth (not even continuous) <sup>48</sup> and hence, they are not "splines" in the standard sense of the literature, where a <sup>49</sup> spline is understood as a smooth piece-wise polynomial. The upper bound (3.6) <sup>50</sup> is actually true if we restrict to certain splines of  $S_{k,m}$  (say B-splines) (see (10.1) <sup>51</sup> in Chapter 2 of [12]). For the sake of completeness let us describe in detail the J.E. Figueroa-López

space  $S_{k,m}$  as well as give estimates for the constants (3.1-3.2). Let  $Q_j$  be the Legendre polynomials of order j on  $\mathbb{L}^2([-1,1], dx)$ . The space  $\mathcal{S}_{k,m}$  is generated by the orthonormal functions

$$\hat{\varphi}_{i,j}(x) := \sqrt{\frac{2j+1}{x_i - x_{i-1}}} Q_j\left(\frac{2x - (x_i + x_{i-1})}{x_i - x_{i-1}}\right) \mathbf{1}_{(x_{i-1}, x_i)}(x),$$

for i = 1, ..., m and j = 0, ..., k, and where  $a = x_0 < \cdots < x_m = b$  are equallyspaced points. It is well-known that  $|Q_j(x)| \leq 1$  and  $|Q'_j(x)| \leq Q'_j(1) = \frac{j(j+1)}{2}$ . Then, fixing  $\Delta_x := x_i - x_{i-1} = \frac{b-a}{m}$ , we have that

$$\hat{\varphi}_{i,j}'(x) = 2\sqrt{2j+1}\,\Delta_x^{-3/2}\,Q_j'\left(\frac{2x-(x_i+x_{i-1})}{x_i-x_{i-1}}\right)\mathbf{1}_{(x_{i-1},x_i)}(x),$$

$$\|\hat{\varphi}_{i,j}'\|_1 \le 2\sqrt{2j+1}\Delta_x^{-3/2} \int_{x_{i-1}}^{x_i} \sup_u |Q_j'(u)| dx \le \sqrt{2j+1}\Delta_x^{-1/2}(j)(j+1).$$

It is now clear that

$$D_2(\mathcal{S}_{k,m}) \le \max_{i,j} \left\{ \|\varphi'_{i,j}\|_1^2 \right\} \le \frac{(k+1)^2 k^2 (2k+1)}{b-a} \, m.$$

In a similar manner one can check that

$$D_1(\mathcal{S}_{k,m}) \le \frac{(k+1)^2(2k+1)}{b-a} m.$$

### 3.3. Rate of Convergence for Smooth Functions Via Splines

As a consequence of the variance and bias term estimates given in the previous two parts, we now estimate the rate of convergence on D of the projection estimators (1.9), on the regular piece-wise polynomials  $\{S_{k,m}\}_{m\geq 1}$ , assuming that the Lévy density s is in the Besov class  $\mathcal{B}^{\alpha}_{\infty}(L^p([a, b]))$  with  $p \geq 2$  and  $\alpha < k + 1$ . It turns out that under the stated conditions, projection estimators converge at a rate at least as good as  $T^{-2\alpha/(2\alpha+1)}$ . The following result is valid provided that, for each time horizon T, the mesh of the sampling times  $\pi_T$  is smaller than the critical mesh  $\delta_T$  introduced in Lemma 3.1. In Section 4, we will see that this rate is actually the best possible even under continuous sampling.

**Proposition 3.5.** Let  $\hat{m}_T := [T^{1/(2\alpha+1)}]$  and let  $\Theta(R, L)$  be the class of Lévy densities s such that  $\|s \cdot \chi_D\|_{\infty} < R$ , and such that the restriction of s to D := [a, b] is a member of  $\mathcal{B}^{\alpha}_{\infty}(L^p([a, b]))$  with  $|s|_{\mathcal{B}^{\alpha}_{\infty}(L^p)} < L$  and  $p \geq 2$ . Then,

(3.8) 
$$\limsup_{T \to \infty} T^{2\alpha/(2\alpha+1)} \sup_{s \in \Theta(R,L)} \mathbb{E}\left[ \|s - \hat{s}_T\|^2 \right] < \infty,$$

where for each T, the estimator  $\hat{s}_{T}$  is given by (1.7) and (1.9) with  $S = S_{k,\hat{m}_{T}}$ ,  $k > \alpha - 1$ , and a mesh  $\bar{\pi}_{T}$  smaller than  $\delta_{T}$ .

<sup>47</sup> Proof. From the two previous parts, there exists a constant K (depending on  $k, a, b, \alpha, R, p, L$ ) such that

$$\|s - s_m^{\perp}\| \le K m^{-\alpha} \quad \text{and} \quad \mathbb{E} \|s^{\perp} - \hat{s}_m^{\pi}\|^2 \le K \frac{m}{T},$$

for  $m \in \mathcal{M}_T := \{m': T > Km'\}$  and  $\bar{\pi} < \delta_T$ . Then for a constant M and for large enough T,

$$\sup_{\epsilon \in \Theta(R,L)} \mathbb{E} \left[ \|s - \hat{s}_{T}\|^{2} \right] \le M \left\{ \left[ T^{1/(2\alpha+1)} \right]^{-2\alpha} + \left[ T^{1/(2\alpha+1)} \right] T^{-1} \right\}.$$

The limit (3.8) is now clear.

s

**Example 3.6.** If s has continuous bounded derivative on  $D := [a, b] \subset \mathbb{R} \setminus \{0\}$ (hence,  $s \in \mathcal{B}^{\alpha}_{\infty}(L^{\infty}([a,b]))$ ), for any  $\alpha < 1$ ), then one can construct regular histogram estimators converging to s on D at a rate faster than  $T^{-1/2}$  if one selects the number of classes approximately equal to  $T^{1/2}$  and the mesh of the partition  $\pi$ smaller than  $\delta_{T}$ .

#### 3.4. About the Critical Mesh

The critical mesh, introduced in Lemma 3.1, gives a bound on the mesh of the sam-pling frequency needed to estimate in a simple way the rate of convergence of the variance term (see Proposition 3.4). Of course, any hope for a feasible implementa-tion of this estimation scheme will require an explicit estimate of this critical mesh. In the compound Poisson case (when  $\nu(\mathbb{R}\setminus\{0\}) < \infty$ ), it turns out that  $\delta_{\tau} = o(\frac{1}{\tau})$ suffices. In the general case, we have the following result, which tell us, in partic-ular, that the sampling frequency needs to be higher when one wishes to estimate the Lévy density closer to the origin.

**Proposition 3.7.** Let  $\rho > 0$  such that  $a\rho > 1$ . Then, there exists  $T_0(\rho) > 0$  and k > 0 such that

$$\sup_{y \in D} \left| \frac{1}{\Delta} \mathbb{P} \left[ X_{\Delta} \ge y \right] - \nu([y, \infty)) \right| < k \frac{1}{T}$$

for all  $T > T_0$  and  $\Delta < T^{-\frac{1}{\rho}T}$ .

*Proof.* As in the proof of Lemma 3.1, we can obtain

$$\sup_{y \in D} \left| \frac{1}{t} \mathbb{P} \left[ X_t \ge y \right] - \nu([y, \infty)) \right| \le \frac{1}{t} \mathbb{P} \left[ X_t^{\varepsilon} \ge a \right] + 2c \frac{1}{T} + \nu([a, \infty)) \mathbb{P} \left[ |X_t^{\varepsilon}| \ge \eta \right]$$

$$+ \, \lambda_{\varepsilon} \, \mathbb{P} \, \left[ |X_t^{\varepsilon}| \geq \eta \right] + \nu([a,\infty)) \, \lambda_{\varepsilon} t + \lambda_{\varepsilon}^2 t,$$

valid for T > 1/a,  $\eta = \frac{1}{T}$ , and  $0 < \varepsilon < a - \eta$  (here  $c := \sup_{a - \eta \le x \le b + \eta} s(x)$ ). Fix  $\varepsilon > 0$  sufficiently small so that  $\rho < \frac{1}{\varepsilon}$ . Let us recall that there exists  $y_0 := y_0(\rho)$ such that

 $\mathbb{P}\left[|X_t^{\varepsilon}| \ge y\right] \le \exp\{\rho y_0 \log y_0\} \exp\{\rho y - \rho y \log y\} t^{y\rho}$ 

for all  $t < \frac{y}{y_0(\rho)}$  (see e.g. [28]). In particular, when  $y = \eta = \frac{1}{T}$  and  $t < T^{-\frac{1}{\rho}T}$ , for T sufficiently large that  $T^{-\frac{1}{\rho}T+1} < \frac{1}{w(\rho)}$ ,

$$\mathbb{P}\left[|X_t^{\varepsilon}| \ge \eta\right] \le kT^{-1}.$$

Similarly, when y = a and  $t < T^{-\frac{1}{\rho}T}$ ,

$$\frac{1}{t} \mathbb{P}\left[X_t^{\varepsilon} \ge a\right] \le kt^{a\rho-1} < kT^{-\frac{1}{\rho}T(a\rho-1)} < kT^{-1}, \tag{49}$$

if  $T > \frac{\rho}{a\rho-1}$ . This proves the result since  $\varepsilon$  is fixed.

З

**Remark 3.8.** The estimate of the critical mesh given in Proposition 3.7 can be improved substantially. Indeed, in a forthcoming paper, we will show that it suffices that  $\Delta = o(T^{-1})$ .

#### 4. Minimax Risk of Estimation for Smooth Lévy Densities

In this section, we show that the rate of convergence  $O(T^{-2\alpha/(2\alpha+1)})$  attained by projection estimators is the best possible, in the sense that there is no estimator  $\hat{s}_T^*$  that can converge to s faster than  $T^{-2\alpha/(2\alpha+1)}$ , for any  $s \in \Theta$ , even assuming continuous-time sampling. In order to prove this, we will assess the long-run behavior of the minimax risk on  $\Theta$ , roughly defined as

$$\inf_{\hat{s}} \sup_{s \in \Theta} \mathbb{E}_{s} \left[ d\left(s, \hat{s}\right) \right],$$

where the infimum is taken over all possible estimators  $\hat{s}$ , and  $d(s, \hat{s})$  measures the distance between  $\hat{s}$  and s.

Traditionally, the performance of nonparametric estimators is gauged by comparing the rate of convergence of the estimator in question to the rate of convergence of the minimax risk when the available data increases. The rates of convergence of minimax risks are available in most of the traditional nonparametric problems. For instance, Ibragimov and Has'minskii [19] and Barron et. al. [2] provided this kind of asymptotics for the problem of density estimation based on i.i.d. random variables, while Kutoyants [22] and Reynaud-Bouret [25] considered the problem of intensity estimation of a *finite* Poisson point processes. This last set-up is relevant for our problem since the jumps of a Lévy process can be associated with a (possibly infinite) Poisson point process on  $\mathbb{R}_+ \times \mathbb{R} \setminus \{0\}$  (see e.g. Theorem 19.2 in [29]). Using this connection, we adapt below a result from [22] to obtain the long-run asymptotics of the minimax risk of estimation of the Lévy density off the origin. The idea of the proof, due to Ibragimov and Has'minskii [19], is based on the statistical toolbox for distributions satisfying the *Local Asymptotic Normality* (LAN) property (see Chapters II and Section IV.5 of [19]).

Let us introduce some notation. Here,  $\ell : \mathbb{R} \to \mathbb{R}$  stands for a *loss function* satisfying the following:

(i)  $\ell(\cdot)$  is nonnegative,  $\ell(0) = 0$  but not identically 0, and  $\ell$  continuous at 0;

(ii)  $\ell$  is symmetric:  $\ell(u) = \ell(-u)$  for all u;

(iii)  $\{u : \ell(u) < c\}$  is a convex set for any c > 0, ;

(iv)  $\ell(u) \exp{\{\varepsilon | u |^2\}} \to 0$  as  $|u| \to \infty$  for any  $\varepsilon > 0$ .

We consider Lévy densities  $s : \mathbb{R} \setminus \{0\} \to \mathbb{R}_+$  that are k times differentiable on an interval  $[a, b] \subset \mathbb{R} \setminus \{0\}$  and satisfy (3.5) for all  $x, y \in [a, b]$ . For given  $k \in \mathbb{N}$  and  $\beta \in (0, 1]$ , we denote such a class of functions by  $\Theta_{k+\beta}(L; [a, b])$ . The proof of the result below is presented in the Appendix 6.3.

**Theorem 4.1.** If  $x_0$  is an interior point of the interval [a, b], then

(4.1) 
$$\liminf_{T \to \infty} \left\{ \inf_{\hat{s}_T} \sup_{s \in \Theta} \mathbb{E}_s \left[ \ell \left( T^{\alpha/(2\alpha+1)} \left( \hat{s}_T(x_0) - s(x_0) \right) \right) \right] \right\} > 0,$$

where  $\alpha := k + \beta$ ,  $\Theta := \Theta_{\alpha}(L; [a, b])$  and the infimum is over all the estimators  $\hat{s}_{T}$  of s based on  $\{X(t)\}_{0 \le t \le T}$ .

The previous result can be strengthen to be uniform in  $x_0 \in (a, b)$  and as a consequence, the long-run behavior of the minimax risk under the integrated mean-square distance can be assessed. The proof of the next result is given in Appendix 6.3.

**Corollary 4.2.** Under the notation and conditions of Theorem 4.1, the following two limits hold:

(4.2) 
$$\lim_{T \to \infty} \left\{ \inf_{\hat{s}_T} \inf_{x \in (a,b)} \sup_{s \in \Theta} \mathbb{E}_s \left[ \ell \left( T^{\alpha/(2\alpha+1)} \left( \hat{s}_T(x) - s(x) \right) \right) \right] \right\} > 0,$$

(4.3) 
$$\liminf_{T \to \infty} T^{2\alpha/(2\alpha+1)} \left\{ \inf_{\hat{s}_T} \sup_{s \in \Theta} \mathbb{E}_s \left[ \int_a^b \left( \hat{s}_T(x) - s(x) \right)^2 dx \right] \right\} > 0.$$

**Remark 4.3.** The previous result is also valid for classes slightly smaller than  $\Theta_{\alpha}(L;[a,b])$  such as

$$\Theta = \Theta_{\alpha}(L; [a, b]) \cap \{s : \|s\|_{\mathbb{L}^{\infty}([a, b])} < R\},$$
 16

which is closely related to the Besov class  $\Theta(R, L)$  of (3.8). Indeed,  $\Theta_{\alpha}(L; [a, b])$  is contained in  $\mathcal{B}^{\alpha}_{\infty}(\mathbb{L}^{\infty}([a, b]))$  (see Section 2.9 of [12]), and thus,

(4.4) 
$$\liminf_{T \to \infty} T^{2\alpha/(2\alpha+1)} \left\{ \inf_{\hat{s}_T} \sup_{s \in \Theta(R,L)} \mathbb{E}_s \left[ \int_a^b \left( \hat{s}_T(x) - s(x) \right)^2 dx \right] \right\} > 0.$$

We conclude that there is no reasonable estimator  $\hat{s}_{\tau}$  of s capable of outperforming the rate  $T^{-2\alpha/(2\alpha+1)}$  uniformly on  $\Theta$ : there is always an  $s \in \Theta$  for which

$$T^{2\alpha/(2\alpha+1)} \mathbb{E}_s \left[ \|\hat{s}_T - s\|^2 \right] > B,$$

for some B > 0 and for large enough T. Therefore, the estimator described in Proposition 3.5 achieves the optimum rate of convergence on  $\Theta(R, L)$  from a minimax point of view.

### 5. A Data-Driven Selection Method and Adaptability

The model selection criterion described in Section 3.3, where one tunes up the number of classes m to the "smoothness" of s, has the obvious drawback of requiring (or at least presuming) the smoothness parameter  $\alpha$ . In the literature of nonparametric statistics, one wishes to devise data-driven selection methods that can *adapt* to arbitrary degree of smoothness (see e.g. Birgé and Massart [5] for an extensive exposition of the topic).

A typical approach for adaptive model selection schemes consists of minimizing an unbiased estimator of the risk of estimation. This approach was developed in [13] in the context of Lévy density estimation. Let us briefly discuss the findings there. The key idea comes from the following refinement of (1.6):

(5.1) 
$$\mathbb{E}\left[\|s - \hat{s}^{c}\|^{2}\right] = \|s\|^{2} + \mathbb{E}\left[-\|\hat{s}^{c}\|^{2} + \operatorname{pen}^{c}(\mathcal{S})\right],$$

where  $\hat{s}^c$  is as in (1.9) substituting  $\hat{\beta}^{\pi}(\varphi)$  by the statistics  $\hat{\beta}^c(\varphi)$  of (1.13),  $s^{\perp}$  is the orthogonal projection in (1.5), and pen<sup>c</sup>(S) is defined in terms of an orthonormal basis  $\mathcal{G} := \{\varphi_1, \ldots, \varphi_d\}$  of S by the formula:

50 (5.2) 
$$\operatorname{pen}^{c}(\mathcal{S}) \equiv \frac{2}{T^{2}} \sum_{t \leq T} \sum_{\varphi \in \mathcal{G}} \varphi^{2}(\Delta X_{t}).$$
 50 51

imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009

З

Equation (5.1) shows that the risk of  $\hat{s}^c$  moves "parallel" to the expectation of the observable statistics  $-\|\hat{s}^c\|^2 + \text{pen}^c(\mathcal{S})$ , suggesting the selection of the model that minimizes such statistics. Concretely, given a collection of sieves  $\{S_m, m \in \mathcal{M}\}$ , we З should choose the projection estimator  $\tilde{s}^c \equiv \hat{s}^c_{\hat{m}}$ , where 

$$\hat{m} \equiv \operatorname{argmin}_{m \in \mathcal{M}} \left\{ -\|\hat{s}_m^c\|^2 + \operatorname{pen}^c(\mathcal{S}_m) \right\}.$$

Such an estimator  $\tilde{s}^c$  is called a *penalized projection estimator* (p.p.e.) since the role of pen<sup>c</sup>( $\mathcal{S}$ ) is to penalize large linear models.

In [16], it is shown that the p.p.e.  $\tilde{s}^c$  is adaptive in the class of Besov Lévy densities of Section 3.2 in the sense that  $\tilde{s}^c$  attains the optimal rate of convergence  $O(T^{-2\alpha/(2\alpha+1)})$  without using the knowledge of  $\alpha$ . Unfortunately, the previous approach intrinsically requires continuous-time sampling of the process to determine the jumps  $\Delta X_t$ . However, the analysis could still be useful if one uses the natural discrete-based proxies of  $\beta^c$  and pen<sup>c</sup>, where the jumps  $\Delta X_t$  are replaced by the increments  $X_{t_k} - X_{t_{k-1}}$ . This idea leads to the estimators  $\hat{s}^{\pi}$  in (1.9) and to the statistic

(5.3) 
$$pen^{\pi}(S) = \frac{2}{T^2} \sum_{k=1}^{n} \sum_{\varphi \in \mathcal{G}} \varphi^2 (X_{t_k} - X_{t_{k-1}})$$
<sup>18</sup>
<sup>19</sup>
<sup>20</sup>
<sup>20</sup>

as the penalization term. In the light of the previous arguments, we proposed a discrete-based model selection criterion as follows

(5.4) 
$$\hat{m}^{\pi} \equiv \operatorname{argmin}_{m \in \mathcal{M}} \left\{ -\|\hat{s}_{m}^{\pi}\|^{2} + \operatorname{pen}^{\pi}(\mathcal{S}_{m}) \right\}$$

$$= \operatorname{argmin}_{m \in \mathcal{M}} \left\{ -\sum_{\varphi \in \mathcal{G}_m} \{ \hat{\beta}^{\pi}(\varphi) \}^2 + \operatorname{pen}^{\pi}(\mathcal{S}_m) \right\}$$

where  $\mathcal{G}_m$  is an orthonormal basis of  $\mathcal{S}_m$ ,  $\hat{\beta}^{\pi}$  is given by (1.7), and pen<sup> $\pi$ </sup> is given by (5.3). The resulting estimator

(5.5) 
$$\widetilde{s} := s_{\hat{m}}^{\pi}$$

will be called (discrete-based) penalized projection estimator.

We hope to extend in a future work the adaptability result in [16] for this discretebased p.p.e. In the sequel, we illustrate the performance of these estimators for an infinite-jump activity Lévy process of relevance in the area of mathematical finance.

### 6. An Example: Estimation of Variance Gamma Processes.

#### 6.1. The Model

Variance Gamma processes were proposed in [23] (see also [8]) as substitutes to the Brownian Motion in the Black-Scholes model. Since their introduction, this kind of processes have received a great dealt of attention, even in the financial industry. For an introduction to many basic properties of variance Gamma processes and other related processes, the reader is referred to Knotz et al. [21].

There are two useful representations for this type of processes. A variance Gamma process  $X = \{X(t)\}_{t \ge 0}$  is a Brownian motion with drift, time changed by a Gamma Lévy process. Concretely, 

51 (6.1) 
$$X(t) = \theta U(t) + \sigma W(U(t)),$$

where  $\{W(t)\}_{t>0}$  is a standard Brownian motion,  $\theta \in \mathbb{R}$ ,  $\sigma > 0$ , and  $U = \{U(t)\}_{t>0}$ is an independent Gamma Lévy process with density at time t given by

(6.2) 
$$f_t(x) = \frac{x^{t/\nu - 1} \exp\left(-\frac{x}{\nu}\right)}{\nu^{t/\nu} \Gamma\left(\frac{t}{\nu}\right)}.$$

Notice that E[U(t)] = t and  $Var[U(t)] = \nu t$ ; therefore, the random clock U has a "mean rate" of one and a "variance rate" of  $\nu$ . There is no loss of generality in restricting the mean rate of the Gamma process U to one since, as a matter of fact, any process of the form

$$\theta_1 V(t) + \sigma_1 W(V(t)),$$

where V(t) is an arbitrary Gamma Lévy process,  $\theta_1 \in \mathbb{R}$ , and  $\sigma_1 > 0$ , has the same law as a process of the form (6.1) with suitably chosen  $\theta$ ,  $\sigma$ , and  $\nu$ . This a consequence of the self-similarity<sup>3</sup> property of Brownian motion and the fact that  $\nu$  in (6.2) is a scale parameter.

The process X is itself a Lévy process since Gamma processes are subordinators X(see Theorem 30.1 of [29]). Moreover, it is not hard to check that "statistically" Xis the difference of two Gamma Lévy processes (see e.g. (2.1) of [6]):

(6.3) 
$$\{X(t)\}_{t\geq 0} \stackrel{\mathfrak{D}}{=} \{X_+(t) - X_-(t)\}_{t\geq 0},$$

where  $\{X_{+}(t)\}_{t>0}$  and  $\{X_{-}(t)\}_{t>0}$  are Gamma Lévy processes with respective Lévy measures

$$\nu_{\pm}(dx) = \alpha \exp\left(-\frac{x}{\beta^{\pm}}\right) dx, \text{ for } x > 0.$$
<sup>23</sup>
<sup>24</sup>

Here,  $\alpha = 1/\nu$  and

As a consequence of this decomposition, the Lévy density of X takes the form

$$s(x) = \begin{cases} \frac{\alpha}{|x|} \exp\left(-\frac{|x|}{\beta^{-}}\right) & \text{if } x < 0, \end{cases}$$

(6.4) 
$$s(x) = \begin{cases} |x| \exp\left(-\frac{\beta}{\beta^+}\right) & \text{if } x > 0, \\ \frac{\alpha}{x} \exp\left(-\frac{x}{\beta^+}\right) & \text{if } x > 0, \end{cases}$$

where  $\alpha > 0, \beta^- \ge 0$ , and  $\beta^+ \ge 0$  (of course,  $|\beta^-| + |\beta^+| > 0$ ). As in the case of Gamma Lévy processes,  $\alpha$  controls the overall jump activity, while  $\beta^+$  and  $\beta^$ take respectively charge of the intensity of large positive and negative jumps. In particular, the difference between  $1/\beta^+$  and  $1/\beta^-$  determines the frequency of drops relative to rises, while their sum measures the frequency of large moves relative to small ones.

# 6.2. The Simulation Procedure

The above two representations provide straightforward methods to simulate a vari-ance Gamma model. One way will be to simulate the Gamma Lévy processes  $\{X_{+}(t)\}_{0 \le t \le T}$  and  $\{X_{-}(t)\}_{0 \le t \le T}$  of (6.3) using the series representation method introduced in Rosiński [26]. The other approach is to generate the random time change  $\{U(t)\}_{0 \le t \le T}$  of (6.1), and then construct a discrete skeleton from the incre-ments  $X(i\Delta t) - X((i-1)\Delta t)$ ,  $i \geq 1$ . The increments of X are simply simulated using normal random variables with mean and variances determined by the increments of U. 

<sup>3</sup>namely, 
$$\{W(ct)\}_{t\geq 0} \stackrel{\mathfrak{D}}{=} \{c^{1/2}W(t)\}_{t\geq 0}$$
, for any  $c>0$ . 51

### 6.3. The Numerical Results

In this part we illustrate the performance of the projection estimators (1.9) and the model selection criterion described in Section 5 using simulation experiments. The approximating linear models  $S_m$  considered here are the span of the indicator functions  $\chi_{[x_0,x_1]}, \ldots, \chi_{(x_{m-1},x_m]}$ , where  $x_0 < \cdots < x_m$  is a regular partition of an interval  $D \equiv [a, b]$ , with 0 < a or b < 0. We perform the following numerical experiment. First, we simulate the variance gamma Lévy process with specified (known) parameter settings. Then, we apply the penalized projection estimator defined by (5.4)-(5.5). Finally, to assess the accuracy of the nonparametric estimator, the true parametric model of s is subsequently fit to the nonparametric estimator using a least-square errors method. Concretely, if  $\tilde{s} := \hat{s}_{m\pi}^{\pi}$  is the discrete-based p.p.e. and  $s^{\theta}$  is the function (6.4), where we set  $\theta := (\alpha, \beta^-, \beta^+)$ , then we find

(6.5)

$$\hat{\theta}_{_{NP}} := \operatorname{argmin}_{\theta} \sum_{i=0}^{\hat{m}^{\pi}-1} (\widetilde{s}(\bar{x}_i) - s^{\theta}(\bar{x}_i)^2$$

where  $\bar{x}_i$  is the midpoint of the interval  $[x_i, x_{i+1}]$ . This approach provides a nonparametric based estimators for the parameters of the variance Gamma process.

Notice that, from an algorithmic point of view, the estimation for the variance Gamma model using penalized projection is not different from the estimation for the Gamma process. We can simply estimate both tails of the variance Gamma process separately. However, from the point of view of maximum likelihood estimation (MLE), the problem is numerically challenging. Even though the marginal density functions have "closed" form expressions<sup>4</sup> (see [8]), there are well-documented issues with MLE (see for instance [24]). The likelihood function is highly flat for a wide range of parameters and good starting values as well as convergence are critical. Also, the separation of parameters and the identification of the variance Gamma process from other classes of the generalized hyperbolic Lévy processes is difficult. In fact, difference between subclasses in terms of likelihood is small. It is important to mention that these issues worsen when dealing with "high-frequency" data.

Let us consider a numerical example motivated by the empirical findings of [8] based on daily returns on the S&P stock index from January 1992 to September 1994 (see their Table I). Using maximum likelihood methods, the annualized estimates of the parameters for the variance Gamma model were reported to be  $\hat{\theta}_{ML} = -0.00056256$ ,  $\hat{\sigma}_{ML}^2 = 0.01373584$ , and  $\hat{\nu}_{ML} = 0.002$ , from where we obtain  $\hat{\alpha}_{ML} = 500$ ,  $\hat{\beta}_{ML}^+ = 0.0037056$ , and  $\hat{\beta}_{ML}^- = 0.0037067$ . Figures 1 and 2 show respectively the left- and right- tails of the true Lévy

Figures 1 and 2 show respectively the left- and right- tails of the true Lévy density and the (discrete-based) penalized projection estimator as well as their corresponding best-fit variance Gamma Lévy densities using (6.5), and their marginal probability density functions (pdf) scaled by  $1/\Delta t$  (the reciprocal of the time span between observations). The estimation was based on 5000 simulated increments with  $\Delta t$  equal to one-eight of a day. The figures seem quite comforting. To get a better idea of the performance of the method, Figures 3 and 4 show the sampling distributions of the estimates of  $\alpha^-$  and  $\beta^+$  obtained from applying the least-square method to the penalized projection estimators. The histograms are based on 1000 samples of size 5000 with  $\Delta t = 1/8$  of a day. This experiment shows clear, though З

 <sup>&</sup>lt;sup>4</sup> More concretely, the density is terms of Bessel special functions of third kind. For more information, see also Section 4.1 in Knotz et al. [21].

not critical, underestimation of the parameter  $\alpha$  and overestimation of the parameters  $\beta$ 's. A simple method of moments (based on the first four moments) yields better results (see Figures 5 and 6). Nonparametric methods are not free-lunches and usually the gain in robustness is paid by a loss in efficiency.

To illustrate the seriousness of applying an efficient estimation method to a misspecified model let us consider a close relative of the variance Gamma process: the CGMY model in [6]. This is defined as a pure-jump Lévy process with Lévy density of the form

$$e_{-}(x) = \int \frac{\alpha^{-}}{|x|^{\nu+1}} \exp\left(-\frac{|x|}{\beta^{-}}\right) \quad \text{if } x < 0,$$

(6.6)

 $s_m(x) = \begin{cases} \frac{\alpha^+}{x^{\nu+1}} \exp\left(-\frac{x}{\beta^+}\right) & \text{if } x > 0, \end{cases}$ where  $\nu > 0$ . In the case when  $\alpha^- = 0$  and  $\nu = 0$ , we recover a Gamma Lévy process, for which MLE are widely available. Let us take  $\alpha^+ = \beta^+ = 1$  and  $\nu = .1$ .

We can estimate the parameter v using a Zolotarev type estimator. This can be done so since the CGMY Lévy process is a tempered stable Lévy process, whose short-term increments behave like stable processes (see Rosínski [27] for details).

The following table shows the sampling average and standard deviations of the estimators of  $\alpha^+$ ,  $\beta^+$ , and v by two methods based on 100 simulation runs. The first method estimates v using the Zolotarev's estimator  $\hat{v}$ , then computes the piece-wise constant p.p.e.  $\tilde{s}$  of (5.5), and finally, estimate  $\alpha^+$  and  $\beta^+$  via the LSE method (6.5) replacing  $s^{\theta}$  by the Lévy density  $s_m$  of (6.6) with  $\theta = (\alpha^+, \beta^+)$  and fixing  $\alpha_- = 0$ and  $v = \hat{v}$ . The second method assumes (erroneously) that the underlying model is a Lévy gamma process and performs maximum likelihood estimation.

| ſ | $\Delta t$ | Penalized Projection Estimators/Least-Squares Fit |                      |                       | Misspecified Gamma MLE |                       |  |
|---|------------|---|----------------------|-----------------------|------------------------|-----------------------|--|
| ſ |            | $\hat{\alpha}^+_{NP}$                             | $\hat{\beta}^+_{NP}$ | $\hat{v}_{Zolotarev}$ | $\hat{\alpha}^+_{MLE}$ | $\hat{\beta}^+_{MLE}$ |  |
|   | .01        | <b>1.03</b> (0.15)                                | <b>0.97</b> (0.14)   | $0.09 \ (0.0002)$     | <b>1.2</b> (0.08)      | <b>0.89</b> (0.079)   |  |
|   |            |   |                      | TABLE 1               |                        |                       |  |

Sampling mean in **bold** and standard errors in parenthesis of the estimators of  $\alpha^+$ ,  $\beta^+$ , and  $\upsilon$ in the CGMY model with theoretical values  $\alpha^- = 0$ ,  $\beta^+ = 1$ ,  $\alpha^+ = 1$ , and v = .1. Sample size is 100 paths.

The results above shows that sometime a modestly efficient robust nonparametric method is preferably to a very efficient estimation method.

#### **Appendix A: Technical Proofs**

Proof of Lemma 3.1. The idea is to exploit the well-known decomposition of the Lévy process as a compound Poisson process  $X^{\varepsilon}$  plus an independent Lévy process  $X^{\varepsilon} := X - \widetilde{X}^{\varepsilon}$  with compactly supported Lévy measure  $\nu_{\varepsilon}(dx) := \mathbf{1}_{\{|x| \le \varepsilon\}} \nu(dx),$ for a suitable chosen  $\varepsilon > 0$ . Concretely, here

$$\widetilde{X}_t^{\varepsilon} = \sum_{i=1}^{N_t} \xi_i,$$

for a homogeneous Poisson process  $\{N_t\}_{t>0}$  with intensity  $\lambda_{\varepsilon} := \nu(\{|x| > \varepsilon\})$  and for independent random variables  $\{\xi_i\}$  with distribution  $\frac{1}{\lambda_{\varepsilon}} \mathbf{1}_{\{|x| > \varepsilon\}} \nu(dx)$ . Clearly,

$$\frac{1}{t} \mathbb{P} \left[ X_t \ge y \right] = \frac{1}{t} \mathbb{P} \left[ X_t^{\varepsilon} \ge y \right] e^{-\lambda_{\varepsilon} t} + \mathbb{P} \left[ X_t^{\varepsilon} + \xi_1 \ge y \right] e^{-\lambda_{\varepsilon} t} (\lambda_{\varepsilon})$$

$$\sum_{n=1}^{\infty} \mathbb{I}\left[ \mathbf{x} \varepsilon + \sum_{n=1}^{n} \varepsilon \right] = \lambda_{n} t \times n (n-1)$$

$$+\sum_{n=2}\mathbb{P}\left[X_t^{\varepsilon}+\sum_{i=1}\xi_i\geq y\right]e^{-\lambda_{\varepsilon}t}\lambda_{\varepsilon}^nt^{n-1}.$$

З

Then, we have

$$\left|\frac{1}{t} \mathbb{P}\left[X_t \ge y\right] - \nu([y,\infty))\right| \le \frac{1}{t} \mathbb{P}\left[X_t^{\varepsilon} \ge y\right] + \left|\lambda_{\varepsilon} \mathbb{P}\left[X_t^{\varepsilon} + \xi_1 \ge y\right] - \nu([y,\infty))\right|$$

 $+\nu([y,\infty))\lambda_{\varepsilon}t+\lambda_{\varepsilon}^{2}t.$ 

The second term on the right hand side of this inequality can itself be decomposed as follows:

$$|\lambda_{\varepsilon} \mathbb{P} \left[ X_t^{\varepsilon} + \xi_1 \ge y \right] - \nu([y, \infty))| \le \int_{y-\eta}^{y+\eta} s(x) \, dx + \lambda_{\varepsilon} \mathbb{P} \left[ |X_t^{\varepsilon}| \ge \eta \right]$$

$$+\nu([y,\infty))\mathbb{P}\left[|X_t^{\varepsilon}| \ge \eta\right]$$

for each  $\eta > 0$  such that  $a - \eta > \varepsilon$ . Since s is bounded off the origin, there exists a k > 0 such that

$$\int_{y-\eta}^{xy+\eta} s(x)dx \le k\eta$$
15
16

for all  $y \in D$ . Fix  $0 < \eta < \frac{1}{T} \land a$  and  $0 < \varepsilon < a - \eta$ . Then,

$$\sup_{y \in D} \left| \frac{1}{t} \mathbb{P} \left[ X_t \ge y \right] - \nu([y,\infty)) \right| \le \frac{1}{t} \mathbb{P} \left[ X_t^{\varepsilon} \ge a \right] + \frac{k}{T} + \nu([a,\infty)) \mathbb{P} \left[ |X_t^{\varepsilon}| \ge \eta \right]$$

$$+\lambda_{\varepsilon} \mathbb{P}\left[|X_t^{\varepsilon}| \ge \eta\right] + \nu([a,\infty)) \lambda_{\varepsilon} t + \lambda_{\varepsilon}^2 t.$$

Finally, since  $\lim_{t\to 0} \frac{1}{t} \mathbb{P} \left[ X_t^{\varepsilon} \ge a \right] = 0$  and  $\lim_{t\to 0} \mathbb{P} \left[ |X_t^{\varepsilon}| \ge \eta \right] = 0$ , we can choose  $\delta_T > 0$  sufficiently small to make each of the terms smaller than 1/T when  $t < \delta_T$ .

Proof of Theorem 4.1.

(i) Fix a Lévy density  $s_0 \in \Theta_\alpha(L/2; [a, b])$  such that  $s_0(x) > 0$ , for all  $x \in \mathbb{R}_0 :=$  $\mathbb{R}\setminus\{0\}$ , and a constant  $\kappa > 0$ . Also, let  $g: \mathbb{R} \to \mathbb{R}_+$  be a symmetric function with compact support  $\mathbb{K}_g$ , satisfying (3.5) with L/2 (instead of L). Moreover, the support of  $x \to g(\kappa(x-x_0))$ , denoted by K, does not contain the origin and also,

$$s_0(x) - \kappa^{-\alpha} g\left(\kappa(x - x_0)\right) > 0, \quad \forall x \in \mathbb{R}_0.$$

Let

$$s_{\theta}(x) := s_0(x) + \theta T^{-\frac{\alpha}{2\alpha+1}} g\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right), \ x \in \mathbb{R}_0,$$

and notice that  $s_{\theta} \in \Theta$  whenever  $|\theta| < \kappa^{-\alpha}$ .

(ii) Without loss of generality we assume that  $\mathbb{K} \cap [-1,1] = \emptyset$ . We follow the notation in [29] (Section 33). Let  $\mathbb{P}_{\theta}^{(T)}$  be the distribution (on  $\mathbb{D}[0,T]$ ) of a Lévy process  $\{X(t)\}_{0 \le t \le T}$  with Lévy density  $s_{\theta}$  (the other two parameters of the generating triplet remain constant). We proceed to prove that  $\left\{ \mathbb{P}_{\theta}^{(T)} : \theta \in (-\kappa^{-\alpha}, \kappa^{-\alpha}) \right\} \text{ is LAN at } \theta = 0 \text{ (see e.g. Definition II.2.1 in [19].}$ By Theorems 33.1 and 33.2 in [29],  $\mathbb{P}_{\theta}^{(T)} \approx \mathbb{P}_{0}^{(T)}$  and the likelihood function,  $J \equiv (T)$ 

$$L_{\theta}(\omega) := \frac{d\mathbb{P}_{\theta}^{(T)}}{d\mathbb{P}_{0}^{(T)}}(\omega)$$
 is given by

$$\frac{\alpha}{2} = \frac{\alpha}{2}$$

$$L_{\theta}(\omega) := \exp\left\{\int_{0}^{T} \int_{\mathbb{K}} \ln\left[1 + \frac{\theta T^{-\frac{\alpha}{2\alpha+1}}}{s_{0}(x)}g\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_{0})\right)\right]\xi(dt, dx; \omega)$$

$$- \theta T^{1-\frac{\alpha}{2\alpha+1}} \int_{\mathbb{K}} g\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right) dx \bigg\},$$

imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009

З

where  $\xi(dt, dx; \omega)$  is the random measure on  $\mathbb{R}_+ \times \mathbb{R}_0$  associated with the jumps of  $\omega \in \mathbb{D}[0,T]$ ; that is,

$$\xi(A;\omega):=\#\{(t,x):\Delta w_t:=w_t-w_{t^-}=x\},\quad A\subset\mathbb{R}_+\times\mathbb{R}_0.$$

Under  $\mathbb{P}_0^{(T)}$ ,  $\xi$  is a Poisson random measure with mean measure  $s_0(x)dxdt$ . We denote  $\bar{\xi}(dt, dx; \omega) := \xi(dt, dx; \omega) - s_0(x)dxdt$ . The likelihood  $L_{\theta}(\omega)$  can be written as follows:

$$L_{ heta}(\omega) = \exp\left\{ heta\Delta_{_T} - rac{ heta^2}{2}\sigma_{_T}^2 + r_{_T}( heta)
ight\},$$

where

$$\Delta_T = T^{-\frac{\alpha}{2\alpha+1}} \int_0^T \int_{\mathbb{K}} s_0^{-1}(x) g\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right) \bar{\xi}(dt, dx),$$

$$\sigma_T^2 = T^{1 - \frac{2\alpha}{2\alpha + 1}} \int_{\mathbb{K}} s_0^{-1}(x) g^2 \left(\kappa T^{\frac{1}{2\alpha + 1}}(x - x_0)\right) dx,$$

$$r_{\scriptscriptstyle T}(\theta) = -\frac{\theta^2}{2} T^{-\frac{2\alpha}{2\alpha+1}} \int_0^T \!\!\!\!\int_{\mathbb{K}} s_0^{-2}(x) g^2 \left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right) \bar{\xi}(dt,dx)$$

$$+ \int_0^T \!\!\!\int_{\mathbb{K}} R\left(\theta T^{-\frac{\alpha}{2\alpha+1}} s_0^{-1}(x) g\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right)\right) \xi(dt,dx),$$

and  $R(u) := \ln(1+u) - u + \frac{u^2}{2}$ . We want to prove that there are nomalizing constants  $\varphi_T > 0$  such that

$$\mathcal{L}_{\mathbb{P}_{0}^{(T)}}\left(\varphi_{T}\Delta_{T}\right) \xrightarrow{\mathfrak{D}} \mathcal{N}(0,1), \ \varphi_{T}^{2}\sigma_{T}^{2} \to 1, \ \text{and} \ r_{T}(\theta) \xrightarrow{\mathbb{P}_{0}^{(T)}} 0$$

as  $T \to \infty$ . To prove the first limit, we invoke the CLT for Poisson integrals by verifying the Liapunov condition (see Theorem 1.1 and Remark 1.2 of [22]). Indeed, for T > 1, we have that

$$T^{-\frac{\alpha(2+\delta)}{2\alpha+1}} \int_0^T \int_{\mathbb{K}} s_0^{-2-\delta}(x) g^{2+\delta} \left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right) (s_0(x)) \, dx dt =$$

$$\kappa^{-1}T^{1-\frac{\alpha(2+\delta)}{2\alpha+1}-\frac{1}{2\alpha+1}} \int_{\mathbb{K}_g} s_0^{-1-\delta} (\kappa^{-1}T^{-\frac{1}{2\alpha+1}}u + x_0) g^{2+\delta}(u) \, du \xrightarrow{T \to \infty} 0.$$

Similarly, for large enough T,

$$\stackrel{T \to \infty}{\longrightarrow} \kappa^{-1} s_0^{-1}(x_0) \int_{\mathbb{K}_g} g^2(u) du.$$

Then,  $\mathcal{L}_{\mathbb{P}_0^{(T)}}(\Delta_T) \xrightarrow{\mathfrak{D}} \mathcal{N}(0, I_0^2)$  with  $I_0^2 := \kappa^{-1} s_0^{-1}(x_0) \int_{\mathbb{K}_q} g^2(u) du$ , and  $\sigma_T^2 \to 0$  $I_0^2$ . We now verify that  $r_{\tau}(\theta)$  vanishes in probability. Notice that the first term of  $r_{\tau}$  converges to 0 since its mean is 0 and its variance vanishes. Similarly, the second term of  $r_T(\theta)$  converges to 0 in probability because its mean and variance both goes to 0. Indeed, using that  $|R(u)| \leq |u|^3/3$ , the 

З

absolute value of its expectation satisfies  $\left| \int_0^T \int_{\mathbb{K}} R\left(\theta T^{-\frac{\alpha}{2\alpha+1}} s_0^{-1}(x) g\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right) \right) (s_0(x)) dx dt \right|$  $\leq \frac{|\theta|^3}{3} T^{1-\frac{3\alpha}{2\alpha+1}} \int_{\mathbb{T}} s_0^{-2}(x) g^3\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_0)\right) dx \xrightarrow{T \to \infty} 0.$ A similar reasoning applies to the variance. Therefore,  $\{\mathbb{P}_{\theta}^{(T)}\}_{\theta \in (-\kappa^{-\alpha}, \kappa^{-\alpha})}$ is Locally Asymptotically Normal (LAN) at  $\theta = 0$  (with the normalizing constants  $\varphi_T := I_0^{-1}$ ). (iii) By Theorem II.12.1 and Remark II.12.2 in [19], if  $\hat{\theta}_T$  is any estimator of  $\theta$ based on  $\{X(t)\}_{0 \le t \le T}$ , then  $\liminf_{T \to \infty} \sup_{|\theta| < \kappa^{-\alpha}} \mathbb{E}_{\theta} \left[ \ell_0 \left( I_0 \left( \hat{\theta}_T - \theta \right) \right) \right] \ge B,$ (A.1)where  $B := \mathbb{E}\left[\ell_0(Z)\chi_{[|Z| < I_0 \kappa^{-\alpha}/2]}\right]$  and  $Z \sim \mathcal{N}(0, 1)$ . Now, let  $\hat{s}_{_T}(\cdot)$  be an arbitrary estimator based on  $\{X(t)\}_{0 \le t \le T}$  and let  $\hat{\theta}_{\tau} := T^{\frac{\alpha}{2\alpha+1}} q^{-1}(0) \left( \hat{s}_{\tau}(x_0) - s_0(x_0) \right).$ Since  $\theta = T^{\frac{\alpha}{2\alpha+1}} g^{-1}(0) \left( s_{\theta}(x_0) - s_0(x_0) \right)$ , we can write  $g(0)\left(\hat{\theta}_{T}-\theta\right) = T^{\frac{\alpha}{2\alpha+1}}\left(\hat{s}_{T}(x_{0})-s_{\theta}(x_{0})\right).$ If we take  $\ell_0(u) := \ell(g(0)I_0^{-1}u), (A.1)$  becomes:  $B \leq \liminf_{T \to \infty} \sup_{|\theta| < \kappa^{-\alpha}} \mathbb{E}_{\theta} \left[ \ell_0 \left( I_0 \left( \hat{\theta}_T - \theta \right) \right) \right]$  $= \liminf_{T \to \infty} \sup_{|\theta| < \kappa^{-\alpha}} \mathbb{E}_{\theta} \left[ \ell \left( T^{\frac{\alpha}{2\alpha+1}} \left( \hat{s}_{T}(x_{0}) - s_{\theta}(x_{0}) \right) \right) \right].$ Since  $\{s_{\theta}: \theta \in (-k^{-\alpha}, k^{-\alpha})\} \subset \Theta$ ,  $\liminf_{T \to \infty} \sup_{s \in \Theta} \mathbb{E}_s \left[ \ell \left( T^{\frac{\alpha}{2\alpha+1}} \left( \hat{s}_T(x_0) - s(x_0) \right) \right) \right] \ge B,$ (A.2)where  $B := 2^{-3/2} \pi^{-1/2} \int_{|z| < I_0 \kappa^{-\alpha}/2} \ell(g(0)I_0^{-1}z) e^{-z^2/2} dz.$ (A.3)This implies (4.1) because the lower bound B does not depend on the family of estimators  $\hat{s}_{\tau}$ . Indeed, for each  $\varepsilon > 0$ , let  $\hat{s}_{\tau}^{(\varepsilon)}$  be such that  $\sup_{s \in \Theta} \mathbb{E}_s \left[ \ell \left( T^{\frac{\alpha}{2\alpha+1}} \left( \hat{s}_T^{(\varepsilon)}(x_0) - s(x_0) \right) \right) \right]$  $< \inf_{\hat{s}_T} \sup_{s \in \Theta} \mathbb{E}_s \left[ \ell \left( T^{\frac{\alpha}{2\alpha+1}} \left( \hat{s}_T(x_0) - s(x_0) \right) \right) \right] + \varepsilon.$ Taking the lim inf as  $T \to \infty$  on both sides, we obtain (4.1) since  $\varepsilon$  is arbitrary. Proof of Corollary 4.2. imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009

(i) Following the same reasoning as in Theorem 4.1, we first prove that for any family of estimators  $\{\hat{s}_T\}_{T>0}$  and arbitrary points  $\{x_T\}_{T>0} \subset (a, b)$ ,

$$(A.4) \qquad \liminf_{T \to \infty} \sup_{s \in \Theta} \mathbb{E}_s \left[ \ell \left( T^{\frac{\alpha}{2\alpha+1}} \left( \hat{s}_{_T}(x_{_T}) - s(x_{_T}) \right) \right) \right] \ge C$$

for some constant C > 0, which is independent of the family of estimators and of the points. Fix a Lévy density  $s_0 \in \Theta_{\alpha}(L/2; [a, b])$  such that  $s_0(x) > 0$ for all  $x \in \mathbb{R}_0 := \mathbb{R} \setminus \{0\}$ , and a constant  $\kappa > 0$ . Again, let  $g : \mathbb{R} \to \mathbb{R}_+$ be a symmetric function with compact support  $\mathbb{K}_g$ , satisfying (3.5) with L/2(instead of L). Moreover, for any  $y \in (a, b)$ , the support of  $x \to g(\kappa(x - y))$ does not contain the origin and

$$s_0(x) - \kappa^{-\alpha} g\left(\kappa(x-y)\right) > 0, \quad \forall x \in \mathbb{R}_0.$$

Let

$$s_{\theta,T}(x) := s_0(x) + \theta T^{-\frac{\alpha}{2\alpha+1}} g\left(\kappa T^{\frac{1}{2\alpha+1}}(x-x_T)\right), \ x \in \mathbb{R}_0,$$

for  $|\theta| < \kappa^{-\alpha}$ . Let  $\mathbb{P}_{\theta}^{(T)}$  be the distribution (on  $\mathbb{D}[0,T]$ ) of a Lévy process  $\{X(t)\}_{0 \le t \le T}$  with Lévy density  $s_{\theta,T}$ . Following the proof of Theorem 4.1,  $\{\mathbb{P}_{\theta}^{(T)} : \theta \in (-\kappa^{-\alpha}, \kappa^{-\alpha})\}$  is LAN at  $\theta = 0$  with the normalizing constants

$$\varphi_{\scriptscriptstyle T} := \kappa^2 \left( \int_{\mathbb{K}_g} s_0^{-1} (\kappa^{-1} T^{-\frac{1}{2\alpha+1}} u + x_{\scriptscriptstyle T}) g^2 \left( u \right) du \right)^{-2},$$

where  $\mathbb{K}_g$  denotes the support of g and it is being assumed that  $[-1,1] \cap \cup_{y \in [a,b]} \{y + \kappa^{-1} \mathbb{K}_g\} = \emptyset$ . Observe that there is an m > 0 for which  $\inf_{T \ge 1} \varphi_T \ge m$ .

(ii) By Theorem II.12.1 and Remark II.12.2 in Ibragimov & Has'minskii (1981), for any  $\delta > 0$ ,

(A.5) 
$$\liminf_{T \to \infty} \sup_{|\theta| < \delta \varphi_T} \mathbb{E}_{\theta} \left[ \ell_0 \left( \varphi_T^{-1} \left( \hat{\theta}_T - \theta \right) \right) \right] \ge C,$$

where  $C := \mathbb{E}\left[\ell_0(Z)\chi_{[|Z|<\delta/2]}\right]$  and  $Z \sim \mathcal{N}(0,1)$ . Since  $\ell_0(|y|)$  is increasing in y,

(A.6) 
$$\liminf_{T \to \infty} \sup_{|\theta| < \delta \varphi_T} \mathbb{E}_{\theta} \left[ \ell_0 \left( m^{-1} \left( \hat{\theta}_T - \theta \right) \right) \right] \ge C.$$

Now, setting,

$$\hat{\theta}_{_{T}} := T^{\frac{\alpha}{2\alpha+1}} g^{-1}(0) \left( \hat{s}_{_{T}}(x_{_{T}}) - s_0(x_{_{T}}) \right),$$

it follows that

$$\sup_{s\in\Theta} \mathbb{E}_s \left[ \ell \left( T^{\frac{\alpha}{2\alpha+1}} \left( \hat{s}_{_T}(x_{_T}) - s(x_{_T}) \right) \right) \right] \geq \sup_{|\theta| < \delta \varphi_T} \mathbb{E}_{\theta} \left[ \ell \left( g(0) \left( \hat{\theta}_{_T} - \theta \right) \right) \right].$$

Taking limit as  $T \to \infty$ , (A.4) is obtained with

(A.7) 
$$C = 2^{-3/2} \pi^{-1/2} \int_{|z| < \delta/2} \ell(g(0) \, m \, z) e^{-z^2/2} dz.$$

J.E. Figueroa-López

142



imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009



imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009



imsart-coll ver. 2008/08/29 file: Figueroa.tex date: March 25, 2009



FIG 6. Sampling Distribution for the Estimator of  $\beta^+$  obtained by the Method of Moments.



### References

З

| 24 |       |  | 24 |
|----|-------|--|----|
| 25 | [1]   | BARNDORFF-NIELSEN O. E. and SHEPHARD N. (2001). Modelling by Lévy processess for financial eco-  | 25 |
| 26 |       | nomics. Levy Processes. Theory and Applications by Barndorff-Nielsen, O.E., Mikosch, T. and<br>Resnick S I 283–318                                 | 26 |
| 27 | [2]   | BARRON A., BIRGÉ L., and MASSART P. (1995). Model selection via penalization. <i>Technical Report</i> ,  | 27 |
| 28 |       | 54, Université Paris-Sud.  | 28 |
| 29 | [3]   | BARRON A., BIRGÉ L., and MASSART P. (1999). Risk bounds for model selection via penalization.<br>Probability Theory and Related Fields 113 301-413 | 29 |
| 30 | [4]   | BIRGÉ L., and MASSART P. (1994). Minimum contrast estimation on sieves. <i>Technical Report</i> , 34,  | 30 |
| 31 |       | Université Paris-Sud.  | 31 |
| 32 | [5]   | BIRGÉ L., and MASSART P. (1997). From model selection to adaptive estimation. In <i>Festschrift for Lucien Le Cam</i> , 55–87.                     | 32 |
| 33 | [6]   | CARR P., GEMAN, H., MADAN, D., and YOR M. (2002). The fine structure of asset returns: An empirical  | 33 |
| 34 |       | investigation. Journal of Business, 305–332.   | 34 |
| 35 | [7]   | CARR P., GEMAN, H., MADAN, D., and YOR M. (2003). Stochastic volatility for Lévy processes. In <i>Mathematical Finance</i> , <b>13</b> , 345–382.  | 35 |
| 36 | [8]   | CARR P., MADAN, D., and CHANG E. (1998). The variance Gamma process and option pricing. Eu-  | 36 |
| 37 |       | ropean Finance Review, 2, 79–105.  | 37 |
| 38 | [9]   | CARR P., and WU, L. (2005). Time-changed Levy processes and option pricing. <i>Journal of Financial</i>  | 38 |
| 39 | [10]  | CHUNG K. I. (2001) A course in Probability Theory Academic Press   | 39 |
| 10 | [11]  | CONT R., and TANKOV P. (2003). Financial modelling with Jump Processes. Chapman & Hall.  | 40 |
| 40 | [12]  | DEVORE R. A., and LORENTZ G. G. (1993). Constructive Approximation. Springer-Verlag.   | 40 |
| 41 | [13]  | FIGUEROA-LÓPEZ J. E. (2004). Nonparametric estimation of Lévy processes with a view towards  | 41 |
| 42 |       | mathematical finance. PhD Thesis. Georgia Institute of Technology, Atlanta, GA 30332, April  | 42 |
| 43 | [4.4] | 2004. http://etd.gatech.edu.No.etd-04072004-122020.  | 43 |
| 44 | [14]  | FIGUEROA-LOPEZ J. E. (2007). On the non-parametric estimation of some time-changed Levy models.<br>In preparation.                                 | 44 |
| 45 | [15]  | FIGUEROA-LÓPEZ J. E. (2008). Small-time moment asymptotics for Lévy processes. Statistics and  | 45 |
| 46 |       | Probability Letters. DOI:10.1016/j.spl.2008.07.012.  | 46 |
| 47 | [16]  | FIGUEROA-LÓPEZ J. E. and HOUDRÉ C. (2006). Risk bounds for the non-parametric estimation of Lévy   | 47 |
| 18 | [17]  | processes. IMS Lecture Notes - Monograph Series. High Dimensional Probability, <b>51</b> , 96–116.   | 19 |
| 10 | [1/]  | of Lévy processes. Available at arXiv:0809.0849v1 [math.PR].   | 40 |
| 49 | [18]  | GRENANDER U. (1981). Abstract Inference. John Wiley & Sons.  | 49 |
| 50 | [19]  | IBRAGINOV I. A. and HAS'MINSKII R. Z. (1981). Statistical Estimation. Asymptotic Theory. Springer-   | 50 |
| 51 |       | Verlag, Berlin, New York, Heidelberg.  | 51 |

З

| J.E. Figueroa-López |  |
|---------------------|--|
|---------------------|--|

| 1        | [20] | JACOD J. (2007). Asymptotic properties of power variations of Lévy processes. <i>ESAIM:P&amp;S</i> , <b>11</b> , 172, 106   | 1      |
|----------|------|---|--------|
| 2        | [21] | KNOTZ S., KOZUBOWSKI T. J., and PODGÓRSKI K. (2001). The Laplace Distribution and Generaliza-   | 2      |
| 3<br>4   |      | tions: A revisit with Applications to Communications, Economics, Engineering, and Finance.<br>Birkhauser Boston   | 3<br>4 |
| 5        | [22] | KUTOYANTS Y. A. (1998). Statistical Inference for Spatial Poisson Processes. Springer.  | 5      |
| 6        | [23] | MADAN D. B. and SENETA E. (1990). The variance Gamma model for share market returns. <i>Journal</i> of <i>Business</i> <b>63</b> , 511–524  | 6      |
| 7        | [24] | PRAUSE K. (1999). The generalized hyperbolic model: estimation, financial derivatives, and risk   | 7      |
| 8        |      | measures. Ph.D. Thesis. University of Freiburg.   | 8      |
| 9        | [25] | REYNAUD-BOURET P. (2003). Adaptive estimation of the intensity of inhomogeneous Poisson pro-  | 9      |
| 10       | [26] | Rosiński J. (2001). Series representations of Lévy processes from the perspective of point processes.   | 10     |
| 11       |      | In Lévy processes-Theory and Applications, 401–415.   | 11     |
| 12       | [27] | ROSIŃSKI J. (2007). Tempering stable processes. Stochastic processes and their applications, 117, 677–707   | 12     |
| 13       | [28] | RÜSCHENDORF L. and WOERNER J. (2002). Expansion of transition distributions of Lévy processes   | 13     |
| 14       | [00] | in small time. Bernoulli, <b>8</b> , 81–96.   | 14     |
| 15       | [29] | SATO K. (1999). Levy Processes and Infinitely Divisible Distributions. Cambridge University Press.<br>SCHOUTENS W (2003) Lévy Processes: Pricing Financial Derivatives John Wiley | 15     |
| 16       | [31] | TODOROV V. (2005). Econometric analysis of jump-driven stochastic volatility models. Duke Uni-  | 16     |
| 17       |      | versity, Working Paper.   | 17     |
| 18       | [32] | WOERNER J. (2003). Variational sums and power variation: a unifying approach to model selection   | 18     |
| 19       |      | and commation in semimartingate models. Diatistics and Decisions, 21, 47 00.  | 19     |
| 20       |      |   | 20     |
| 21       |      |   | 21     |
| 22       |      |   | 22     |
| 23       |      |   | 23     |
| 24       |      |   | 24     |
| 25       |      |   | 25     |
| 26       |      |   | 26     |
| 27       |      |   | 27     |
| 28       |      |   | 28     |
| 29       |      |   | 29     |
| 30       |      |   | 30     |
| 31       |      |   | 31     |
| 32       |      |   | 32     |
| 33       |      |   | 33     |
| 34       |      |   | 34     |
| 35       |      |   | 35     |
| 27       |      |   | 30     |
| 20       |      |   | 31     |
| 30<br>20 |      |   | 30     |
| 39<br>40 |      |   | 39     |
| 40       |      |   | 40     |
| 40       |      |   | 41     |
| 43       |      |   | 43     |
| 44       |      |   | 44     |
| 45       |      |   | 45     |
| 46       |      |   | 46     |
| 47       |      |   | 47     |
| 48       |      |   | 48     |
| 49       |      |   | 49     |
| 50       |      |   | 50     |
| 51       |      |   | 51     |

|  | On the Estimation of Symmetric  |
|--|---|
| D  | istributions under Peakedness Order   |
|  | Constraints   |
|  | Javier Rojo <sup>1,*</sup> and José Batún-Cutz  |
|  | Rice Universitu and Universidad Autónoma de Yucatán. Mérida. México   |
|  | Abstract: Consider distribution functions $F$ and $G$ and suppose that $F$ is more peaked about $a$ than $G$ is about $b$ . The problem of estimating $F$ or $G$ , or both, when $F$ and $G$ are symmetric, arises quite naturally in applications. The empirical distribution functions $F_n$ and $G_m$ will not necessarily satisfy the order constraint imposed by the experimental conditions. Rojo and Batun-Cutz (2007) proposed some estimators that are strongly uniformly consistent when both $m$ and $n$ tend to infinity. However the estimators fail to be consistent when only either $m$ or $n$ tend to infinity. Here estimators are proposed that circumvent these problems and the asymptotic distribution of the estimators is delineated. A simulation study compares these estimators in terms of Mean Squared Error and Bias behavior with their competitors. |
| Conte  | nts   |
| <ol> <li>Intt</li> <li>Nev</li> <li>2.1</li> <li>2.2</li> <li>2.3</li> <li>2.4</li> <li>Asy</li> <li>3.1</li> <li>3.2</li> <li>4 Exa</li> <li>5 Sim</li> <li>6 Con</li> <li>Acknow</li> <li>Reference</li> </ol> | roduction   |
| L. Int<br>The co<br>tions t  | roduction<br>oncept of stochastic order was pioneered by Lehmann (1955), and applica<br>o hypotheses testing were discussed in Lehmann (1959), henceforth referred<br>FSH-1. Lehmann and Rojo (1992) provided characterizations of stochastic   |

imsart-coll ver. 2008/08/29 file: Rojo\_Batun.tex date: March 25, 2009

ordering in terms of the maximal invariant with respect to the group of mono-tone transformations, and connections with other partial orderings were provided. Since the publication of TSH-1, there has been a large number of papers discussing various types of stochastic orders and their properties. Thus, one finds a large liter-ature on stochastic orders in Economics (e.g. first-, second-, third-order stochastic dominance), reliability (e.g. IFR, IFRA, NBU, etc.), and applied probability (e.g. Laplace transform and dispersive orders). Marshall and Olkin (2007) and Shaked and Shantikumar (2007) are excellent references to the literature on stochastic or-ders.

The attention to this area of statistics and applied probability is well deserved. These concepts arise naturally in many applications in engineering, survival analysis, biology, economics, etc.

In corrosion engineering, for example, the times until pitting of metals immersed in a corrosive environment are measured under different solution corrosivities to discern the impact of the solution acidity on the pitting corrosion times. Shibata and Takeyama (1977) present data which strongly supports the belief that the times until pitting should be shorter in some sense, for the more corrosive environment. In toxicity studies, cells are grown in environments containing different levels of toxic materials (e.g. Arenaz et al (1992)). Invariably, the data supports the intu-itive notion that the stronger the toxic solution is, the shorter the lifetimes of the organisms.

Another set of examples arises from clinical trials. This is illustrated by a clinical trial run to evaluate the efficiency of maintenance chemotherapy for acute myelon-geneous leukemia (AML). The trial was conducted at Stanford University (Embury et al (1977)). After reaching a state of remission through treatment by chemother-apy, the patients who entered the study were randomized into two groups. The first group received maintenance chemotherapy; the second group did not. One would then expect that in this case, the survival times in the control group would be stochastically smaller than those in the first group.

Stochastic ordering, together with failure rate ordering, and monotone likelihood
 ratio ordering, are examples of *location* orderings. There are situations, however,
 when the interest lies in comparing distributions based on their *spread* rather than
 on their location.

Various concepts of spread, concentration, or dispersion have appeared in the literature. For example, Brown and Tukey (1946), Fraser (1957), Bickel and Lehmann (1979), Lehmann (1988), Doksum (1969), and Shaked (1980), define F to be more dispersive than G, denoted as  $F >_d G$ , if, for every u > v,

(1.1) 
$$F^{-1}(u) - F^{-1}(v) \ge G^{-1}(u) - G^{-1}(v).$$

Shaked (1982), Bartoszewicz (1985a, 1985b, 1986), Oja (1981), and Rojo and He (1991), among others, have discussed various characterizations and properties of the dispersive order. Doksum (1969) utilized this concept to study power properties of rank tests, and showed that the power of certain rank tests is isotonic with respect to this order. Rojo (1995b, 1999) considered the problem of estimating the quantile function  $F^{-1}$  and the distribution function F when  $F <_d G$ , and the asymptotic theory of the resulting estimators was delineated. Rojo and Wang (1994) also showed that the power of tests based on L-statistics is isotonic with respect to the dispersive order. For other properties of the dispersive order, and connections with other partial orderings, see Bickel and Lehmann (1979), Proschan (1965), Karlin (1968), Shaked (1980, 1982), and Schweder (1982). When F and G

are assumed symmetric, (1.1) can be seen to be equivalent to

$$F^{-1}(u) - F^{-1}(1/2) \ge (\le) G^{-1}(u) - G^{-1}(1/2)$$

depending on whether  $u \ge (\le) 1/2$ .

Birnbaum (1948) proposed a different concept of dispersion based on the distribution functions rather than on the quantile functions. According to Brinbaum, the distribution function F is more peaked about the point a than the distribution function G is about the point b if, for all  $x \ge 0$ ,

(1.2) 
$$F((x+a)^{-}) - F(-x+a) \ge G((x+b)^{-}) - G(-x+b),$$

where  $h(x^{-}) = \lim_{\epsilon \downarrow 0} h(x - \epsilon)$ . We will write  $F >_p G$  whenever (1.2) holds. It is easy to see that the condition (1.2) is equivalent to

$$F(x^{-}) \ge G(x^{-}) \quad \text{for } x \ge 0$$

(1.3) 
$$\begin{array}{ccc} F(x) \ge G(x) & \text{for } x \ge 0 \\ F(x) \le G(x) & \text{for } x < 0. \end{array}$$

whenever F and G are symmetric about the point 0.

When F and G are continuous, it is easy to see that (1.2) is equivalent to requiring that |X - a| be stochastically smaller than |Y - b|, and, although in general  $F <_d G \not\Rightarrow F >_p G$  and  $F >_p G \not\Rightarrow F <_d G$ , it is easy to verify that  $F <_d G \Rightarrow F >_p G$ , when F and G are symmetric and continuous. When a and b in (1.2) are, respectively, the means of F and G, the condition (1.2) implies the obvious order on the variances of F and G.

An interesting example from statistical genetics, discussed in Rojo *et al* (2007), illustrates the importance of this concept in applications. Haseman-Elston (1972) proposed a regression model to assess the effect of a candidate gene on a phenotype when using sib-paired data. There have been some modifications of the initial proposal (see e.g. Elston et al. (2000)). The original model, Haseman-Elston (1972), represents the expected value of the squared phenotypic differences as a linear function of the proportion of alleles shared identical-by-descent (IBD) at the locus of interest. Let  $\lambda_i$  represent the proportion of alleles shared identical by descent ( $\lambda_i = 0, \frac{1}{2}$ , or 1). The Haseman and Elston (1972) regression model may then be written as follows:  $E(X_i|\lambda_i) = \alpha + \beta \lambda_i$ , where  $X_i$  represents the squared sib-pair difference for the  $i^{th}$  sib-pair conditional on  $\lambda_i$ . Writing  $Z_{1i} = \theta + g_{1i} + \varepsilon_{1i}$  and  $Z_{2i} = \theta + g_{2i} + \varepsilon_{2i}$  where  $Z_{1i}$  and  $Z_{2i}$  represent, respectively, the phenotype values for siblings one and two, and where  $\theta$  is the population mean, and  $g_{ij}$  and  $\varepsilon_{ij}$  are the genetic and the residual effects, respectively, the model is then represented as

$$E(X_j|\lambda_j) = \eta_{\varepsilon}^2 + 2(1-\lambda_j)\eta_g^2$$
<sup>40</sup>
<sup>41</sup>
<sup>42</sup>

where,  $\eta_{\varepsilon}^2 = E((\varepsilon_{1i} - \varepsilon_{2i})^2)$  and  $\eta_g^2$  represents the variance in the trait due to allelic variation at the locus of interest. As a consequence of linkage between the candidate gene and the phenotype, siblings sharing two alleles IBD at the locus of interest will tend to be more similar than siblings sharing one allele IBD, and siblings sharing one allele IBD will in turn be more similar than siblings sharing no alleles IBD. It is then clear that phenotypical similarity of sibs within the same pair is being measured in terms of the spread of the distribution of the differences of the siblings' phenotypical measurements. 

Existing sib-paired data illustrates very clearly that the distribution functions of sib-pair differences are symmetrically distributed. This will happen, for example, if 

З 

 $(X - \mu_X, Y - \mu_Y)$  has the same distribution as  $(\mu_X - X, \mu_Y - Y)$ , as it happens under the assumption of a bivariate normal distribution, and if the means  $\mu_X$  and  $\mu_Y$  are equal, then the sib-pair differences are symmetrically distributed. When the candidate gene is linked to the phenotype of interest, the cumulative distributions of the differences within sib-pairs are ordered by peakedness. This is illustrated by sib-paired data on plasma Lipoprotein (a) data. Figure 1 shows the empirical 

FIG 1. Empirical distribution functions of phenotypic differences for the sib-pair data.

distribution functions for plasma Lipoprotein (a) differences within sib-pairs for a sample of Caucasian individuals from the Dallas metroplex area. The pairs of siblings were classified into groups according to the number of shared alleles identical by descent.

Note that the assumptions of symmetry and peakedness are close to being satisfied, but the plots also show areas where these characteristics do not hold. We will illustrate our estimators later in section 4, by computing them for this example.

The points a and b about which peakedness of F and G will obtain, will be assumed known throughout this work. In the linkage example to be considered in section 4, the assumption of known a and b can be justified under the assumption of bivariate normality of the siblings' phenotypes with equal means. This is a common assumption in the literature. Thus, irrespective of whether a and b are known or unknown, the difference of the phenotypes is always symmetric about zero. Dropping the assumption of bivariate normality of the sib-pairs phenotypes, existing models, see *e.g.* Liu (1988) Table 15.7, yield a zero mean for the phenotypic differences. We, therefore, will assume that a and b are zero.

The goals of this paper are to develop estimators for symmetric F and G, which satisfy (1.2), and to delineate their asymptotic theory.

<sup>50</sup> Under the assumption that F and G are discrete distributions satisfying (1.2), El <sup>51</sup> Barmi and Rojo (1997) provided the nonparametric maximum likelihood estimators <sup>51</sup>

| 1  | of $F$ and $G$ and tests were given to test the hypothesis of homogeneity of $F$ and $G$     |
|----|--|
| 2  | against the alternative that $F$ and $G$ satisfy (1.2). Rojo, Batun, and Durazo (2007)       |
| 3  | proposed estimators for continuous $F$ and $G$ , when (1.2) holds and the case of            |
| 4  | censored data was also considered, but without the symmetry assumption. Rojo and             |
| 5  | Batun-Cutz (2007) proposed estimators for symmetric $F$ and $G$ when $(1.2)$ holds           |
| 6  | using results from Schuster (1975), and the asymptotic theory was delineated for the         |
| 7  | case when both n and $m \to \infty$ . El Barmi and Mukerjee (2008), following the ideas in   |
| 8  | Rojo (2004) and Rojo and Batun (2007), proposed estimators which are consistent              |
| 9  | for $F(G)$ and their asymptotic theory was developed. Unfortunately, the proofs of           |
| 10 | their asymptotic results for the estimators of $F$ and $G$ depend on letting <b>both</b> $n$ |
| 11 | and $m$ increase to infinity. The purpose of this paper is to consider modifications         |
| 12 | of the estimators proposed by Rojo and Batun-Cutz (2007) that yield consistent               |
| 13 | estimators for $F(G)$ when only $n(m) \to \infty$ . The asymptotic distribution theory is    |
| 14 | considered and a simulation study compares the estimators to the estimator of El             |
| 15 | Barmi and Mukerjee (2008).   |
|    |  |

The organization of this paper is as follows: Sections 2 proposes the estimators and finite sample properties are discussed. Section 3 delineates the asymptotic the-ory showing that the estimators are strongly and uniformly consistent and their asymptotic theory is developed. Section 4 illustrates the new estimators using the sib-pair data, and section 5 discusses the results of computer simulations which compare the bias and mean squared error of the new estimators with the bias and mean squared error of the estimators of Rojo and Batun-Cutz (2007) and El Barmi and Mukerjee (2008).

Although the estimators proposed in Rojo and Batun-Cutz (2007) have larger absolute bias than the estimators proposed here, the selection of the better estimators based on Mean Squared Error (MSE) behavior is not as clear. Whereas the new estimators have smaller MSE in a neighborhood of zero, the estimators of Rojo and Batun-Cutz have smaller MSE in the tails of the distributions, and the region of the support of the distribution where the latter estimators behave better seems to increase as the tail-heaviness of the distributions increase.

#### 2. New Estimators and Their Finite Sample Properties

Let  $X_1, \ldots, X_n$  and  $Y_1, \ldots, Y_m$  be independent random samples from the symmetric distributions (about 0) F and G respectively, and let  $F_n$  and  $G_m$  be the empirical distribution functions based on  $X_1, \ldots, X_n$  and  $Y_1, \ldots, Y_m$ . Suppose than  $F >_p G$ . Rojo and Batun-Cutz (2007) considered the problem of the estimation of F and G under the peakedness restriction and proposed the following strongly uniformly consistent estimators

(2.1) 
$$F_{n,m}^1 = \Phi_1(\Phi_2(F_n, \Phi_1(G_m)))$$

(2.2) 
$$F_{n,m}^2 = \Phi_2(\Phi_1(F_n), \Phi_1(G_m)),$$

where  $\Phi_1$  and  $\Phi_2$  are operators defined by

$$\Phi_1(f)(x) = \frac{1}{2}(f(x) + 1 - f(-x^-)), \text{ and}$$

$$\Phi_2(f,g)(x) = \begin{cases} \max\{f(x),g(x)\} & \text{if } x \ge 0\\ \min\{f(x),g(x)\} & \text{if } x < 0. \end{cases}$$

J. Rojo and J. Batún-Cutz

Note that the operator  $\Phi_1$  symmetrizes the function f, Schuster (1975), and the operator  $\Phi_2$  imposes the "stochastic order" restriction (see, e.g., Lo (1987), Rojo and Ma (1996), and Rojo (2004)). Unfortunately the estimators  $F_{n,m}^i$ , for i = 1, 2do not converge to F when only  $n \to \infty$ . This follows since, for example, for  $F_{n,m}^2$ when x > 0 and  $\varepsilon > 0$ , 

$$\lim_{n \to \infty} P[F_{n,m}^2(x) - F(x) > \varepsilon] \ge P[\Phi_1(G_m(x)) - F(x) > \varepsilon] > 0.$$

This is a drawback of  $F_{n,m}^2$  that is also shared by  $F_{n,m}^1$  and  $G_{n,m}^i$  for i = 1, 2,and the strong uniform consistency of these estimators requires that both m and ntend to infinity. To circumvent this problem, new estimators are proposed here.

# 2.1. Definition of the New Estimators

Let  $\widehat{F}_n = \Phi_1(F_n)$  and  $\widehat{G}_m = \Phi_1(G_m)$  be the symmetrized empirical distribution functions (Schuster, 1975). Then the empirical distribution function, and the symmetrized empirical distribution function of the combined samples are defined as follows:

$$C_{n,m} = \frac{n}{m+n}F_n + \frac{m}{n+m}G_m$$
 and <sup>19</sup>

(2.3) 
$$\widehat{C}_{n,m} = \Phi_1(C_{n,m}) = \frac{n}{m+n}\widehat{F}_n + \frac{m}{n+m}\widehat{G}_m$$

respectively. Then our new estimators for F and G are

25 (2.4) 
$$\widehat{F}_{n,m}^1 = \Phi_1(\Phi_2(F_n, C_{n,m})),$$
 25 (2.4)

(2.5) 
$$\widehat{G}_{n,m}^1 = \Phi_1(\Phi_2^*(G_m, C_{n,m})),$$

32 (2.6) 
$$\widehat{F}_{n,m}^2 = \Phi_2(\Phi_1(F_n), \Phi_1(C_{n,m})), \text{ and}$$
  
33

where

Note that the estimators  $\widehat{F}_{n,m}^1$  and  $\widehat{G}_{n,m}^1$  first impose the constraint of "stochastic order" by requiring that the estimator of F(G) be larger (smaller) than  $C_{n,m}$ for  $x \ge 0$  and smaller (larger) than  $C_{n,m}$  for x < 0. The second requirement of symmetry is then imposed by the operator  $\Phi_1$ . By contrast, the estimators  $\hat{F}_{n,m}^2$ and  $\widehat{G}_{n,m}^2$ , first impose the constraint of symmetry and then, through the operator  $\Phi_2$ , the constraint of "stochastic order" is imposed.

El Barmi and Mukerjee (2008) proposed estimators for F and G when  $F <_p G$ . In our notation and making the appropriate change for the case  $F >_p G$ , their estimator for F is given, for  $x \ge 0$ , by

$$F_{nm}^{*}(x) = \frac{1}{2}(1 + \max\left\{F_{n}(x) - F_{n}(-x^{-}), C_{nm}(x) - C_{nm}(-x^{-})\right\}).$$

imsart-coll ver. 2008/08/29 file: Rojo\_Batun.tex date: March 25, 2009

This estimator is the same as our estimator  $\widehat{F}_{n,m}^2$  since for  $x \ge 0$ ,  $\widehat{F}_{n,m}^2(x) = \max \left\{ \frac{1}{2} (1 + F_n(x) - F_n(-x^-)), \frac{1}{2} (1 + C_{nm}(x) - C_{nm}(-x^-)) \right\}$ 

$$= \frac{1}{2} + \frac{1}{2} \max \left\{ F(x) - F(-x^{-}) C(x) - C(-x^{-}) \right\}$$

$$= \frac{1}{2} + \frac{1}{2} \max \{ F_n(x) - F_n(-x^*), C_{nm}(x) - C_{nm}(-x^*) \} \}$$
  
=  $F_{nm}^*(x).$ 

Therefore, by symmetry,  $\hat{F}_{n,m}^2 = F_{nm}^*$ .

### 2.2. Bias Functions

The operator  $\Phi_1$  does not introduce any bias in the "symmetrization" procedure. In fact, it is well known that  $\hat{F}_n$  and  $\hat{G}_m$  are unbiased estimators for F and G, and have smaller variance than  $F_n$  and  $G_m$  respectively. However, the operators  $\Phi_2$  and  $\Phi_2^*$  do introduce bias when estimating F and G. The bias function of the estimators are discussed next and compared to the estimator provided by El Barmi and Mukerjee (2008).

For  $x \ge 0$  define  $F_n^+(x) = \frac{1}{n} \sum_{i=1}^n I_{[-x \le X_i \le x]}$ ,  $F_{nm}^{+*} = \max\{F_n^+, \frac{nF_n^+ + mG_m^+}{n+m}\}$  and finally, let  $F_{nm}^* = \frac{1}{2}(1 + F_{nm}^{+*})$ ;  $G_m^+$ ,  $G_{n,m}^{+*}$  and  $G_{n,m}^*$  are defined similarly. The estimator  $F_{nm}^*$  is the estimator for F studied by El Barmi and Mukerjee (2008) following ideas of Rojo (2004). Note that for  $x \ge 0$ ,

$$E(F_{nm}^{*}(x)) = \frac{1}{2} + \frac{1}{2}E(F_{nm}^{+*}(x))$$

$$= \frac{1}{2} + \frac{1}{2}E\{F_n^+(x) + \max\{0, \frac{m}{m+n}(G_m^+(x) - F_n^+(x))\}\}$$

$$= \frac{1}{2} + \frac{1}{2}E(F_n^+(x)) + \frac{m}{2(m+n)}E\{\max(0, G_m^+(x) - F_n^+(x))\}$$

and since  $\frac{1}{2} + \frac{1}{2}E(F_n^+(x)) = F(x)$ ,

(2.8) 
$$Bias(F_{nm}^{*}(x)) = \frac{m}{2(m+n)} E\{\max(0, G_{m}^{+}(x) - F_{n}^{+}(x))\}.$$

Note that  $Bias(F_{nm}^*(x)) \to 0$  as  $\frac{n}{m} \to \infty$ . Since our estimator  $\widehat{F}_{nm}^2$  defined by (2.6) turns out to be equal  $F_{nm}^*$ , then its bias function is also given by (2.8).

Now consider the estimator  $\widehat{F}_{nm}^1$  given by (2.4). For  $x \ge 0$ ,

$$\widehat{F}_{n,m}^{1}(x) = \Phi_{1}(\max(F_{n}(x), C_{nm}(x)))$$

$$= \frac{1}{2} \left\{ 1 + \max(F_{n}(x), C_{nm}(x)) - \min(F_{n}(x, x^{-}), C_{nn}(x, x^{-})) \right\}$$

$$43$$

$$= \frac{1}{2} \{1 + \max(F_n(x), C_{nm}(x)) - \min(F_n(-x^-), C_{nm}(-x^-)))\}$$

$$= \frac{1}{2} \left\{ 1 + F_n(x) - F_n(-x^-) + \max(0, C_{nm}(x) - F_n(x)) \right\}$$

$$+ \max(0, F_n(-x^-) - C_{nm}(-x^-)) \}.$$

Thus, 
$$E(\widehat{F}_{n,m}^1(x)) = F(x) + \frac{1}{2}E(\max(0, C_{nm}(x) - F_n(x))) + \frac{1}{2}E(\max(0, F_n(-x^-) - 51)) + \frac{1}{2}E(\max(0, F_n(-x^-) -$$

imsart-coll ver. 2008/08/29 file: Rojo\_Batun.tex date: March 25, 2009

 $= \frac{m}{2(m+n)} \{ E(\max(0, G_m(x) - F_n(x))) \}$ 

This result will also follow from the fact that  $\hat{F}_{nm}^1 >_p \hat{F}_{nm}^2 = F_{nm}^*$ . Next consider the estimator  $F_{nm}^2$  defined in equation (2.2) and in Rojo and

 $+E(\max(0, F_n(-x^-) - G_m(-x^-))))$ 

 $Bias(\widehat{F}_{n,m}^{1}(x)) = \frac{1}{2}E(\max(0, \frac{m}{n+m}(G_{m}(x) - F_{n}(x))))$ 

Then, for  $x \ge 0$ ,

$$Bias(F_{nm}^{1}(x)) = \frac{1}{4}E\left(\max\{0, G_{m}^{+}(x) - F_{n}^{+}(x) - F_{n}((-x)^{-}) - F_{n}(x) + 1\}\right)$$

$$+\frac{1}{4}E\left(\max\{0, G_m^+(x) - F_n^+(x) - 1 + F_n(x) + F_n((-x)^-)\}\right).$$

The last expression is then seen to be equal to

$$\frac{1}{4}E(\max\{\max(0, G_m^+(x) - F_n^+(x) - F_n((-x)^-) - F_n(x) + 1),$$

$$\max(G_m^+(x) - F_n^+(x) + F_n((-x)^-) + F_n(x) - 1, 2(G_m^+ - F_n^+)))$$

$$\geq \frac{1}{4}E(\max(0, 2(G_m^+(x) - F_n^+(x)))) = Bias(F_{n,m}^*).$$

imsart-coll ver. 2008/08/29 file: Rojo\_Batun.tex date: March 25, 2009

+  $\frac{1}{4}E(\max(0, (-1+2F_n((-x)^-)+G_m^+(x))))).$ 

$$= \frac{1}{2}(1 + F_n(x) - F_n((-x)^-)) + \frac{1}{2}\max(0, \frac{1}{2}(1 - 2F_n(x) + G_m^+(x))) + \frac{1}{2}\max(0, \frac{1}{2}(1 - 2F_n(x))) + \frac{1}{2}\max(0, \frac{1}{2}(1 - 2F_n(x))) + \frac{1}{2}(1 - 2F_n(x))) + \frac{1}{2}\max(0, \frac{1}{2}(1 - 2F_n(x))) + \frac{1}{2}(1 - 2F_n(x)) + \frac{1}{2}(1 - 2F_n(x))) + \frac{1}{2}(1 - 2F_n(x)) + \frac{1}{2}(1 - 2F_n(x))) + \frac{1}{2}(1 - 2F_n(x)) + \frac{1}{2}(1 - 2F_n(x$$

$$= \frac{1}{2}(1 + F_n(x) - F_n((-x)^-)) + \frac{1}{2}\max(0, \frac{1}{2}(1 - 2F_n(x) + G_m^+(x))) + \frac{1}{2}\max(0, \frac{1}{2}(-1 + 2F_n(x) - F_n(x))) + \frac{1}{2}\exp(0, \frac{1}{2}(-1 + 2F_n(x))) + \frac{1}{2}\exp(0, \frac{1}{2}(-1 + 2F_n(x)$$

Batun-Cutz (2007):

Ther

refore,  

$$E(F_{nm}^{1}(x)) = F(x) + \frac{1}{4}E(\max(0, (1 - 2F_{n}(x) + G_{m}^{+}(x))))$$

$$+\frac{1}{2}\max(0,\frac{1}{2}(-1+2F_n((-x)^-)+G_m^+(x))).$$

follows easily that 
$$E(F_{nm}^2(x)) = F(x) + \frac{1}{2}E(\max(0, G_m^+(x) - F_n^+(x)))$$
 and

It follows easily that 
$$E(F_{nm}^2(x)) = F(x) + \frac{1}{2}E(\max(0, G_m^+(x) - F_n^+(x)))$$
 and hence  
 $Bias(F_{nm}^2(x)) = \frac{1}{2}E(\max(0, G_m^+(x) - F_n^+(x))) > Bias(F_{nm}^*)$ , for  $x \ge 0$ .  
Finally, consider the estimator  $F_{nm}^1$  given in Rojo and Batun-Cutz (2007). For  
 $x \ge 0$ 

Finally, consider the estimator 
$$F_{nm}$$
 given in Rojo and Batun-Cutz (2007).  $x \ge 0$ 

$$F_{nm}^{1}(x) = \frac{1}{2} \left\{ 1 + \max(F_{n}(x), \frac{1}{2}(1 + G_{m}(x) - G_{m}((-x)^{-}))) \right\}$$

$$(x) = \frac{1}{2} \left\{ 1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-))) \right\}$$

lows easily that 
$$E(F_{nm}^2(x)) = F(x) + \frac{1}{2}E(\max(0, G_m^+(x) - F_n^+(x))) = \frac{1}{2}E(\max(0, G_m^+(x) - F_n^+(x))) > Bias(F_{nm}^*)$$
, for  $x \ge 1$ 

$$Bias(F_{nm}(x)) = \frac{1}{2}E(\max(0, G'_m(x) - F'_n(x))) > Bias(F'_nm), \text{ for } x \ge 0.$$
  
Finally, consider the estimator  $F^1_{nm}$  given in Rojo and Batun-Cutz (2007). F  $x > 0$ 

$$\geq 0$$

$$F_{nm}^{1}(x) = \frac{1}{2} \left\{ 1 + \max(F_{n}(x), \frac{1}{2}(1 + G_{m}(x) - G_{m}((-x)^{-}))) \right\}$$

$$= \frac{1}{2} \left\{ 1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^{-1}))) - \min(F_n(x), \frac{1}{2}(1 + G_n((-x)^{-1}) - G_n((-x)^{-1}))) \right\}$$

$$\frac{1}{2} \left\{ 1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-))) \right\}$$

$$\left\{1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-)))\right\}$$

$$\left\{1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-)))\right\}$$

$$\max(0, G'_m(x) - F'_n(x))) > Bias(F'_{nm})$$
, for  $x \ge 0$ .  
ne estimator  $F^1_{nm}$  given in Rojo and Batun-Cutz (2007). For

$$G_m((-x)^-)))$$

+ max
$$(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-)))$$

$$= \min(E_{n}(x), \frac{1}{2}(1 + G_{m}(x) - G_{m}((-x))))$$

$$-\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \Big\}$$

$$-\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \bigg\}$$

$$\frac{1}{2} \left\{ 1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-))) \right\}$$

$$\frac{1}{2} \left\{ 1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^{-}))) \right\}$$

$$\frac{1}{2} \left\{ 1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-))) \right\}$$

$$2\left( -\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \right)$$

$$-\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \bigg\}$$

$$= \frac{1}{2}(1 + F_n(x) - F_n((-x)^-)) + \frac{1}{2}\max(0, \frac{1}{2}(1 - 2F_n(x) + G_m^+(x)))$$

$$\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x))))$$

$$-\min(F_n(-x), \frac{1}{2}(1+G_m((-x)^-)-G_m((x))))\bigg\}$$

$$2\left(\frac{1}{2}\left(1+G_{m}(x)\right)^{2}-G_{m}(x)\right)\right)$$
  
-min(F<sub>n</sub>(-x),  $\frac{1}{2}(1+G_{m}((-x)^{-})-G_{m}((x))))$ 

$$= \frac{1}{2} \left\{ 1 + \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^{-1}))) - \min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^{-1}) - G_m((-x)))) \right\}$$

$$-\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \bigg\}$$

$$-\min(F_n(-x), \frac{1}{2}(1+G_m((-x)^-)-G_m((x)))) \bigg\}$$

$$\frac{1}{2}(1+F_n(x)-F_n((-x)^-)) + \frac{1}{2}\max(0,\frac{1}{2}(1-2F_n(x)+G_m^+(x)$$

$$\lim_{x \to \infty} (F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x))))$$
  
$$1 + F_n(x) - F_n((-x)^-)) + \frac{1}{2}\max(0, \frac{1}{2}(1 - 2F_n(x) + C_n(x)))$$

$$\frac{1}{2}(1+F_n(x)-F_n((-x)^-)) + \frac{1}{2}\max(0,\frac{1}{2}(1-2F_n(x)+G_m^+(x))) + \frac{1}{2}\max(0,\frac{1}{2}(1-2F_m^+(x)+G_m^+(x))) + \frac{1}{2}\max(0,\frac{1}{2}$$

$$-\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \bigg\}$$

$$-\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \bigg\}$$

$$-\min(F_n(-x), \frac{1}{2}(1 + G_m((-x)^-) - G_m((x)))) \bigg\}$$

$$\left[F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-)))\right]$$

+ max
$$(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-)))$$

$$\max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-))))$$

$$+ \max(F_n(x), \frac{1}{2}(1 + G_m(x) - G_m((-x)^-)))$$

$$F_{nm}^{2}(x) = \max\left\{\frac{1}{2}(1+F_{n}(x)-F_{n}((-x)^{-})), \frac{1}{2}(1+G_{m}(x)-G_{m}((-x)^{-}))\right\}.$$
  
It follows easily that  $E(F_{nm}^{2}(x)) = F(x) + \frac{1}{2}E(\max(0, G_{m}^{+}(x) - F_{n}^{+}(x)))$  and hence

 $+\frac{1}{2}E(\frac{m}{n+m}\max(0,F_n(-x^-)-G_m(-x^-)))$ 

 $\geq \frac{m}{2(m+n)}E(\max(0, G_m^+(x) - F_n^+(x))) = Bias(F_{nm}^*).$ 

d hence

 $C_{nm}(-x^{-}))$  and then, for  $x \ge 0$ 

З

The corresponding inequalities for the case of x < 0 follow by symmetry. Similar results may be obtained for the estimators  $G_{n,m}^1 = \Phi_1(\Phi_2^*(\Phi_1(F_n), G_m)), G_{n,m}^2 = \Phi_2^*(\Phi_1(F_n), \Phi_1(G_m))$ , and  $\widehat{G}_{n,m}^1$  and  $\widehat{G}_{n,m}^2$ . It is easy to see that all the estimators for *F* have positive (negative) bias for x > 0 (x < 0), while the estimators for *G* have negative (positive) bias for x > 0 (x < 0). The following proposition summarizes the results about the bias functions.

**Proposition 1.** Let  $F >_p G$  be symmetric distribution functions, and let  $X_1, \ldots, X_n$ and  $Y_1, \ldots, Y_m$  be independent random samples from F and G respectively. The bias functions of the estimators for F and G given by (2.1), (2.2), (2.4), (2.5), (2.6), and (2.7), satisfy the following properties. For all x,

$$\begin{aligned} (i)|Bias(\widehat{F}_{n,m}^{1}(x))| &\geq |Bias(\widehat{F}_{n,m}^{2}(x))| \\ &= \frac{m}{2(m+n)} E\{\max(0,G_{m}^{+}(|x|) - F_{n}^{+}(|x|)\} = |Bias(F_{n,m}^{*}(x))| \end{aligned}$$

$$(ii) |Bias(F_{n,m}^1(x))| \ge |Bias(F_{n,m}^2(x))| \ge |Bias(\widehat{F}_{n,m}^2(x))|$$

$$(iii) |Bias(\widehat{C}_{n,m}^1(x))| \ge |Bias(\widehat{C}_{n,m}^2(x))| = \frac{m}{2} |F(min(0, E^+(|m|))|)$$

$$(iii) |Bias(\widehat{G}_{n,m}^{1}(x))| \ge |Bias(\widehat{G}_{n,m}^{2}(x))| = -\frac{m}{2(m+n)} E\{\min(0, F_{n}^{+}(|x|) - G_{m}^{+}(|x|))\}$$

$$(iv) |Bias(G_{n,m}^1(x))| \ge |Bias(G_{n,m}^2(x))| \ge |Bias(G_{n,m}^2(x))|.$$

an

# 2.3. Estimators as Projections onto Appropriate Convex Spaces

Recall the definitions of the new estimators given by (2.4) - (2.7). Schuster (1975) showed that the operator  $\Phi_1$  projects its argument to its closest symmetric distribution. That is, letting  $\mathcal{F}$  be the convex set of symmetric distributions about zero, then for an arbitrary distribution H,  $\| \Phi_1(H) - H \|_{\infty} = \inf_{F \in \mathcal{F}} \| H - F \|_{\infty}$ . Rojo and Ma (1996), and Rojo and Batun-Cutz (2007) have shown that the operator  $\Phi_2$  has the property that for arbitrary distributions H and G,  $|\Phi_2(H(x), G(x)) - H(x)| =$  $\inf_{F \in \mathcal{F}^*} |F(x) - G(x)|$ , where  $\mathcal{F}^*$  is the convex set of distributions F satisfying (1.3). Thus, for F and G distribution functions let

 $\mathcal{F}_1 = \{ \text{distribution functions } F \text{ satisfying } (1.3) \text{ with } G \text{ replaced by } C_{n,m} \}$ 

 $\mathcal{F}_1^* = \{ \text{symmetric distributions } F \text{ satisfying } (1.3) \text{ with } G \text{ replaced by } \Phi_1(C_{n,m}) \}$ 

d 
$$\mathcal{F}_2^* = \{ \text{all symmetric at 0 distribution functions} \}.$$

Thus the estimator  $\widehat{F}_{n,m}^2$  first projects  $F_n$  onto  $\mathcal{F}_2^*$  and then projects  $\Phi_1(F_n)$ onto  $\mathcal{F}_1^*$ . By contrast, the estimator  $\widehat{F}_{n,m}^1$  first projects  $F_n$  onto  $\mathcal{F}_1$  to obtain  $\Phi_2(F_n, C_{n,m})$  and then projects the latter onto  $\mathcal{F}_1^*$ . With appropriate changes in the above notation, similar comments hold for the estimators  $\widehat{G}_{n,m}^i$  for i = 1, 2.

# 2.4. Peakedness Order of New and Previous Estimators

<sup>46</sup> By construction, the estimators  $F_{n,m}^i$  and  $\hat{F}_{n,m}^i$ , for i = 1, 2 are more peaked than the estimators  $G_{n,m}^i$  and  $\hat{G}_{n,m}^i$ , respectively. Rojo and Batun-Cutz (2007) showed that  $F_{n,m}^1 >_p F_{n,m}^2$ . The next theorem provides comparisons in terms of peakedness for several of the estimators and provides a simple relationship between  $F_{n,m}^2$  and  $\hat{F}_{n,m}^2$ . **Lemma 1.** Let  $F >_p G$  be symmetric distribution functions, and let  $X_1, \ldots, X_n$ and  $Y_1, \ldots, Y_m$  be independent random samples from F and G respectively. Consider the estimators for F and G given by (2.1), (2.2), (2.4), (2.5), (2.6), (2.7). Then (i)  $\widehat{F}_{n,m}^2 = \frac{n}{n+m}\widehat{F}_n + \frac{m}{n+m}F_{n-m}^2$ (ii)  $\widehat{F}_{n\,m}^1 >_p \widehat{F}_{n\,m}^2 >_p \widehat{G}_{n\,m}^2 >_p \widehat{G}_{n\,m}^1$ (iii)  $F_{n,m}^1 >_p F_{n,m}^2 >_p \widehat{F}_{n,m}^2$ , and  $G_{n,m}^1 <_p G_{n,m}^2 <_p \widehat{G}_{n,m}^2$ . *Proof.* (i) For x > 0,  $\widehat{F}_{n,m}^{2}(x) = \max\{\widehat{F}_{n}(x), \widehat{C}_{n,m}(x)\} = \frac{n}{n+m}\widehat{F}_{n}(x) + \frac{m}{n+m}\max\{\widehat{F}_{n}(x), \widehat{G}_{m}(x)\}$  $= \frac{n}{n+m}\widehat{F}_n(x) + \frac{m}{n+m}F_{n,m}^2(x).$ The result then follows by symmetry. (ii) First we prove that  $\widehat{F}_{n,m}^1 >_p \widehat{F}_{n,m}^2$ . Let  $x \ge 0$ , then  $\widehat{F}_{n,m}^{1}(x) = \frac{1}{2} \left[ \max\{F_{n}(x), C_{n,m}(x)\} + 1 - \min\{F_{n}((-x)^{-}), C_{n,m}((-x)^{-})\} \right]$  $\geq \frac{1}{2} \left[ C_{n,m}(x) + 1 - C_{n,m}((-x)^{-}) \right] = \widehat{C}_{n,m}(x).$ (2.9)Using similar arguments it can be shown that  $\widehat{F}_{n,m}^1(x) \geq \widehat{F}_n(x)$ . Therefore, combining the last inequality and (2.9) we obtain  $\widehat{F}_{n,m}^1(x) \geq \widehat{F}_{n,m}^2(x)$ . The result follows from symmetry. We now prove that  $\widehat{F}_{n,m}^2 >_p \widehat{G}_{n,m}^2$ . For  $x \ge 0$ ,  $\widehat{F}_{n,m}^2(x) = \max\{\widehat{F}_n(x), \widehat{C}_{n,m}(x)\} \ge 0$  $\widehat{C}_{n,m}(x) \geq \widehat{G}_{n,m}^2(x)$ . The result follows by symmetry. Since for  $x \ge 0$ ,  $\widehat{G}_{n,m}^1(x) \le \widehat{C}_{n,m}(x)$  and  $\widehat{G}_{n,m}^1(x) \le \widehat{G}_m(x)$ . Then  $\widehat{G}_{n,m}^2 >_p \widehat{G}_{n,m}^1$ 

Finally consider (iii). The result that  $F_{n,m}^1 >_p F_{n,m}^2$  follows from Rojo and Batun-Cutz (2007). The result that  $F_{n,m}^2 >_p \widehat{F}_{n,m}^2$  follows from the arguments used to prove that  $\operatorname{Bias}(F_{n,m}^2) \geq \operatorname{Bias}(\widehat{F}_{n,m}^2)$ .

Note that (i) implies that for  $x \ge 0$ ,  $Bias(F_{n,m}^2(x)) = \frac{m+n}{m}Bias(\widehat{F}_{n,m}^2(x))$ , so that  $|Bias(F_{n,m}^2(x))| = \frac{m+n}{m} |Bias(\widehat{F}_{n,m}^2(x))|$  for all x, thus providing a more accurate description of the result about bias given in proposition 1.

#### 3. Asymptotics

by symmetry.

This section discusses the strong uniform convergence of the estimators and their asymptotic distribution theory. One important aspect of the asymptotic results for the estimators  $\hat{F}_{n,m}^i$   $(\hat{G}_{n,m}^i)$ , i = 1, 2 discussed here is that they hold even when only n(m) tends to infinity. This is in sharp contrast with the results of Rojo and Batun-Cutz (2007) and those of El Barmi and Mukerjee (2008). We discuss the strong uniform convergence first. 

З

З

(3.2)

(3.3)

3.1. Strong Uniform Convergence The following theorem provides the strong uniform convergence of the estimators  $F_{n,m}^i$   $(G_{n,m}^i)$ , i = 1, 2. The results use the strong uniform convergence of the symmetrized  $\hat{F}_n$  ( $\hat{G}_m$ ) to F (G) as  $n \to \infty$  ( $m \to \infty$ ), Schuster (1975). **Theorem 3.1.** Let F and G be symmetric distribution functions with  $F >_p G$ , and let  $X_1, \ldots, X_n$  and  $Y_1, \ldots, Y_m$  be independent random samples from F and G respectively. Then, (i)  $\widehat{F}_{n,m}^i$ , for i = 1, 2, converge uniformly with probability one to F as  $n \to \infty$ . (ii)  $\widehat{G}_{n,m}^i$  for i = 1, 2 converge uniformly with probability one to G as  $m \to \infty$ . *Proof.* (i) Consider  $\widehat{F}_{n,m}^2$  first. Then, for  $x \ge 0$ , (3.1)  $\widehat{F}_{n,m}^2(x) - F(x) = \widehat{F}_n(x) - F(x) + \frac{m}{n+m} \max\{0, \widehat{G}_m(x) - \widehat{F}_n(x)\}.$ But, since  $F(x) \ge G(x)$ ,  $\widehat{G}_m(x) - \widehat{F}_n(x) < \widehat{G}_m(x) - G(x) + F(x) - \widehat{F}_n(x).$ Hence  $\max\{0, \widehat{G}_m(x) - \widehat{F}_n(x)\} \leq \max\{0, \widehat{G}_m(x) - G(x) + F(x) - \widehat{F}_n(x)\}$  $< |\widehat{G}_{m}(x) - G(x)| + |\widehat{F}_{n}(x) - F(x)|$ and therefore, the left side of (3.1) is bounded above by  $|\widehat{F}_n(x) - F(x)| + (\frac{m}{m+n}) \{ |\widehat{G}_m(x) - G(x)| + |\widehat{F}_n(x) - F(x)| \}$ Since  $\widehat{F}_n$ , and  $\widehat{G}_m$  are strongly and uniformly consistent for F and G, then as  $n \to \infty$ , with probability one,  $\sup_{x \ge 0} |\widehat{F}_{n,m}^2(x) - F(x)| \to 0,$ regardless of whether  $m \to \infty$  or not. When x < 0 the result follows by symmetry. Let us now consider the case of  $\widehat{F}_{n,m}^1$ . For  $x \ge 0$  $\widehat{F}_{n,m}^{1}(x) - F(x) = \widehat{F}_{n}(x) - F(x) + \frac{1}{2} \frac{m}{n+m} \left[ \max\{0, G_{m}(x) - F_{n}(x) \} \right]$  $-\min\{0, G_m(-x^-) - F_n(-x^-)\}].$ Since  $F(x) \ge G(x)$  and  $F(-x) \le G(-x)$ , then it follows that  $\max\{0, G_m(x) - F_n(x)\} - \min\{0, G_m(-x^-) - F_n(-x^-)\}$ is bounded above by

 $\max\{0, G_m(x) - G(x) + F(x) - F_n(x)\} - \min\{0, G_m(-x^-) - G(-x) + F(-x) - F_n(-x^-)\}$ 

and, therefore, the left side of (3.2) is bounded above by
### J. Rojo and J. Batún-Cutz

(3.4)

$$|\widehat{F}_n(x) - F(x)| + \frac{1}{2} \frac{m}{m+m} (|G_m(x) - G(x)| + |F(x) - F_n(x)|$$

+  $|G_m(-x^-) - G(-x)| + |F(-x) - F_n(-x^-)|).$ 

Taking the supremum over x in (3.4), and then letting  $n \to \infty$ , the result follows, whether  $m \to \infty$  or not, from the strong uniform convergence of  $F_n$ ,  $G_m$ , and  $F_n$ to F, G, and F respectively. The result for x < 0 follows by symmetry.

(iii) The proof for the strong uniform convergence of  $\widehat{G}_{n,m}^2$  to G, when only  $m \to \infty$ is similar. We sketch the proof. For x < 0

$$\widehat{G}_{n,m}^2(x) - G(x) = \widehat{G}_m(x) - G(x) + \frac{n}{n+m} \max\{0, \widehat{F}_n(x) - \widehat{G}_m(x)\}.$$

Therefore, since F(x) < G(x) for x < 0,  $\max\{0, \widehat{F}_n(x) - \widehat{G}_m(x)\}$  is bounded above by

$$\max\{0, \hat{F}_n(x) - F(x) + G(x) - \hat{G}_m(x)\} \le |\hat{F}_n(x) - F(x)| + |G(x) - \hat{G}_m(x)|.$$

When  $m \to \infty$ , the result follows, regardless of whether  $n \to \infty$  or not, from the strong uniform convergence of  $\widehat{F}_n$  and  $\widehat{G}_m$  and using a symmetry argument to handle the case of x > 0.

(iv) This case is omitted as it follows from similar arguments.

# 3.2. Weak Convergence

Consider first the point-wise asymptotic distribution for  $\widehat{F}_{n,m}^i$ , i = 1, 2. Recall that

$$\sqrt{n}(\widehat{F}_n(x) - F(x)) \xrightarrow{W} N\left(0, \frac{F(-|x|)(2F(|x|) - 1)}{2}\right).$$

Therefore, when  $n/m \to \infty$ , and using (3.1)-(3.4), Slutsky's theorem and the central limit theorem for  $\widehat{F}_n$ , we get the following result:

(3.5) 
$$\sqrt{n}(\widehat{F}_{nm}^{i}(x) - F(x)) \xrightarrow{W} N\left(0, \frac{F(-|x|)(2F(|x|) - 1)}{2}\right).$$

Thus, under these conditions,  $\hat{F}^i_{n,m}$ , i = 1, 2, are  $\sqrt{n}$ -equivalent and have the same asymptotic distribution as the symmetrized  $\widehat{F}_n$  which happens to have the same asymptotic limit as in the case when G is completely known. Note that this result assumes only that  $n/m \to \infty$  and hence the result holds if m is fixed and  $n \to \infty$ . This is in sharp contrast with the results of El Barmi and Mukerjee (2008) that require that both n and m tend to infinity. Similar results hold for the estimators  $\widehat{G}_{n,m}^{i}$ , i = 1, 2. These are summarized in the following theorem.

**Theorem 3.2.** Suppose that  $F >_p G$  and let  $X_1, \ldots, X_n$  and  $Y_1, \ldots, Y_m$  be random samples from F and G respectively. Then for i = 1, 2,

(i) If 
$$n/m \to \infty$$
 then

$$\sqrt{n}(\widehat{F}^i_{nm}(x) - F(x)) \quad \stackrel{\mathcal{D}}{\to} \quad N\left(0, \frac{F(-|x|)(2F(|x|) - 1)}{2}\right).$$

imsart-coll ver. 2008/08/29 file: Rojo\_Batun.tex date: March 25, 2009

З

(ii) If  $m/n \to \infty$  then

$$\sqrt{n}(\widehat{G}^i_{nm}(x) - G(x)) \quad \overset{\mathcal{D}}{\to} \quad N\left(0, \frac{G(-|x|)(2G(|x|) - 1)}{2}\right).$$

We now turn our attention to the weak convergence of the processes

$$\left\{\sqrt{n}\left(\widehat{F}_{nm}^{i}(x) - F(x)\right) : -\infty < x < \infty\right\},\,$$

and

$$\left\{\sqrt{n}\left(\widehat{G}_{nm}^{i}(x) - G(x)\right) : -\infty < x < \infty\right\},$$
<sup>10</sup>
<sup>11</sup>
<sup>12</sup>
<sup>12</sup>

for i = 1, 2. Only the results for  $\widehat{F}_{n,m}^i$ , i = 1, 2 will be discussed in detail as the results for  $\widehat{G}_{n,m}^i$ , i = 1, 2 can be obtained by similar arguments. Although the processes  $\left\{\sqrt{n}\left(\widehat{F}_{nm}^i(x) - F(x)\right) : -\infty < x < \infty\right\}$  for i = 1, 2 are correlated, we are only interested in their marginal behavior. For that purpose let  $\{W_1(x) : -\infty < x < \infty\}$  denote a mean zero Gaussian process with covariance function

(3.6) 
$$E(W_1(x)W_1(y)) = \begin{cases} \frac{1}{2}(1-F(y))(F(x)-F(-x)) & \text{if } |y| > |x| \\ \frac{1}{2}F(x)(F(-y)-F(y)) & \text{if } |y| < |x|, \end{cases}$$

and let  $\{W_2(x) : -\infty < x < \infty\}$  denote a mean zero Gaussian process with covariance function

$$(3.7) \quad E(W_2(x)W_2(y)) = \begin{cases} \frac{1}{2}(1-G(y))(G(x)-G(-x)) & \text{if } |y| > |x| \\ \frac{1}{2}G(x)(G(-y)-G(y)) & \text{if } |y| < |x|. \end{cases}$$

We have the following result:

**Theorem 3.3.** Under the conditions of the previous Theorem,

(i) If 
$$n/m \to \infty$$
, then

$$\{\sqrt{n}(\widehat{F}_{nm}^i(x) - F(x)), -\infty < x < \infty\} \xrightarrow{W} \{W_1(x) : -\infty < x < \infty\}, and$$

(ii) If  $m/n \to \infty$ , then

$$\{\sqrt{n}(\widehat{G}_{nm}^{i}(x) - G(x)), -\infty < x < \infty\} \xrightarrow{W} \{W_{2}(x) : -\infty < x < \infty\}.$$

*Proof.* The proof follows easily by the continuous mapping Theorem after observing that the weak limit of  $\{\sqrt{n}(\hat{F}_n(x) - F(x))\}, -\infty < x < \infty\}$  is the process  $\{W_1(x) : -\infty < x < \infty\}$ , together with the fact that, using (3.1),

(3.8) 
$$\widehat{F}_{n,m}^2(x) - F(x) = \widehat{F}_n(x) - F(x) + \frac{m}{n+m} \max\{0, \widehat{G}_m(x) - \widehat{F}_n(x)\},\$$

with  $\|\sqrt{n}\frac{m}{n+m}\{\max 0, \widehat{G}_m - \widehat{F}_n\}\|_{\infty} \to 0$  with probability one, where  $\|\cdot\|_{\infty}$  denotes the sup norm. Similar arguments yield the results for the other processes.

The asymptotic theory for  $\widehat{F}_{n,m}^2$  was discussed by El Barmi *et al* (2008) for the case that both *n* and *m* go to infinity and hence their result does not include our result here when *m* is bounded and  $n \to \infty$ . When  $n/m \to c$  with  $0 \le c < \infty$ , and F(x) > G(x) for all x > 0 the weak limit of  $\{\sqrt{n}(\widehat{F}_{nm}^i(x) - F(x)), -\infty < x < \infty\}$  is  $\{W_1(x) : -\infty < x < \infty\}$ , for i = 1, 2, which is the weak limit of the

J. Rojo and J. Batún-Cutz

process  $\{\sqrt{n}(F_{n,2}(x) - F(x)), -\infty < x < \infty\}$  discussed in Rojo and Batun-Cutz (2007). Let  $\{Z(x), -\infty < x < \infty\}$  represent the weak limit of the empirical process  $\{\sqrt{n}(F_n(x) - F(x)), -\infty < x < \infty\}$ . That is  $\{Z(x), -\infty < x < \infty\}$  is a mean zero Gaussian process with covariance function E(Z(t)Z(s)) = F(s)(1-F(t)) for  $s \leq t$ . When  $n/m \to c$  with  $0 \le c < \infty$ , and F(x) = G(x) for all x the weak limits of  $\{\sqrt{n}(\widehat{F}_{nm}^{i}(x) - F(x)), -\infty < x < \infty\}$  for i = 1, 2 follow from the results in Rojo (2004) as follows: **Theorem 3.4.** Let F(x) = G(x) for all x and let  $n/m \to c$  for  $0 \le c < \infty$ . Let  $\{W_i(x), -\infty < x < \infty\}$ , for i = 1, 2 be the mean zero Gaussian processes with covariance functions given by (3.6) and (3.7), respectively. Let  $W_i^*(x) =$  $W_i(|x|)sgn(x)$ , for i = 1, 2. Then (i) The process  $\sqrt{n}(\widehat{F}_{n,m}^2 - F(x)), -\infty < x < \infty$ } converges weakly to the process  $\{\max(W_1^*(x), \frac{\sqrt{c}W_2^*(x) + cW_1^*(x)}{1+c}), -\infty < x < \infty\}$  with  $W_1^* \stackrel{\mathcal{D}}{=} W_2^*$  and independent of the process  $\{\max(W_1^*(x), \frac{\sqrt{c}W_2^*(x) + cW_1^*(x)}{1+c}), -\infty < x < \infty\}$ (ii) The process  $\sqrt{n}(\hat{F}_{n,m}^1 - F(x)), -\infty < x < \infty$ } converges weakly to the process  $\{H(|x|)sgn(x), -\infty < x < \infty\}$  where  $H(x) = \frac{1}{2}\{\max\{Z_1(x), \frac{c}{1+c}Z_1(x) + C_1(x)\}\}$  $\frac{\sqrt{c}}{1+c}Z_{2}(x)\} - \min\{Z_{1}(-x), \frac{c}{1+c}Z_{1}(-x) + \frac{\sqrt{c}}{1+c}Z_{2}(-x)\}, \text{ and } \{Z_{i}(x), -\infty < x < \infty\}, i = 1, 2 \text{ are independent copies of the process } \{Z(x), -\infty < x < \infty\}.$ *Proof.* (i) Consider  $\hat{F}_{n,m}^2$  first. When F(x) = G(x) for all x, it follows from (3.8) that, for  $x \ge 0$ ,  $\sqrt{n}(\widehat{F}_{n,m}^{2}(x) - F(x)) = \max\{\sqrt{n}(\widehat{F}_{n}(x) - F(x)), \sqrt{n/m}\frac{m}{n+m}\sqrt{m}(\widehat{G}_{m}(x) - G(x))\}$ +  $\frac{n}{n+m}\sqrt{n}(F(x)-\widehat{F}_n(x))\}.$ (3.9)By the independence of  $\widehat{F}_n$  and  $\widehat{G}_m$  and their weak convergence to  $W_1$  and  $W_2$ , it follows that the bivariate process  $\left\{\sqrt{n/m}\frac{m}{n+m}\sqrt{m}(\widehat{G}_m(x) - G(x)), \frac{n}{n+m}\sqrt{n}(F(x) - \widehat{F}_n(x)), -\infty < x < \infty\right\}$ converges weakly to the process  $\{\frac{\sqrt{c}}{1+c}W_2(x), \frac{c}{1+c}W_1(x), -\infty < x < \infty\}$ . Since for x < 0,  $\widehat{F}_{n,m}^2(x) - F(x) \stackrel{\mathcal{D}}{=} F(-x) - \widehat{F}_{n,m}^2(-x)$ , the result then follows for  $0 < c < \infty$  from the continuous mapping theorem after observing that the mapping  $h(x,y) = (\frac{1+c}{c}y, x+y)$  is continuous, and then applying it to (3.9) to get the result. The case of c = 0 follows immediately since it then follows that the second term on the right side of (3.9) converges to zero in probability. 

(*ii*) Note that for x > 0

$$\frac{n}{m+n}(F_n(x) - F(x)) + \frac{m}{n+m}(G_m(x) - F(x))\}$$

$$+ \frac{1}{2}\min\{F_n(-x) - F(-x), 48$$

$$\frac{2}{m+n}(F_n(-x) - F(-x)) + \frac{m}{n+m}(G_m(-x) - F(-x))\}.$$
<sup>50</sup>
<sub>51</sub>

imsart-coll ver. 2008/08/29 file: Rojo\_Batun.tex date: March 25, 2009

Since the function  $h(x, y, z, w) = \frac{1}{2} [\max\{x, x+z\} - \min\{y, y+w\}]$  is continuous,

by the continuous mapping theorem we obtain  

$$\sqrt{n}(\widehat{F}_{nm}^1(x) - F(x)) \xrightarrow{W} \frac{1}{2} \left[ \max\{Z_1(x), \frac{c}{1+c}Z_1(x) + \frac{\sqrt{c}}{1+c}Z_2(x) \} \right]$$

$$\sqrt{n(F_{nm}^{1}(x) - F(x))} \rightarrow \frac{1}{2} \left[ \max\{Z_{1}(x), \frac{1}{1+c}Z_{1}(x) + \frac{1}{1+c}Z_{2}(x) \} \right]$$

(3.10) 
$$-\min\{Z_1(-x), \frac{c}{1+c}Z_1(x) + \frac{\sqrt{c}}{1+c}Z_2(-x)\} = H(x).$$

The result then follows after considering the case x < 0 and following a similar argument.

It has been observed, e.g. Rojo (1995a), Rojo (2004), and Rojo and Batun-Cutz (2007), that weak convergence of the processes of interest fails to hold when the underlying distributions F and G coincide at some point  $x_0$  and are unequal in some neighborhood to the right of  $x_0$ . That is the case here as well. Suppose that  $F(x_0) = G(x_0)$  for  $x_0 > 0$  and F(x) > G(x) for  $x \in (x_0, x_0 + \nu), \nu > 0$ . If  $\frac{m}{n} \to c$ ,  $0 < c \leq \infty$ , as  $m, n \to \infty$ , it follows that

(3.11) 
$$\sqrt{n}(\widehat{F}^1_{nm}(x_0) - F(x_0)) \xrightarrow{\mathcal{D}} H(|x_0|) sgn(x_0),$$

with H(x) defined as in (*ii*) of theorem 3.4 with  $(Z_1(x_0), Z_2(x_0))$  a zero-mean bivariate normal distribution vector with covariance  $(1 - F(x_0))F(x_0)$ .

However, for  $x \in (x_0, x_0 + \nu)$  the sequence  $\sqrt{n}(\hat{F}_{nm}^1(x) - F(x))$  converges in distribution to the distribution given in (3.5). Then it can be seen, using arguments as in Rojo (1995a), that the process  $\{\sqrt{n}(\hat{F}_{nm}^1(x) - F(x)) : -\infty < x < \infty\}$  is not tight and hence cannot converge weakly.

We finish this section with results that provide the weak convergence of the processes  $\{\sqrt{n}(F_{n,m}^i(x) - F(x)), -\infty < x < \infty\}$  for i = 1, 2, in the case that F(x) = G(x) for all x. **Theorem 3.5.** Let  $n/m \to c$  with  $0 \le c < \infty$ , and F(x) = G(x) for all x. (i) The process  $\{\sqrt{n}(F_{n,m}^2(x) - F(x)), -\infty < x < \infty\}$  converges weakly to  $\{sgn(x)\max\{sgn(x)W_1(x), sgn(x)\sqrt{c}W_2(x), -\infty < x < \infty\},\$ where  $W_1$  and  $W_2$  are independent copies of the mean zero Gaussian process with covariance function defined by (3.6). (ii) The process  $\{\sqrt{n}(F_{n,m}^1(x) - F(x)), -\infty < x < \infty\}$  converges weakly to  $\frac{1}{2}\{\max\{Z(xsgn(x)),\sqrt{c}W(xsgn(x))\}$  $-sgn(x)\min\{Z(-xsgn(x)), \sqrt{c}W(-xsgn(x))\}; -\infty < x < \infty\},$ where Z and W are independent mean zero Gaussian process with covariance functions defined by E(Z(s)Z(t)) = F(s)(1 - F(t)) for s < t, and (3.6) respectively. *Proof.* (i) The result follows from the independence of  $\{\sqrt{n}(\widehat{F}_n^*(x) - F(x)), -\infty < \infty\}$  $x < \infty$  and  $\{\sqrt{m}(G_m^*(x) - G(x)), -\infty < x < \infty\}$ , their weak convergence to  $W_1$ and  $W_2$ , and the continuous mapping theorem after observing that  $\sqrt{n}(F_{n}^2(x) - F(x)) = sgn(x)\max\{sgn(x)\sqrt{n}(\widehat{F}_n^*(x) - F(x)),$ 

$$sgn(x)(\sqrt{\frac{n}{m}}\sqrt{m}(\widehat{G}_m^*(x)-G(x))\}.$$

imsart-coll ver. 2008/08/29 file: Rojo\_Batun.tex date: March 25, 2009

(*ii*) Consider the case of x > 0 and write

$$\sqrt{n}(F_{n,m}^{1}(x) - F(x)) = \frac{\sqrt{n}}{2} \{1 - 2F(x) + \max(F_{n}(x), \widehat{G}_{m}(x)) - \min(F_{n}(-x), \widehat{G}_{m}(-x))\}$$

$$= \frac{1}{2} \{ \max\{\sqrt{n}(F_n(x) - F(x)), \sqrt{\frac{n}{m}}\sqrt{m}(\widehat{G}_m(x) - G(x)) \}$$

$$-\min\{\sqrt{n}(F_n(-x) - F(-x)), \sqrt{\frac{n}{m}}\sqrt{m}(\widehat{G}_m(-x) - G(-x))\}\}.$$

For x < 0, a similar argument leads to

$$\sqrt{n}(F_{n,m}^{1}(x) - F(x)) = \frac{1}{2} \{ \min\{\sqrt{n}(F_{n}(x) - F(x)), \sqrt{\frac{n}{m}}\sqrt{m}(\widehat{G}_{m}(x) - G(x)) \}$$

$$-\max\{\sqrt{n}(F_n(-x) - F(-x)), \sqrt{\frac{n}{m}}\sqrt{m}(\widehat{G}_m(-x) - G(-x))\}\},\$$

The result then follows by the continuous mapping theorem after letting  $n/m \to c$ with  $\sqrt{n}(F_n(x) - F(x))$  and  $\sqrt{m}(\widehat{G}_m(x) - G(x))$  independent and converging weakly to Z and W respectively.

# 4. Example with Sib-pair Data: An Illustration

In this section the estimator  $\widehat{F}_{n,m}^2$  is illustrated by using the sib-paired data for the Caucasian population in the Dallas metroplex area. As can be observed from Figure 2, the new estimated distribution functions now satisfy both the constraint of symmetry and the constraint of peakedness. Thus, since siblings with two alleles identical by descent are more similar than those siblings sharing only one allele identical by descent, the distribution function denoted by IBD2 is more peaked about zero than the other two distribution functions. Similar comments apply to the other comparisons.

### 5. Simulation Work

Monte Carlo simulations were performed to study the finite-sample properties of the estimators  $\hat{F}_{nm}^1$  and  $\hat{F}_{nm}^2$  defined by (2.4) and (2.6) respectively. We consider various examples of underlying distributions (Normal, Cauchy, Laplace, mixtures of normals, and T), and sample sizes (n = 10, 20, 30 for F and m = 10, 20, 30 for G). Each simulation consisted of 10,000 replications.

Figures 3 and 4 show the bias functions for the four estimators considered here. Figure 3 considers  $F \sim Cauchy(0, 1)$  and  $G \sim Cauchy(0, 2)$ , and Figure 4 considers the case with  $F \sim Laplace(0, 1)$  and  $G \sim Laplace(0, 1.5)$ . As shown in Proposition 1, the estimator  $\widehat{F}_{n,m}^2$  has uniformly the smallest absolute bias. These Figures are representative of the results that we obtained. One result that holds in all of our simulations is that  $|Bias(F_{n,m}^1(x))| \geq |Bias(\widehat{F}_{n,m}^1(x))|$  for all x. Unfortunately, we are unable to prove this conjecture.

Turning our attention to comparing the estimators in terms of the Mean Squared
Error (MSE) Figures 5 - 10 show the ratio of the MSE of the empirical distribution to the MSE of each of the four estimators considered here. These plots are
representative of all the examples considered. As it can be seen from the plots, the

З



FIG 2. Order restricted estimators for the sib-pair data incorporating peakedness.

empirical distribution function is dominated by the estimators in every case and for all x. Whereas the estimators  $\hat{F}_{n,m}^i$  behave better than the estimators  $F_{n,m}^i$ , i = 1, 2 in a neighborhood of zero, the roles are reversed on the tails of the underlying distribution. What is observed is that the region of the support of F where  $\hat{F}_{n,m}^i$  dominate  $F_{n,m}^i$ , i = 1, 2 shrinks as the tails of the distributions get heavier, and when the distribution G is far from F. Thus, there is no clear choice among the four estimators, unless the tail is of special interest, in which case the estimator  $F_{n,m}^2$  seems to be the clear choice.

#### 6. Conclusions

З

Estimators were proposed for the distribution functions F and G when it is known that  $F >_p G$ , and F and G are symmetric about zero. The estimator for F(G) was seen to be strongly uniformly consistent when only n(m) goes to infinity and the asymptotic theory of the estimators was delineated without requiring that both nand m go to infinity. Finite sample properties of the estimators were considered and it was shown that the estimator  $\hat{F}_{n,m}^2$  has the uniformly smaller absolute bias of the four estimators considered here. The choice of which estimator is best in terms of mean squared error (mse), however, is not clear. Although the estimators  $\hat{F}_{n,m}^i$ for i = 1, 2 have smaller mse than the estimators  $F_{n,m}^i$ , i = 1, 2 in a neighborhood of zero, the tails are problematic for  $\hat{F}_{n,m}^i$  and the estimators  $F_{n,m}^i$  tend to have smaller mse as demonstrated by the simulation study.























| 1      | Acl          | knowledgements  | 1      |
|--------|--------------|---|--------|
| 2      | - Th         |   | 2      |
| 3      | The          | work of the first author was partially supported by NSF Grant DMS-053246,   | 3      |
| 4<br>F | INSA         | A Grant H98230-06-1-0099, and NSF REU Grant DMS-0552590. The second   | 4<br>F |
| 5      | auth         | nor acknowledges support by project PROMEP/103.5/08/2988. The authors are   | 5      |
| 0      | grat         | eful to Drs Rudy Guerra of Rice University and Jonathan Cohen of University   | 6      |
| 1      | of 'I        | exas Southwestern Medical Center who made their data accessible. Their work   | (      |
| 8      | was          | supported by NIH grant RO1 HL-53917.  | 8      |
| 9      |              |   | 9      |
| 10     | Rof          | orences   | 10     |
| 11     | Itel         |   | 11     |
| 12     | [1]          | ARENAZ P. BITTICKS I. PANELL K. H. and GARCÍA S. (1992). Genotoxic potential of crown ethers  | 12     |
| 13     | [+]          | in mammalian cells; induction of sister chromatid exchanges. <i>Mutation Research</i> , <b>280</b> , 109–115.   | 13     |
| 14     | [2]          | BARTOSZEWICZ, J. (1985a). Moment inequalities for order statistics from ordered families of distri-   | 14     |
| 15     | [2]          | butions. Metrika, <b>32</b> , 383–389.<br>BAPTOSZEWICZ I (1985b) Dispersive ordering and monotone failure rate distributions. Adv. Appl.  | 15     |
| 16     | [0]          | Prob., 17, 472–47.  | 16     |
| 17     | [4]          | BARTOSZEWICZ, J. (1986). Dispersive ordering and the total time on test transformation. Statist.  | 17     |
| 18     | [5]          | Prob. Lett., 4, 285–288.<br>BICKEL P. L. and LEHMANN F. L. (1970). Descriptive statistics for nonparametric models. IV  | 18     |
| 19     | [0]          | Spread. In Contributions to Statistics, Jaroslav Hájek Memorial Volume, ed. J. Jureckova. Dor-  | 19     |
| 20     |              | drecht, Riedel, 33–40.  | 20     |
| 21     | [6]          | BIRNBAUM, Z. W. (1948). On random variables with comparable peakedness. Ann. Math. Statist.,  | 21     |
| 22     | [7]          | BROWN, G. and TUKEY, J. W. (1946). Some distributions of sample means. Ann. Math. Statist.,   | 22     |
| 23     |              | 7, 1–12.  | 23     |
| 24     | [8]          | DOKSUM, K. A. (1969). Starshaped transformations and the power of rank tests. The Annals of<br>Mathematical Statistica 40, 1167, 1176   | 24     |
| 25     | [9]          | EL BARMI, H. and ROJO, J. (1997). Likelihood ratio test for peakedness in multinomial populations.  | 25     |
| 26     |              | J. Nonparam. Statist., 7, 221–237.  | 26     |
| 27     | [10]         | EL BARMI, H. and MUKERJEE, H. (2008). Peakedness and peakedness ordering in symmetric distri-<br>butions. <i>Lowrnal of Multivariate Analysis</i> doi:10.1016/j.imva2008.06.011 (In Press)            | 27     |
| 28     | [11]         | ELSTON, R. C., BOXBAUM, S. and OLSON, M. (2000). Haseman and Elston Revisited. <i>Genetic Epi</i> -   | 28     |
| 29     | [10]         | <i>dem.</i> , <b>19</b> , 1–17.   | 29     |
| 30     | [12]         | EMBURY, S. H., ELIAS, L., HELLER, P. H., HOOD, C. E., GREENBERG, P. L. and SCHRIER, S. L. (1977).<br>Remission maintenance therapy in acute myelogenous lukemia. <i>Western Journal of Medicine</i> . | 30     |
| 31     |              | <b>126</b> , 267–272.   | 31     |
| 32     | [13]         | FRASER, D. A. S. (1957). Nonparametric Methods in Statistics. Wiley, New York.  | 32     |
| 33     | [14]         | HASEMAN, J. K. and ELSTON, R. C. (1972). The investigation of linkage between a quantitative trait<br>and a marker locus. <i>Behav. Genet.</i> , <b>2</b> , 2–19.                                     | 33     |
| 34     | [15]         | KARLIN, S. (1968). Total Positivity. Stanford University Press, CA.   | 34     |
| 35     | [16]         | LEHMANN, E. L. (1955). Ordered families of distributions. Ann Statist., 26, 399–419.  | 35     |
| 36     | [17]<br>[18] | LEHMANN, E. L. (1959). Testing Statistical Hypotheses. Wiley, New York.   | 36     |
| 31     | [19]         | LEHMANN, E. L. and Rojo, J. (1992). Invariant directional orderings. Ann Statist., <b>20</b> , 2100–2110.   | 37     |
| 38     | [20]         | LIU, B. H. (1988). Statistical Genomics, Linkage, Mappings, and QTL Analysis. CRC Press,  | 38     |
| 39     | [91]         | New York.<br>Lo S H (1987) Estimation of distribution functions under order restrictions. Statistics and  | 39     |
| 40     | [21]         | Decisions, 5, 251–262.  | 40     |
| 41     | [22]         | MARSHALL, A. W. and OLKIN, I. (2007). Life Distributions: Structure of Nonparametric, Semi-   | 41     |
| 42     | [23]         | parametric, and Parametric Families. Springer Science+Bussines Media, LCC, New York.  | 42     |
| 43     | [20]         | <b>33</b> , 303–312.  | 43     |
| 44     | [24]         | OJA, H. (1981). On location, scale, skewness, and kurtosis of univariate distributions. Scand. J.   | 44     |
| 45     | [25]         | Statist., 8, 154–168.<br>PROSCHAN, F. (1965). Peakedness of distributions of convex combinations. Ann. Math. Statist  | 45     |
| 46     | [20]         | <b>36</b> , 1703–1706.  | 46     |
| 47     | [26]         | Rojo, J. and HE, G. Z. (1991). New properties and characterizations of the dispersive ordering.   | 47     |
| 48     | [97]         | Statist. & Prob. Lett., 11, 365–372.<br>ROIO J. and WANG J. (1994). Test based on L-statistics to test the equality in dispersion of two  | 48     |
| 49     | [41]         | probability distributions. Statistics and Probability Letters. 21, 107–113.   | 49     |
| 50     | [28]         | ROJO, J. (1995a). On the weak convergence of certain estimators of stochastically ordered survival  | 50     |
| 51     |              | functions. J. Nonparam. Statist., 4, 349–363.   | 51     |

| 1      | [29]         | ROJO, J. (1995b). Nonparametric quantile estimation under order constraints, J. Nonparam.  | 1      |
|--------|--------------|--|--------|
| 2      | [30]         | Statist., 5, 185-200.<br>BOID I and MA Z (1996). On the estimation of stochastically ordered survival functions $I$  | 2      |
| 3      | [50]         | Statist. Compu. and Simu., 55, 1–21.   | 3      |
| 4      | [31]         | ROJO, J. (1998). Estimation of the quantile function of an IFRA distribution. Scand. J. Statist.,  | 4      |
| 5<br>6 | [32]         | <b>25.2</b> , 293–310.<br>ROJO, J. (1999). On the estimation of a survival function under a dispersive order constraint. J.  | 5      |
| 7      | [99]         | Nonparam. Statist., <b>11</b> , 107–135.   | 7      |
| 8      | ျပပျ         | First Erich L. Lehmann Symposium - Optimality, (J. Rojo, ed), IMS LNMS Vol 44, 37–61.  | ,<br>8 |
| 9      | [34]         | ROJO, J. and BATUN-CUTZ, J. (2007). Estimation of symmetric distributions subjects to peakedness   | 9      |
| 10     |              | order. Series in Biostatistics Vol 3, Advances in Statistical Modeling and Inference Chapter 13, 649–670   | 10     |
| 11     | [35]         | ROJO, J., BATUN-CUTZ, J. and DURAZO, R. (2007). Inference under peakedness restrictions. <i>Statistica</i>   | 11     |
| 12     | [a a]        | Sinica 17.3, 1165–1189.  | 12     |
| 13     | [36]         | SCHUSTER, E. (1975). Estimating the distribution function of a symmetric distribution. <i>Biometrika</i> , 62, 3, 631–635.   | 13     |
| 14     | [37]         | SCHWEDER, T. (1982). On the dispersion of mixtures. Scand. J. Statist., 9, 165–169.  | 14     |
| 15     | [38]         | SHAKED, M. (1980). On mixtures from exponential families. J. R. Statist. Soc. B., 42, 192–198.   | 15     |
| 16     | [39]<br>[40] | SHAKED, M. (1982). Dispersive orderings of distributions. J. Appl. Prob., 19, 510–520.<br>SHAKED, M. and SHANTIKUMAR, J. G. (2007) Stochastic Orders. Springer Science+Bussines Media, | 16     |
| 17     |              | LCC, New York.   | 17     |
| 18     | [41]         | SHIBATA, T. and TAKEYAMA, T. (1977). Stochastic theory of pitting corrosion. Corrosion, 33.7, 243.   | 18     |
| 19     |              |  | 19     |
| 20     |              |  | 20     |
| 21     |              |  | 21     |
| 22     |              |  | 22     |
| 23     |              |  | 23     |
| 24     |              |  | 24     |
| 25     |              |  | 25     |
| 26     |              |  | 26     |
| 27     |              |  | 27     |
| 28     |              |  | 28     |
| 29     |              |  | 29     |
| 30     |              |  | 30     |
| 31     |              |  | 31     |
| 32     |              |  | 32     |
| 33     |              |  | 33     |
| 34     |              |  | 34     |
| 35     |              |  | 35     |
| 30     |              |  | 30     |
| 20     |              |  | 31     |
| 30     |              |  | 30     |
| 40     |              |  | 40     |
| 40     |              |  | 40     |
| 42     |              |  | 42     |
| 43     |              |  | 43     |
| 44     |              |  | 44     |
| 45     |              |  | 45     |
| 46     |              |  | 46     |
| 47     |              |  | 47     |
| 48     |              |  | 48     |
| 49     |              |  | 49     |
| 50     |              |  | 50     |
| 51     |              |  | 51     |

# A Functional Generalized Linear Model with Curve Selection in Cervical **Pre-cancer Diagnosis using Fluorescence** Spectroscopy

# Hongxiao Zhu<sup>1</sup> and Dennis D. Cox<sup>2</sup>

## Rice University

Abstract: A functional generalized linear model is applied to spectroscopic data to discriminate disease from non-disease in the diagnosis of cervical precancer. For each observation, multiple functional covariates are available, and it is of interest to select a few of them for efficient classification. In addition to multiple functional covariates, some non-functional covariates are also used to account for systematic differences caused by these covariates. Functional principal components are used to reduce the model to multivariate logistic regression and a grouped Lasso penalty is applied to the reduced model to select useful functional covariates among multiple curves.

### Contents

| 1  | Introduction   |
|----|--|
| 2  | Functional Generalized Linear Model with Curve Selection               |
| 3  | Simulation Study   |
| 4  | Real Data Application–Fluorescence Spectral Curve Selection and Cervi- |
|    | cal Pre-Cancer Diagnosis   |
| 5  | Determining the Related Parameters                                     |
| 6  | Discussion   |
| Aŗ | ppendix: Proof of Proposition 1  |
| Ac | knowledgements   |
| Re | eferences  |

# 1. Introduction

Classification with functional data is a challenging problem due to the high dimensionality of the observation space. One solution is to reduce the dimension and use the reduced features for classification, such as the work of Hall et al. [6], Zhao et al. [19] and Ferré and Villa [5]). Another way is to use generalized linear regression by treating the class labels as responses and functional observations as predictors, which was proposed by James [8] and Müller and Stadtmüller [11]. Ratcliffe et al.

<sup>1</sup>Department of Statistics, Rice University, 6100 Main St. MS-138, Houston, Texas 77005, U.S.A., email: hxzhu@stat.rice.edu

AMS 2000 subject classifications: 60K35; secondary 60K37

<sup>&</sup>lt;sup>2</sup>Department of Statistics, Rice University, 6100 Main St. MS-138, Houston, Texas 77005, U.S.A., email: dcox@stat.rice.edu 

Keywords and phrases: Functional generalized linear model, curve selection, grouped lasso, fluorescence spectroscopy, cervical cancer



FIG 1. Left Panel: Fluorescence spectral curves at different excitation wavelengths. Right Panel: The image plot of fluorescence spectroscopy data (EEM).

[14] and Leng and Müller [9] applied this type of modeling to medical and gene expression data, respectively. Our basic concern in this study is the case when there are multiple functions per observation in a classification problem, and we wish to perform a curve selection to select few important curves and perform classification based on the selected curves.

The example that motivated our work is fluorescence spectroscopy data being investigated for cervical pre-cancer diagnosis. Fluorescence spectroscopy is an optical technique proposed for cervical pre-cancer screening. As a non-invasive, lowcost diagnosis tool, it provides a promising alternative to the existing methods for early-stage cancer diagnosis. One important step in this type of diagnosis is to discriminate the diseased observations from normal based on the high dimensional functional data — the fluorescence spectral measurements. In many clinical studies, several different spectra can be produced and used simultaneously for diagnosis ([12]), which makes the classification difficult since introducing more spectra not only provides more information but also more noise. Among these multiple spectral curves, it is suspected that some spectral curves contain more disease related information and hence are more "important" than others (see [3]). Furthermore, in order to produce an inexpensive commercial device, we would like to measure as few spectra as is necessary. This makes it beneficial to use statistical analysis to find out those curves that are good enough for diagnosis and remove the unnecessary ones, which can improve the diagnostic accuracy and reduce the cost.

The data studied in this paper are from a clinical study in which multiple fluores-cence spectra were measured at the same sites where biopsies were taken for patho-logical diagnosis. Each observation consists of several spectral curves measured in the following way: an excitation light at a certain fixed excitation wavelength is produced to illuminate the cervical tissue. The excitation light is absorbed by various endogenous fluorescent molecules in tissue, resulting in emission of fluorescent light. The emitted fluorescent light is measured by an optical detector and the spectrum is obtained as one smooth curve. The excitation light is varied at several different wavelengths and gives multiple spectral curves for each measurement. The 

imsart-coll ver. 2008/08/29 file: Hongxiao.tex date: March 25, 2009

left panel of Figure 1 shows the plot of all spectral curves from one measurement. Each measurement contains 16 spectral curves measured at excitation wavelengths З ranging from 330 nm to 480 nm with increments of 10 nm. Each spectral curve contains fluorescence intensities recorded on a range of emission wavelengths between 385nm and 700nm. If we use a color plot to represent the intensities, we can stack all the 16 spectra and obtain an image as shown in the right panel of Figure 1. We call such fluorescence spectroscopy measurements excitation-emission matrices (EEMs).

This study aims to select a subset of spectral curves from the 16 available curves for the purpose of classification. We will look at the problem from the functional data analysis ([13]) point of view and propose a functional generalized linear model, which will select among multiple functional predictors and perform binary classification. The proposed model allows both functional predictors and non-functional predictors. The non-functional predictors are variables associated with the measurements which may cause systematic difference in spectra, such as tissue type of the measurement site, or the menopausal status of patients.

The structure of this paper is as follows: Section 2 introduces the functional generalized linear model with curve selection and Section 3 provides a simulation study. The real data application to the fluorescence spectroscopy data is presented in Section 4, and details on determining related parameters are discussed in Section 5. A more general discussion is given in Section 6.

### 2. Functional Generalized Linear Model with Curve Selection

Consider *n* i.i.d. observations where each observation contains *J* functions. For i = 1, ..., n and j = 1, ..., J, let  $x_{ij}(t)$  denote the *j*th function observed from the *i*th observation, where  $E[x_{ij}(t)] = \mu_j(t)$ . Note that the *J* functions within each observation can be rather arbitrary hence we assume different mean function  $\mu_j(t)$  for each  $x_{ij}(t)$ . In addition to functional data, we assume there is a non-functional vector  $z_i$  associated with each observation. Suppose the responses we observed are binary variables  $y_i$ . Similarly to James [8] and Müller and Stadtmüller [11], we propose a functional generalized linear model to connect the binary responses with the predictors. Let  $p_i = Pr(y_i = 1|z_i, x_{ij}(t), j = 1, ..., J)$  and

(2.1) 
$$p_i = g^{-1}(\eta_i)$$

(2.2) 
$$\eta_i = \alpha_0 + z_i^T \alpha + \sum_{j=1}^J \int_{T_j} \beta_j(t) (x_{ij}(t) - \mu_j(t)) dt$$

where  $T_j$  is the domain of  $x_{ij}(t)$ ,  $\alpha_0$  is the univariate intercept,  $\alpha$  is a vector of coefficients for the non-functional predictors, and the  $\beta_j(t)$ 's are the functional regression coefficients. For convenience, we center  $x_{ij}(t)$  at its mean in the integrand. Here the link function  $g(\cdot)$  is a one-to-one continuous function. To perform curve selection, we propose the following constraint on the functional regression coefficients:

(2.3) 
$$\sum_{j=1}^{J} ||\beta_j||_{L^2} < s$$

where  $||f||_{L^2} = (\int f^2(t)dt)^{1/2}$ , s is a pre-defined constant. Note that (2.3) is a combined constraint of  $L^2$  norm and  $l^1$  norm. This is an extension of the groupwise variable selection in the multivariate setting proposed by Yuan and Li [18].

Because of the properties of this combined constraint, we expect  $\beta_j \equiv 0$  for a number of j's, depending on the value of s.

Due to the infinite dimensionality of functional data, multivariate methods can not be used directly for solving the above proposed model. One can discretize  $x_{ij}(t)$ on a finite grid and transform the problem to a multivariate regression model, but the number of grid points is an issue and there will be high correlation between contiguous grid points because of the "functional" properties of  $x_{ij}(t)$ . A natural choice is to apply standard functional dimension reduction methods to reduce the dimension first and solve the problem on the reduced space. If we assume  $\forall j, x_{ij}(t) \in$  $\mathcal{H}_j$  for some separable Hilbert space  $\mathcal{H}_j$ , and  $E[x_{ij}(t)] = \mu_j(t)$ , we can expand  $x_{ij}(t) - \mu_j(t)$  on a set of orthonormal basis  $\{\phi_k^j\}_{k=1}^{\infty}$ 

(2.4) 
$$x_{ij}(t) - \mu_j(t) = \sum_{k=1}^{\infty} c_{ijk} \phi_k^j(t)$$

and a truncated version of (2.4) can be used to approximate  $x_{ij}(t)$  since  $\sum_{k=1}^{\infty} |c_{ijk}|^2 < \infty$ . And similarly, we assume  $\beta_j(t) \in \mathcal{H}_j, \forall j$ , and this gives

(2.5) 
$$\beta_j(t) = \sum_{k=1}^{\infty} b_{jk} \phi_k^j(t)$$

Note that the orthonormal basis  $\{\phi_k^j\}_{k=1}^{\infty}$  can be chosen to be a known basis such as a Fourier basis or a wavelet basis. If in addition, we assume  $x_{ij}(t) \in L_2[\Omega \times T_j]$  for the domain  $T_j$  and the underlying sample space  $\Omega$ , i.e.,  $\int_{T_j} E[x_{ij}(t)^2] < \infty, \forall j$ , Mercer's theorem and Karhunen-Loève theorem ([2]) suggest taking the orthonormal basis to be the eigenfunctions of the covariance operator K, where K is defined by

(2.6) 
$$Kx(t) = \int x(s)k(s,t)ds, \qquad k(s,t) = Cov(x(s),x(t)).$$

In this case, the coefficients  $\{c_{ijk}, k = 1, ..., \infty\}$  are called functional principal component scores of the functional data. Using the functional principal component method is different from using a known basis in that the eigenbasis functions need to be estimated. Various estimating methods are proposed as in Ramsay and Silverman [13], and in Hall, Müller and Wang [7].

Once the functional principal component scores or the orthonormal basis coefficients have been estimated, we can reduce equation (2.2) to

(2.7) 
$$\eta_{i} = \alpha_{0} + z_{i}^{T} \alpha + \sum_{j=1}^{J} \sum_{k=1}^{\delta_{j}} c_{ijk} b_{jk}$$

where  $\delta_j$  is the truncation parameter for the *j*th functional predictor. We thus transfer the functional regression to multivariate regression. The constraint condition (2.3) will be reduced to

(2.8) 
$$\sum_{j=1}^{J} ||b_j||_2 < t$$
46
47
47
48

where  $b_j = (b_{j1}, \ldots, b_{j\delta_j})$  and  $|| \cdot ||_2$  stands for the Euclidean norm. Curve selection can thus be performed through selecting variables in (2.7) using the grouped

H. Zhu and D.D. Cox

Lasso type constraint (2.8), i.e., if one curve  $x_i(t)$  is selected, then the coefficients  $b_{jk}, k = 1, \ldots, \delta_j$ , will all be non-zero. The Lasso (Least Absolute Shrinkage and Se-lection Operator) was first proposed by Tibshirani [16] for model selection in linear regression models. The basic idea was to find a subset of the covariates with non-zero coefficients by applying an  $l_1$  constraint to the regression coefficients based on the ordinary least square estimation. Yuan and Lin [18] extended the regular Lasso to cases where the covariates can be grouped, such as multi-factor ANOVA. They combine the  $l_1$  and  $l_2$  constraints so that the resulting model selects variables at the group level and is invariant under group-wise orthogonal transformation. To solve our problem based on the reduced model (2.7) and (2.8), we borrow the algorithm proposed by Meier et al. [10], where they extend the group-wise lasso regression of Yuan and Lin [18] to a logistic regression setup. Assume the link function in (2.1)is a logit link, i.e., 

(2.9) 
$$\log(\frac{p_i}{1-p_i}) = \eta_i$$
14
15
16

The estimate can be obtained by minimizing the convex function

$$l( heta) = -l( heta) + \lambda \sum_{j=1}^{J} s(\delta_j) ||b_j||_2$$

(2.10) 
$$Q_{\lambda}(\theta) = -l(\theta) + \lambda \sum_{j=1} s(\delta_j) ||b_j||_2$$

where  $\theta = \{\alpha_0, \alpha, b_j, j = 1, ..., J\}$ , and  $l(\cdot)$  is the log-likelihood function:

(2.11) 
$$l(\theta) = \sum_{i=1}^{n} \{y_i \eta_i - \log(1 + \exp(\eta_i))\}$$
<sup>24</sup>  
25  
26

Here  $s(\delta_j)$  is used to rescale the penalty with respect to the dimensionality of  $b_j$ , usually taken to be  $\sqrt{\delta_j}$ , and  $\lambda > 0$  is the tuning parameter to control the amount of penalty. Note that in the model of Meier et al. [10], they only allow one unpenalized term, i.e., only the intercept term is unpenalized. In our proposed model, in addition to the intercept  $\alpha_0$ , we allow the coefficients  $\alpha$  of nonfunctional predictors to be unpenalized. Meier et al. stated the attainability of the minimum of the optimization problem in their paper and provided a proof. Actually, some conditions must be satisfied for the attainability to hold. Here we provide a general sufficient condition for the minimum of Equation (2.10) to be attained.

**Proposition 1.** For  $0 < \sum_{i=1}^{n} y_i < n, \lambda > 0, s(\delta_j) > 0, \forall j$ , assume the design matrix X formed by

$$X = \begin{pmatrix} 1 & z_1^T & c_{111} & \dots & c_{11\delta_1} & \dots & \dots & c_{1J1} & \dots & c_{1J\delta_J} \\ \vdots & & & & & & \\ 1 & z_n^T & c_{n11} & \dots & c_{n1\delta_1} & \dots & \dots & c_{nJ1} & \dots & c_{nJ\delta_J} \end{pmatrix}$$

is an n by m matrix of rank m, 
$$n \ge m$$
. Assume the maximum likelihood estimator  
for the logistic regression (with log-likelihood in Equation (2.11)) exists. Then the  
Equation (2.10) has an unique minimizer  $\theta^*$ .

The proof for Proposition 1 is in the Appendix. Meier et al. [10] proposed a Block Coordinate Gradient Descent algorithm to solve the group lasso logistic regression and provided an R package called *grplasso*. We will use this package to perform curve selection based on reduced model in equations (2.7) and (2.8). The initiation of the algorithm is the same as in *grplasso*.

imsart-coll ver. 2008/08/29 file: Hongxiao.tex date: March 25, 2009

# 3. Simulation Study

To verify the performance of the proposed method in classification problems with multiple functional covariates, we generate n = 1000 i.i.d. observations. Each observation contains one non-functional covariate and three functional covariates. The non-functional covariate is generated from uniform (0, 1) distribution. And the three functional covariates are generated using the first 4 cosine basis functions on the domain [0, 1], i.e., using basis  $\phi_0(t) = 1, \phi_k(t) = \sqrt{2}\cos(k\pi t), k = 1, \dots, 3$ . For each functional covariate, the 4 coefficients of the cosine basis are generated independently from a normal distribution with some fixed mean and variance 0.5. We set the coefficient functions for the first and third functional covariates to be zero and set the coefficient function for the second to be non-zero. Figure 2 shows the plot of both non-functional covariates and functional covariates for the first 50 observations. The binary responses  $y_i$  are generated by sampling from a Bernoulli distribution with success probability  $p_i = (1 + \exp(-\eta_i))^{-1}$ , where  $\eta_i$  is computed from equation (2.2) using numerical integration. The proportion of 1's among the binary  $y_i$ 's is 57.3%. The data are randomly split into a training set and a test set with 800 observations in the training set and 200 observations in the test set.

To apply the proposed model to these data, one can choose different types of orthonormal basis for dimension reduction. Since the data are generated using cosine basis, we will show the results of using the cosine basis so that the estimated coefficients can be compared with their known true values. We have also tried using functional principal components, and the curve selection and prediction results are very similar to that of using cosine basis.

For the choice of cosine basis, we reduce the dimension of the functional pre-dictors using the first 4 cosine basis functions. The group-wise lasso regression algorithm of Meier et al. [10] is then applied to the reduced scores. Figure 3 shows the estimation paths for the regression coefficients as a function of  $\lambda$ . Note that for the estimated coefficient function  $\hat{\beta}_j$ , we plotted their  $L^2$  norm, i.e.,  $||\hat{\beta}_j|| =$  $\sqrt{\int_{T_i} \hat{\beta}_j(t)^2 dt}$ , where the function  $\hat{\beta}_j$  are obtained through inverse transform of the estimated coefficients  $\hat{b}_i$ . From Figure 3, we see that for a large range of  $\lambda$ , i.e.,  $15.7 < \lambda < 115$ , the method correctly picked out the non-zero coefficient function  $\beta_2$ . The values of  $\beta_2(t)$  at 6 selected  $\lambda$ 's is plotted in Figure 4 in comparison with the true  $\beta_2(t)$ . Table 1 shows the estimated coefficients (in form of the cosine basis scores  $\hat{b}_i$  compared with the true values under the 6 selected  $\lambda$ 's. From Table 1, we see that as the penalty parameter  $\lambda$  increases, the magnitudes of the estimated coefficients shrink toward 0. When  $\lambda = 0$ , the estimates are equal to the maxi-mum likelihood estimates, which gives nonzero estimates to all coefficients. When  $\lambda$  ranges from 22.4 to 89.6, the coefficients corresponding to the first and third curve are exactly 0, and the coefficients corresponding to the second curve are nonzero. For  $\lambda > 14.1$ , the estimates are (almost all) closer to 0 than the true values. We believe that these shrinkage effects are caused by the continuous-shrinkage property of Ridge and Lasso penalty (see Tibshirani [16]). It has been suggested that there may be large bias in the estimators related to the inconsistency of the original Lasso under certain conditions, i.e., that the Lasso does not satisfy the "oracle proper-ties" (Fan and Li [4], Zhao and Yu [20]). Some modifications have been proposed to overcome the drawbacks of Lasso and make the estimators satisfy the oracle properties (see Zou [22]). In this study, we only focus on the curve selection and prediction, but more research can be done on the consistency of the grouped-Lasso regression under the functional data setup. 

imsart-coll ver. 2008/08/29 file: Hongxiao.tex date: March 25, 2009



FIG 2. Data plot of both non-functional covariates and functional covariates for the first 50 observations used in simulation.

To perform prediction on test set, the estimated coefficient function  $\hat{\beta}_i(t), j =$ 1, 2, 3 are plugged into the test set using (2.2) and the estimated success probability  $\hat{p}_i$  are computed for each observation, from which we can plot a ROC curve(see Zweig & Campbell [23]) for each  $\lambda$ . From each ROC curve, we pick a point that maximizes the sum of sensitivity and specificity, and this point will be used as the optimal classification point. The misclassification rate at the optimal point and the corresponding area under the ROC curves are computed at different values of  $\lambda$ and plotted in Figure 5. From Figure 5, we find that when  $\lambda$  is around 22.4, the prediction on the test set gives the best sensitivity (93%) and specificity (73%) and an fairly large area under ROC curve (0.88), and the corresponding misclassification rate is 16%. 

imsart-coll ver. 2008/08/29 file: Hongxiao.tex date: March 25, 2009



FIG 3. Estimated paths of coefficient vector at different  $\lambda$  values

# 4. Real Data Application–Fluorescence Spectral Curve Selection and Cervical Pre-Cancer Diagnosis

Totally 717 EEM measurements were made on 306 patients, and each measurement contains 16 spectral curves. Measurements were taken from different sites on the cervix and may include repeated measurements at the same site. All the measure-ments were made using the same instrument (called FastEEM3) in the same clinic (British Columbia Cancer Agency, Vancouver, CA). Data were split into a training set and a test set with 396 measurements in the training set and 321 in the test set. The proportions of diseased cases within each set are 0.21, 0.20, respectively. Two non-functional covariates are considered in this study: the colposcopic tissue type of the measurements, and the menopausal status. Colposcopic tissue type is a binary variable indicating two types of tissue — squamous and columnar, which is obtained prior the fluorescence spectroscopy measurements. Menopausal status of a patient is a categorical variable which has three levels: pre-, peri- and post-menopause. The first 5 functional principal components are chosen as the scores extracted from each functional predictor, which reduce the data to a total of 80 scores. To reduce bias, the test set scores (the scores of orthonormal basis) are computed based on only information obtained from the training set. For example, the eigenfunctions used for computing functional principal components scores of the test set are estimated from the training set. 

<sup>49</sup> The grouped lasso logistic regression is used to pick the excitation wavelengths <sup>50</sup> as  $\lambda$  decreases from 6 to 0. Due to the large number of curves, the plot of coefficient <sup>51</sup> path is hard to visualize. In Figure 6, we summarize the excitation spectral curves



imsart-coll ver. 2008/08/29 file: Hongxiao.tex date: March 25, 2009

|            |             | E               | stimated c       | oefficients      | at differen      | t $\lambda$ values | 3             |
|------------|-------------|-----------------|------------------|------------------|------------------|--------------------|---------------|
| Coef       | True Values | $\lambda = 118$ | $\lambda = 89.6$ | $\lambda = 22.4$ | $\lambda = 14.1$ | $\lambda = 5.3$    | $\lambda = 0$ |
| $\alpha_0$ | 0.5         | 0.3             | 0.3              | 0.39             | 0.42             | 0.46               | 0.5           |
| $\alpha$   | 1           | 0.63            | 0.64             | 0.82             | 0.87             | 0.97               | 1.06          |
| $b_{11}$   | 0           | 0               | 0                | 0                | 0                | 0.03               | 0.15          |
| $b_{12}$   | 0           | 0               | 0                | 0                | 0                | -0.04              | -0.17         |
| $b_{13}$   | 0           | 0               | 0                | 0                | 0                | 0.04               | 0.18          |
| $b_{14}$   | 0           | 0               | 0                | 0                | 0                | 0                  | -0.0          |
| $b_{21}$   | 1           | 0               | 0.13             | 0.58             | 0.67             | 0.79               | 0.9           |
| $b_{22}$   | 2           | 0               | 0.31             | 1.43             | 1.67             | 2.01               | 2.29          |
| $b_{23}$   | -3          | 0               | -0.42            | -1.92            | -2.24            | -2.66              | -3.02         |
| $b_{24}$   | -1          | 0               | -0.18            | -0.84            | -0.99            | -1.21              | -1.41         |
| $b_{31}$   | 0           | 0               | 0                | 0                | 0                | 0.02               | 0.03          |
| $b_{32}$   | 0           | 0               | 0                | 0                | 0.01             | 0.07               | 0.13          |
| $b_{33}$   | 0           | 0               | 0                | 0                | 0.04             | 0.34               | 0.56          |

0.01

0.09

0.14

the excitation curves are ordered by 340 > 460 > 420 > 410 according to the order of being selected. The estimated coefficients at different values of  $\lambda$  are used to predict in the test set, from where we can evaluate the performance of different  $\lambda$  values. Due to the fact that the total proportion of diseased cases is small, the misclassification rate is not an ideal criterion for evaluating the prediction result (see [21], page 22 for details). To reduce the risk of false negatives, we wish to keep the sensitivity high enough and sacrifice some specificity. Hence for each fixed  $\lambda$ , we pick a point from the empirical ROC curve using the criterion that the sum of the sensitivity and specificity is maximized. The Figure 7 shows the area under the curves and the optimal sum of the sensitivity and specificity at different values of  $\lambda$ . At  $\lambda = 1.64$ , the sum is maximized at 1.44 with sensitivity 87% and specificity 57%, and the corresponding area under ROC curve is 0.77, and misclassification rate is 38%.

Since the main purpose of the above analysis is for curve selection rather than classification, once the functional covariates are selected, different classifiers can be applied to perform classification based on the selected subset of curves. In addition to logistic regression, we also performed classification with 3 other classifiers using the selected curves. By choosing  $\lambda = 1.64$ , we selected function predictor curves at excitations: 330, 340, 360, 370, 410, 420, 460 and 480, and used the first 5 functional principal components to reduce the dimension. We refitted the logistic model without penalty and compared the prediction results on the test set with 3 other classifiers in Table 2. The corresponding ROC curves are plotted in Figure 8. From Figure 8, we find that logistic regression, k-nearest neighbor(KNN) and linear discriminant analysis(LDA) provide similar ROC curves. The highest sum of sensitivity and specificity is 1.43, obtained by KNN, which is only slightly smaller than the grouped lasso results at  $\lambda = 1.64$ . The LDA method provides the same specificity with logistic regression but higher sensitivity.

З

### 5. Determining the Related Parameters

 $b_{34}$ 

In our proposed model, two types of parameters need to be determined: the tuning parameter  $\lambda$  and the truncation parameters  $\delta_j$ , j = 1, ..., J. In this section, we discuss how to determine these parameters. The choice of tuning parameter  $\lambda$  is important for prediction. In Meier et al. [10]

З



formance. However, there are also cases where only a small number of observations are available and splitting out a test set is not possible. In this case, we can adopt model selection criteria such as AIC, practical  $C_p$  or BIC. AIC tends to select a model with optimal prediction, while BIC tends to identify the true sparse model if the true model is included in the candidate set (see Yang [17]). In the grouped Lasso linear regression model, Yuan and Lin [18] proposed an approximation to the degrees of freedom and used a  $C_p$  criterion for selecting the tuning parameter  $\lambda$ . Whether this criterion can be extended to logistic regression case for selecting  $\lambda$  is an open question.

In addition to the tuning parameter  $\lambda$ , the truncation parameter  $\delta_j$  in equation (2.7) is also one concern of the study. In the real application of Section 4, we let

| TABLE | <b>2</b> |
|-------|----------|
|-------|----------|

The classification results using 4 different methods on the selected curves. Auc: Area under ROC curve. MisR: Misclassification rate. Sens: Sesitivity. Speci: Specificity. Sum: The sum of sensitivity and specificity. Logistic: logistic regression. KNN: k-nearest neighbor. LDA: linear discriminant analysis. SVM: support vector machine.

| Method   | Auc  | MisR | Sens | Speci | Sum  |
|----------|------|------|------|-------|------|
| Logistic | 0.76 | 31%  | 71%  | 68%   | 1.39 |
| KNN      | 0.68 | 27%  | 68%  | 74%   | 1.43 |
| LDA      | 0.75 | 31%  | 75%  | 68%   | 1.42 |
| SVM      | 0.64 | 28%  | 48%  | 79%   | 1.26 |



FIG 7. Prediction results at different  $\lambda$  values.

 $\delta_j \equiv \delta$  and reported the curve selection and prediction results with  $\delta = 5$ . To find out whether other choices of  $\delta$  are better for prediction, we compute the prediction results for the test set at different number of  $\delta$  but fixing  $\lambda = 1.64$ . The quantitative prediction results are plotted in Figure 9. From Figure 9, we can see that using 11 functional principal components, the area under the ROC curve are maximized at 0.780, and the sum of sensitivity and specificity are maximized at 1.47, with a relatively small misclassification rate 31%. The sensitivity and specificity reach 81%, 66%, respectively.

It is also suspected that the optimal  $\lambda$  may interact with  $\delta$  so determining one by fixing the other may be suboptimal. In our study, we also have tried to determine both the parameters by training the model under different combinations of them, and predicting on the test set. It turns out that at around  $\lambda = 1.64$  the prediction results of the model is better than other choices of  $\lambda$ , and this is quite stable across different choice of  $\delta$ , especially for  $\delta$  greater than 3. In Figure 10, We plot the area under the ROC curve for 11 different  $\delta$  and for appropriately selected  $\lambda$ values across a meaningful range, i.e.,  $\lambda = (5, 3, 1.64, 1.5, 1, 0.27)$ . It shows that the line with  $\lambda = 1.64$  stays on the top for  $\delta$ 's larger than 3. The reason for the small interaction between  $\lambda$ 's and  $\delta$ 's can be the following: the orthogonal basis approximation tends to be accurate with only a few components. For example, in functional principle components, over 97% of the variability will be counted in the first priciple component score for all excitation curves. Later components only add details to the model but does not change the likelihood dramatically. Therefore the minimum of Equation (2.10) as a function of  $\lambda$  does not change much when  $\delta$ changes. But this is not true for non-orthogonal basis approximation methods such as B-spline. 

<sup>46</sup> Note that choosing  $\delta_j \equiv \delta$  is just a convenient choice, which has the advantage <sup>47</sup> that it leaves only two parameters to determine and cross validation is feasible <sup>48</sup> for determining these parameters. However, it also brings in the risk of loosing <sup>49</sup> information. In general, one may use different truncation parameters if there are <sup>50</sup> large differences on the properties of the curves such as smoothness. If all curves <sup>51</sup> are obtained through similar sources and are similar in shape and other above

imsart-coll ver. 2008/08/29 file: Hongxiao.tex date: March 25, 2009

З



FIG 8. ROC curves obtained when training 4 different classifiers based on selected curves and predicting on the test set.

mentioned properties, it would be safe to choose a common  $\delta$ . As an alternative, since the step of estimating  $\{c_{ijk}, k = 1, \ldots, \delta_j, j = 1, \ldots, J\}$  can be independent of the group-wise Lasso step, one can use approximation criteria such as error sum of squares(SSE) to determine the truncation parameters for each curve. For example, if using functional principal component, we can choose a level of approximation(e.g., let the percentage of variabilities explained to be greater than 99%) and select the number of eigenfunctions to achieve this. However, better approximation does not necessarily give better prediction.

### 6. Discussion

We have proposed a functional logistic regression model to perform classification and curve selection. This model automatically selects among the functional covariates through the grouped Lasso variable selection. The proposed model gives information about which curves will be selected if we are willing to use a subset of the functional covariates for classification. For example, under penalty  $\lambda = 5$ , the best four functional predictors selected in our real data application are curves at excitations 340, 410, 420 and 460. The selected functional covariates can then be used with different classifiers for accurate classification.

There are several aspects that can be studied in more detail. Firstly, the basis expansion step can be combined more tightly with the grouped Lasso regression step using techniques similar to Müller and Stadtmüller [11]. It is necessary to investigate the consistency properties of the estimated coefficient function  $\beta_i(t)$ 's,



such as the oracle property. The algorithm of Meier et al. [10] requires that the tuning parameter  $\lambda$  to be predefined on a grid of values, where they proposed a way to find the range of  $\lambda$  of interest. This method, although faster, makes it difficult to find a precise  $\lambda$  value that is optimal for prediction purposes. Efficient algorithms for searching  $\lambda$  is of great importance especially when functional data is involved.

Alternative methods for curve selection can be formulated through the Bayesian paradigm. Bayesian variable selection models can be derived for selecting variables at a group level and thus can be used for curve selection as well.

# Appendix: Proof of Proposition 1

The proof of Proposition 1 uses a result stated in the following lemma.

**Lemma 1.** Let  $f : \mathbb{R}^n \to \mathbb{R}$  be a strictly convex function with a minimizer  $\tilde{x}$ , and 49  $let g : \mathbb{R}^n \mapsto [0, \infty)$  be a convex function. Then f + g has a unique minimizer  $x^*$  in  $\mathbb{R}^n$ .

*Proof.* Let h(x) = f(x) + g(x). It is easy to show that h(x) is strictly convex from

The Area Under ROC Curve v.s. Num. of PC

З



FIG 10. The area under the ROC curves for 6 different  $\lambda$  values and 11 different choice of FPC's.

the definition. We claim that the existence of a minimizer  $\tilde{x}$  of f implies that h is coercive, which means  $h(x) \to \infty$  as  $||x|| \to \infty$ . The coerciveness and strict convexity of h implies the existence of a unique minimizer  $x^*$ .

To show that h is coercive, it is sufficient to show that f is coercive (since  $g \ge 0$ ). The minimizer  $\tilde{x}$  of f is the unique minimizer of f by strict convexity. Also, f is convex hence is continuous on  $\mathbb{R}^n$  (see [15], page 82). Thus  $\forall r > 0, \forall x$  such that  $||x - \tilde{x}|| > r$ , we claim

$$f(x) > \frac{b}{r} ||x - \tilde{x}|| + f(\tilde{x})$$

where  $b = \inf\{f(x) : ||x - \tilde{x}|| = r\} - f(\tilde{x})$ . Note that b exists and b > 0 by continuity of f. To show this inequality, let  $x_0 = r(x - \tilde{x})/(||x - \tilde{x}||) + \tilde{x}$ , so that  $x_0$  lies on the line formed by x and  $\tilde{x}$ , with  $||x_0 - \tilde{x}|| = r$  and  $||x - x_0|| = ||x - \tilde{x}|| - r$ . Thus  $f(x_0) - f(\tilde{x}) \ge b$  by the definition of b. Now let  $\alpha = r/||x - \tilde{x}||$ . We see that  $x_0 = \alpha x + (1 - \alpha)\tilde{x}$ . By strict convexity of f,

$$f(x_0) < \alpha f(x) + (1 - \alpha)f(\tilde{x})$$

Thus

$$\frac{b}{r}||x-\tilde{x}|| + f(\tilde{x}) \le (f(x_0) - f(\tilde{x}))\frac{||x-\tilde{x}||}{r} + f(\tilde{x})$$

=

$$< (\alpha f(x) + (1 - \alpha)f(\tilde{x}) - f(\tilde{x}))\frac{||x - \tilde{x}||}{r} + f(\tilde{x})$$

49 Since  $||x - \tilde{x}|| \ge ||x|| - ||\tilde{x}||$ ,  $||x|| \to \infty$  implies  $||x - \tilde{x}|| \to \infty$ , which implies 50  $f(x) \to \infty$  by the above inequality and the facts that  $b > 0, r > 0, f(\tilde{x})$  finite. 51 Therefore, f is coercive, and so is h.

Since h is coercive, we have  $h(x) \to \infty$  as  $||x|| \to \infty$ . Therefore, if we pick an arbitrary point  $x_1 \in \mathbb{R}^n$ , there exists a constant  $\delta > 0$  such that  $h(x) > h(x_1)$  for all  $||x - x_1|| > \delta$ . Since the domain  $||x - x_1|| \le \delta$  is compact and h(x) is strictly convex on it, h(x) has a unique minimizer in  $||x - x_1|| \le \delta$ , which we denote as  $x^*$ . (A strictly convex real valued function defined on a compact domain has a unique minimum on its domain.) This  $x^*$  is also the global minimizer since  $h(x) > h(x_1) \ge h(x^*)$  on  $||x - x_1|| \ge \delta$ .

Proof of Proposition 1: Based on results in Lemma 1, we let f to be  $-l(\theta)$  and g to be  $\lambda \sum_{j=1}^{J} s(\delta_j) ||b_j||_2$ , therefore our objective function in Equation (2.10) is the sum of f and g, where  $\theta = \{\alpha_0, \alpha, b_j, j = 1, ..., J\}$ , and  $l(\theta) = \sum_{i=1}^{n} y_i \eta_i - \log(1 + \exp(\eta_i))$  with  $\eta_i = \alpha_0 + z_i^T \alpha + \sum_{j=1}^{J} \sum_{k=1}^{\delta_j} c_{ijk} b_{jk}$ . Firstly, we show that  $-l(\theta)$  is strictly convex. It is sufficient to show that its

Firstly, we show that  $-l(\theta)$  is strictly convex. It is sufficient to show that its Hessian is positive definite. Since the Hessian takes the form

$$\nabla^2_{\theta}(-l(\theta)) = X^T D X$$

where  $D = \text{diag}\{\exp(\eta_i)/(1 + \exp(\eta_i))^2, i = 1, ..., n\}$ . It is positive definite since X is of rank m (full rank). Secondly, since the maximum likelihood estimator exists,  $-l(\theta)$  has an unique minimizer. The existence of maximum likelihood estimator for logistic regression requires some conditions for the design matrix X. Basically, the n rows of X can not be completely separated or quasi-completely separated in  $\mathbb{R}^m$ . See [1] for details. In practice, as long as we can find a numerical solution for the MLE at  $\lambda = 0$ , we would believe that the maximum likelihood estimator exists. Finally, let  $g(b) = \lambda \sum_{j=1}^{J} s(\delta_j) ||b_j||_2, b^T = (b_1^T, \ldots, b_J^T)$ . It is easy to see that g(b)is convex by the triangle inequality. Therefore by Lemma 1,  $Q_{\lambda}(\theta)$  has a unique minimizer  $\theta^*$ .

# Acknowledgements

This research was supported by the National Cancer Institute grant PO1-CA82710 and by the National Science Foundation grant DMS0505584. We thank the referees for constructive comments, and thank Dr. Wotao Yin for helpful discussions on the convex optimization.

### References

| [1] | ALBERT, A and ANDERSON, J.A. (1984). On the existence of maximum likelihood estimates in logic | $_{\rm stic}$ |
|-----|--|---------------|
|     | regression models. Biometrika, <b>71</b> , 1–10.   |               |
|     |  |               |

- [2] ASH, R. B. (1975). Topics in Stochastic Processes. Academic Press, New York.
- [3] CHANG, S. K., FOLLEN, M., MALPICA, A., UTZINGER, U., STAERKEL, G., COX, D., ATKINSON, E. N., MACAULAY, C. and RICHARDS-KORTUM, R. (2002). Optimal excitation wavelengths for discrimination of cervical neoplasia. *IEEE Transactions on Biomedical Engineering*, 49, 1102–1110.
- [4] FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. J. Amer. Statist. Assoc., 96, 1348–1360.
- [5] FERRÉ, L. and VILLA, N. (2006). Multilayer perceptron with functional inputs: an Inverse Regression Approach. the Scandinavian Journal of Statistics, 33, 807–823.
- [6] HALL, P., POSKITT, D. S. and PRESNELL, B. (2001). A Functional data-analytic approach to signal discrimination. *Technometrics*, 43, 157–214.
- [7] HALL, P., MÜLLER, H. and WANG, J. (2006). Properties of principal component methods for functional
   and longitudinal data analysis. Ann. Statist, 34, 1493–1517.

З

<sup>[8]</sup> JAMES, G. M. (2002). Generalized linear models with functional predictors. J. R. Stat. Soc. Ser. B,
64, 411-432.
[9] LENG, X. and Müxers, H. (2007). Chariford in an informational hold and hold and

 <sup>50
 [9]</sup> LENG, X. and MÜLLER, H. (2005). Classification using functional data analysis for temporal gene
 50

 51
 expression data. Bioinformatics, 22, 68–76.
 51

# H. Zhu and D.D. Cox

| 1        | [10] | MURD I CHER & and DÜWLYLYDY D (2008) The mean Lasse for legistic remeasion I D Stat  | 1      |
|----------|------|--|--------|
| 2        | [10] | Soc. Ser. B, <b>70</b> , 53–71.  | 1<br>2 |
| 2        | [11] | MÜLLER, H. and STADTMÜLLER, U. (2005). Generalized functional linear models. Ann. Statist, 33,   | 2      |
| 4        | [19] | 774–805.<br>Ramanulam N. Mitchell M. F. Mahadevan A. Thomsen S. Maldica, A. Weicht T. Atvin  | 4      |
| 5        | [12] | SON, N. and RICHARDS-KORTUM, R. (1996). Spectorscopic diagnosis of cervical intraepithelial neopla-  | 5      |
| 6        |      | sia(CIN) in vivo using laser induced fluorescence spectra at multiple excitation wavelengths. Lasers   | 6      |
| 7        | [19] | Surg. Med., 19, 63–67.   | 7      |
| ,<br>8   | [13] | RATCLIFFE, S. J., HELLER, G. Z. and LEADER, L. R. (2002). Functional data analysis with application  | 8      |
| 9        |      | to periodically stimulated foetal heart rate data. II: Functional logistic regression. Statistics in   | 9      |
| 10       | [15] | Medicine, <b>21</b> , 1115–1127.<br>ROCKAPELLAR R. T. (1970). Conver Analysis, Princeton University Press  | 10     |
| 11       | [16] | TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. J. R. Stat. Soc. Ser. B,  | 11     |
| 12       |      | <b>58</b> , 267–288.   | 12     |
| 13       | [17] | YANG, Y. (2005). Can the strengths of AIC and BIC be shared? A conflict between model identifi-<br>cation and regression estimation. <i>Biometrika</i> , <b>92</b> , 937–950 | 13     |
| 14       | [18] | YUAN, M. and LIN, Y. (2006). Model selection and estimation in regression with grouped variables.  | 14     |
| 15       | [10] | J. R. Stat. Soc. Ser. B, 68, 49–67.  | 15     |
| 16       | [19] | ZHAO, X., MARRON, J. S. and WELLS, M. 1. (2004). The functional data analysis view of longitudinal data. <i>Statistica Sinica</i> , 4, 789–808.                              | 16     |
| 17       | [20] | ZHAO, P. and YU, B. (2006). On model selection consistency of Lasso. J. Machine Learning Re-   | 17     |
| 18       | [01] | search, 7, 2541–2563.  | 18     |
| 19       | [21] | & Sons, Inc., New York.  | 19     |
| 20       | [22] | ZOU, H. (2006). The adaptive Lasso and its oracle properties. J. Amer. Statist. Assoc., 101, 1418-   | 20     |
| 21       | [99] | 1429.<br>ZUERC M. H. and CLUEDERL, C. (1002). Dessive accepting characteristic (DOC) plate, a funda  | 21     |
| 22       | [23] | mental evaluation tool in clinical medicine. <i>Clinical Chemistry</i> , <b>39</b> , 561–577.  | 22     |
| 23       |      |  | 23     |
| 24       |      |  | 24     |
| 25       |      |  | 25     |
| 26       |      |  | 26     |
| 27       |      |  | 27     |
| 28       |      |  | 28     |
| 29       |      |  | 29     |
| 30       |      |  | 30     |
| 31       |      |  | 31     |
| 32       |      |  | 32     |
| 33       |      |  | 33     |
| 34       |      |  | 34     |
| 35       |      |  | 35     |
| 36       |      |  | 36     |
| 37       |      |  | 37     |
| 38       |      |  | 38     |
| 39       |      |  | 39     |
| 40       |      |  | 40     |
| 41       |      |  | 41     |
| 42       |      |  | 42     |
| 43       |      |  | 43     |
| 44       |      |  | 44     |
| 45       |      |  | 45     |
| 46       |      |  | 46     |
| 41       |      |  | 47     |
| 48<br>40 |      |  | 48     |
| 49<br>50 |      |  | 49     |
| 50       |      |  | 50     |
| 51       |      |  | 51     |

|     | Hemodynamic Response Function: A   |
|-----|--|
|     | Frequency Domain Approach  |
|     | Ping Bai <sup>1</sup> , Young Truong <sup>2</sup> and Xuemei Huang <sup>3</sup>  |
|     | University of North Carolina at Chapel Hill  |
|     | <b>Abstract:</b> Hemodynamic response function (HRF) has played an important<br>role in many recent functional magnetic resonance imaging (fMRI) based brain<br>studies where the main focus is to investigate the relationship between stimuli<br>and the neural activity. Standard statistical analysis of fMRI data usually<br>calls for a "canonical" model of HRF, but it is uncertain how well this fits<br>the actual HRF. The objective of this paper is to exploit the experimental<br>designs by modeling the stimulus sequences using stochastic point processes.<br>The identification of the stimulus-response relationship will be conducted in<br>the frequency domain, which will be facilitated by fast Fourier transforms<br>(FFT). The usefulness of this approach will be illustrated using both simulated<br>and real human brain data. Under regularity conditions, it is shown that the<br>estimated HRF possesses an asymptotic normal distribution. |
| Co  | ontents  |
| 1   | Introduction   |
| 2   | A Frequency Domain Method for Estimating HRF   |
| 3   | A Brief Survey of HRF Modeling   |
| 1   | Simulated Numerical Results for HRF estimation   |
| 5   | A Real Data Analysis   |
|     | 5.1 Experiment Paradigm and Data Description   |
|     | 5.2 Analysis and Results   |
|     | 5.3 Implications   |
|     | 5.4 Discussions  |
| 3   | Concluding Remarks   |
| 7   | Sampling Properties of the Estimates   |
|     | 7.1 Point Process  |
|     | 7.2 Stationary Time Series   |
|     | (.3 Cumulants and Spectra $\dots \dots \dots$  |
|     | (.4 Fast Fourier Transforms       207         75       C       1         207       207   |
|     | (.5  Complex Normal  |
|     | (.6 Asymptotics for Periodograms   |
|     | <sup>1</sup> Department of Statistics and Operational Research, University of North Carolina, Chape  |
| Hil | l, NC 27599-3260, email: pingbai@gmail.com   |
|     | <sup>4</sup> Department of Biostatistics, University of North Carolina, Chapel Hill, NC 27599-3260, email  |
| rι  | longwolos, unc. eau<br>3Denantment of Neurolem, University of North Coroling, Charol Hill, NC 27500 2260, and  |

1

191

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

Bai, Truong and Huang

7.7Window Estimates — The Smoothed Periodograms21017.8Estimating the Transfer Function21227.9Estimating the Hemodynamic Response Function2133Acknowledgements2144References2145

### 1. Introduction

Consider a statistical problem in which data are acquired by applying stimuli at times  $\tau_1 < \tau_2 < \cdots$  and simultaneously a varying response Y(t) is recorded. Suppose it is desired to conduct the associated statistical inference based on the model:

(1.1) 
$$Y(t) = \sum_{j} h(t - \tau_j) + \epsilon(t),$$

where  $h(\cdot)$  is an unknown function and  $\epsilon(t)$  is a stationary, zero mean, noise series with power spectrum given by  $s_{\epsilon\epsilon}(\cdot)$ . It is assumed that the function h(t) = 0 for t < 0 and will have finite duration.

This type of problems has played an important role in the fields of psychol-ogy, neurobiology, neurology, radiology, biomedical engineering, and many others, where data acquisition is carried out in functional magnetic resonance imaging (fMRI) experiments. As a noninvasive technique, fMRI allows us to study dynamic physiological processes at a time scale of seconds. The basis of fMRI is the Blood Oxygenation Level Dependent (BOLD) effect [40]. Due to differential magnetic sus-ceptibility of oxygenated (oxygen-rich) hemoglobin and deoxygenated hemoglobin, the BOLD effect reflects the changes in hemodynamics which in turn yields greater MRI intensity when brain activity increases (see [30]). It is this hemodynamic re-sponse to the underlying neuronal activity that makes the fMRI signal in brain areas of activation a blurred and delayed version of the stimuli. Figure 1 shows the recorded BOLD signals (solid line) triggered by a single event (dashed line on the left panel) and a sequence of consecutive of stimuli (dashed line on the right panel) respectively. Both of them show the blur-and-delay effect caused by the hemodynamic response.

In practice, the BOLD effect is modeled through the convolution of the stimulus sequence  $X(\cdot)$  and a hemodynamic response function (HRF)  $h(\cdot)$  given by

(1.2) 
$$BOLD(t) = h \otimes X(t) = \int h(t-u)X(u)du,$$

whose discrete time version is represented by model (1.1). Typically, an fMRI dataset consists of a 3D grid of voxels, each containing a time series of measurements that reflect brain activity. For each of roughly 200,000 voxels that lie inside the brain images, we wish to carry out the estimation of the HRF which will subsequently be applied to infer those voxels that were activated under certain experimental conditions.

The rest of this paper is organized as follows. Our estimate is described in Sec-tion 2, which is based on frequency domain methods applied to the point processes and ordinary time series. A brief survey of HRF modeling is provided in Section 3. Section 4 illustrates the performance of our proposed method through a simulated data analysis. For the purpose of comparison, a real fMRI data set is analyzed using the proposed method and a popular fMRI tool in Section 5. Discussions and concluding remarks are given in Section 6. Proofs are given in the last section of the paper. 

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009



FIG 1. Left Panel: The recorded BOLD signal (solid line) triggered by a single event (dashed line). Right Panel: The recorded BOLD signal (solid line) triggered by a typical block-design sequence (dashed line).

# 2. A Frequency Domain Method for Estimating HRF

Model (1.1) has the structure of linear time invariant system carrying the stimuli X(t) onto an response time series Y(t). These models are generally studied by the frequency domain methods based on cross-spectral analysis (see [10]).

Define the system transfer function by

$$H(f) = \sum_{u} h(u) \exp(-iuf), \qquad f \in \mathbb{R}.$$

Define the finite Fourier transform of Y(t) by

$$\varphi_0^T(f) \equiv \varphi_Y^T(f) = \sum_{t=0}^{T-1} \exp(-ift)Y(t)$$

with a similar definition for  $\varphi_{\epsilon}^{T}(f), f \in \mathbb{R}$ . Also, define

$$\varphi_1^T(f) \equiv \varphi_X^T(f) = \sum_{t=0}^{T-1} \exp(-ift) X(t) = \sum_j \exp(-if\tau_j), \qquad f \in \mathbb{R},$$

the last sum is over the available stimuli before time T-1. It follows from (1.1) that

(2.1) 
$$\varphi_0^T(f) = H(f)\varphi_1^T(f) + \varphi_{\epsilon}^T(f), \qquad f \in \mathbb{R}.$$

Now let  $m_f$  denote the integer  $m \in \{0, 1, \dots, T-1\}$  such that  $2\pi m/T$  is closest to the (angular) frequency  $f \in (0, \pi/2)$ . Let K denote a positive integer. Then, for smooth  $H(\cdot)$ ,

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

З
Bai, Truong and Huang

hold for the 2K + 1 nearest frequencies around f of the form  $2\pi(m_f + k)/T$ . Thus a reasonable estimate of H(f) can be obtained by regressing  $\varphi_0^T(2\pi(m_f + k)/T)$  on  $\varphi_1^T(2\pi(m_f + k)/T)$  for  $k = 0, \pm 1, \ldots, \pm K$ , which is given by

(2.3) 
$$\hat{H}(f) = \hat{s}_{01}(f)/\hat{s}_{11}(f), \quad f \in \mathbb{R}$$

where

(2.4) 
$$\hat{s}_{jj'}(f) = (2K+1)^{-1} \sum_{k=-K}^{K} \tilde{s}_{jj'} \left(\frac{2\pi}{T} (m_f + k)\right),$$

(2.5) 
$$\tilde{s}_{jj'}(f) = (2\pi T)^{-1} \varphi_j^T(f) \overline{\varphi_{j'}^T(f)}, \quad f \in \mathbb{R}, \quad j, j' \in \{0, 1\}.$$

Here  $\overline{a}$  is the conjugation of  $a \in \mathbb{C}$ . This is similar to the linear regression setting and therefore the residual sum of squares (RSS) is given by

(2.6) 
$$\hat{s}_{\epsilon\epsilon}(f) = \frac{2K+1}{2K+1-1} \left( \hat{s}_{00}(f) - \hat{s}_{01}(f) \hat{s}_{11}^{-1}(f) \hat{s}_{10}(f) \right).$$

Note that

$$\hat{s}_{\epsilon\epsilon}(f) \propto \hat{s}_{00}(f) \Big( 1 - \frac{|\hat{s}_{01}(f)|^2}{\hat{s}_{11}(f)\hat{s}_{00}(f)} \Big) = \hat{s}_{00}(f)(1 - R_{01}^2(f)),$$

where

$$|\hat{R}_{01}(f)|^2 = \frac{|\hat{s}_{01}(f)|^2}{\hat{s}_{00}(f)\hat{s}_{11}(f)}, \qquad f \in \mathbb{R},$$
22
23
24

is the squared coherence, which lies between 0 and 1, the closer it is to 1 the stronger is the linear relationship between the two series.

Let  $s_{\epsilon\epsilon}(\cdot)$  denote the power spectrum of the noise series. It can be shown that (see Section 7) the estimate  $\hat{H}(f)$  is asymptotically complex normal with mean H(f)and variance  $s_{\epsilon\epsilon}(f)/K\hat{s}_{11}(f)$ . And  $\hat{H}(f_1)$ ,  $\hat{H}(f_2)$ , ...,  $\hat{H}(f_M)$  are asymptotically independent normal for distinct  $f_1, f_2, \ldots, f_M$  [7].

In practice, we use a smoother estimate known as window estimate by observing that (2.4) can be written more generally as

(2.7) 
$$\hat{s}_{jj'}(f) = \sum_{k \neq 0} b^{-1} W\left( b^{-1} \left( f - \frac{2\pi k}{T} \right) \right) \tilde{s}_{jj'}\left( \frac{2\pi k}{T} \right), \qquad f \in \mathbb{R},$$

where  $W(\cdot)$  is a non-negative function called the weight or window function, and  $b \equiv b_T \searrow 0$  is the smoothing parameter. It has been shown that (2.7) has better sampling properties than (2.4) as an estimate of the cross-spectrum of the bivariate time series. See [10] and Section 7. From now on, our estimate of H(f) will be based on (2.3) and (2.7).

We remark that when j = j', then (2.5) becomes

(2.8) 
$$\tilde{s}_{jj}(f) = (2\pi T)^{-1} \varphi_j^T(f) \overline{\varphi_j^T(f)} = (2\pi T)^{-1} |\varphi_j^T(f)|^2, \quad f \in \mathbb{R},$$

which is the *periodogram* of the series Y(t) when j = 0, or of the series X(t) when j = 1. The periodogram is an important statistic in spectral time series analysis. Under certain conditions,  $\hat{R}_{01}(f)$  is asymptotically normal with mean  $R_{01}(f)$ 

and variance proportional to constant  $(1 - R_{01}^2(f))/(Tb)$ . Moreover, if  $R_{01}(f) = 0$ , then

50  
51 (2.9) 
$$F(f) = \frac{c|\hat{R}_{01}(f)|^2}{1 - |\hat{R}_{01}(f)|^2} \sim F_{2,2c},$$
51 51

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

where  $c = (bT/\gamma) - 1$  and  $\gamma = \int \lambda^2$  with  $\lambda$  being the lag-window generator [39]. This result can be used to test for a response to the stimulus by computing a test statistic for significant activation  $F(f_a)$  at the fundamental frequency of activation  $f_a$ . Under the null hypothesis of no activation, the *F*-statistic at the fundamental frequency of activation,  $F(f_a)$ , has a *F* distribution with 2 and 2*c* degrees of freedom. Large values of  $F(f_a)$  indicate a large effect at the fundamental frequency.

The estimate of the impulse response function  $h(\cdot)$  is then given by

$$\hat{h}(u) = \frac{1}{Q} \sum_{q=0}^{Q} \hat{H}\left(\frac{2\pi q}{Q}\right) \exp\left(i\frac{2\pi u q}{Q}\right),$$

where  $Q \equiv Q_T$  denote a sequence of positive integers tending to  $\infty$  with T. Under certain conditions,  $(\hat{h}(u_1), \ldots, \hat{h}(u_J))$  is asymptotically normal with mean  $(h(u_1), \ldots, h(u_J))$  and covariance matrix

$$\frac{2\pi}{bT}\int W(\lambda)^2 d\lambda \cdot \frac{1}{Q^2}\int \exp\left(i(u_j-u_k)\lambda\right)\frac{s_{\epsilon\epsilon}(\lambda)}{s_{11}(\lambda)}d\lambda, \qquad j,k=1,2,\ldots J.$$

See Section 7 for more details.

### 3. A Brief Survey of HRF Modeling

The basis of model (1.1) is the linearity of BOLD fMRI responses when multiple stimuli are presented in succession. This was first studied by Boynton, et al. [6]. The linearity arises from the fact that a stimulus induces the neural activity in a specific region of the brain. This then brings blood flow changes (hemodynamics) in that region, while BOLD fMRI responses are measured from these blood flow changes. In addition to giving this clear picture of how BOLD fMRI works, the linear transform model is important in two respects. Firstly, the assumption of linearity of the fMRI response and neural activity makes it possible to determine changes in neural activity by the amplitude changes in hemodynamic response. Secondly, this linear transform model also shows that when multiple stimuli are presented in succession, the hemodynamic response would be the summation of the individual responses generated by the single stimulus respectively.

Modeling the relationship between the fMRI response and stimuli is a key step towards detecting fMRI activity. Standard statistical analysis is carried out based on the following model:

(3.1) 
$$Y(t) = \beta \sum_{j} h(t - \tau_j) + \epsilon(t),$$

where the HRF  $h(\cdot)$  is pre-specified and  $\beta$  is a voxel specific parameter, to be utilized for testing fMRI activity [15, 21, 23, 24, 31, 33, 43]. The assumptions made about the shape of the HRF vary among different methods. Some of them are very stringent, while others are relatively more flexible. Typically, a "canonical" HRF is employed to process fMRI data. Some studies have reported variation in the shape of HRF across subjects [2, 4], and within the same subject across regions [11, 32, 38].

Detecting fMRI activity has also been evolved from using block-designs (where the stimulus times  $\tau_i$  occur consecutively to form a block) to event-related fMRI З

(ER-fMRI) designs [41]. In the latter case, stimuli (or events) are applied for short bursts in a stochastic manner. The recorded BOLD fMRI signals measure transient changes in brain activity associated with discrete events. This feature makes ER-fMRI a useful tool to estimate the change in the MR signal triggered by neuronal activity.

As an early study of ER-fMRI, Dale and Buckner [16] correlated the selective averaging (or time averaging) data and the fixed HRF induced data in a voxelwise manner. Friston et al. [24] employed a Poisson function with a fixed parameter for the HRF. In the general linear model (GLM) framework, Friston et al. [22] esti-mated the HRF by two given temporal basis functions. To enhance its flexibility, this idea was extended by Josephs et al. [31] to include more basis functions. These are very important contributions since the basis sets allow one to estimate an HRF of arbitrary shape for different events at different voxels of the brain, and at the same time the inferences can be easily made. Many studies on modeling HRF have since focused on the refinement and improvement of the basis sets idea. For example, Woolrich et al. [44] introduced a technique by applying some constraints to avoid nonsensical HRF, which is a big problem when using simple basis functions. More recently, Lindquist and Wager [34] proposed another method, using three superimposed inverse logistic functions, to model the HRF. This paper also described some most popular HRF modeling techniques, such as smooth finite impulse response (FIR) filter [29], canonical HRF with time and dispersion derivatives [14] and the canonical SPM HRF [25]. A flexible method based on splines has been considered by Zhang et al. [47].

From a Bayesian perspective, Genovese [26] and Gössl et al. [28] proposed to model the HRF by a number of parameters and prior distributions are given to each parameter. See also Woolrich et al. [44] and Lindquist and Wager [34]. Inferences of the parameters were then made at each voxel using Markov Chain Monte Carlo (MCMC) technique. The disadvantage of these methods is the slow performance of general MCMC techniques for the inferences.

The above methods are referred to as the time-domain methods. We now consider the frequency-domain approach. Assuming a periodic stimulus design, fMRI time series analysis can be greatly simplified in the frequency domain, which is more natural as the problem of modeling the relationship between the response and the stimuli is reduced to a few parameters related to the stimulus frequency information. One of the first frequency domain approaches is given by Lang and Zeger [33], who used model (3.1) along with a two-parameter gamma function to model the HRF. The two parameters vary at different voxels and hence the estimated HRF varies from voxel to voxel. It was reported that this approach has an identifiability problem of the parameters. The issue was addressed in Marchini and Riplev [37] using a fixed HRF approach. 

We remark that a common theme among the time-domain methods for testing activity is the two-step procedure: the extra step is required for modeling paramet-rically the temporal noise series. This will affect the power of the test. While the above frequency approaches avoided the noise modeling part, they lack the ability to address: (1) the varying HRF issue, and (2) different types of stimulus designs. Moreover, the linear transformed model of fMRI response has not been tested and some studies [16, 30] reported the presence of non-linearity.

In Section 2, we described a regression approach based on model (1.1) for ad-dressing these problems. The application of point processes to model the stimuli is novel in the HRF modeling literature. The procedure is also greatly simplified. Namely, it enables us to estimate the HRF directly and simultaneously test the lin-

З

earity assumption in a single step. The idea of using point processes can be traced back to the work of Brillinger [8, 9] who applied this to identification problems emerged in neurophysiological neural spike train analysis.

### 4. Simulated Numerical Results for HRF estimation

We illustrate the usefulness of our procedure in a simulated study. Here we use one of the HRF's from the literature [16, 27] to generate the response. The stimulus sequence contains a brief trial (1 second) intermixed with events "on" or "off" (Figure 2a). In this experiment, each trial lasts for 18 seconds and there are sixteen runs. So the "average" (because of the random "on" or "off") frequency of the event is 18/288 = 0.0625. The estimated power spectrum (Figure 2b) and the frequency of the event are precisely estimated. The second peak corresponds to the frequency of every other event. In the first experiment (Figure 3a-b), the response is given by  $Y(t) = a \int h(t-u)X(u)du + \epsilon(t)$  with a = 0.5 and the noise  $\epsilon$  is generated from an AR(1) with coefficient 0.7:  $\epsilon(t) = 0.7\epsilon(t-1) + z(t), z(t) \sim N(0, 3^2)$ . In the second experiment (Figure 3c-d), the noise is generated from an ARMA(2,2):  $\epsilon(t) - 0.8897\epsilon(t-1) + 0.4858\epsilon(t-2) = z(t) - 0.2279z(t-1) + 0.2488z(t-2),$  $z(t) \sim N(0, .3^2)$ . The AR(1) time series model was chosen to represent the default settings in Statistical Paramtric Mapping (SPM) [21] and FMRIB Software Library (FSL) [43], while the ARMA case was mainly for testing ten strengths of our method under other types of correlated structures. The coefficients were selected to illustrate the performance of the procedure under moderate serially correlated noise. Large coherency at the stimulus frequency of 0.0625 indicates that the activation is strong, and there is some linearity in the response and the stimulus series. This is also confirmed by the highly significant F-statistics (Figures 4a-d). The significant level is a dashed line that sits near the bottom of the graph. The variability of the proposed estimated is illustrated in Figures 5a-c using various noise level with SD=0.3, 0.5. We remark that the number of runs (=16) used in these simulations is based on recently published articles.

We further apply the procedure to examine the main application of fMRI to detect regions of activation. These are illustrated in Figures 6, 7. In these experiments, the responses are generated from

$$Y(t) = a \int h(t-u)X(u) \, du + \epsilon(t),$$

with varying a to show contrast of the regions so that the sub-region has a higher value of a. The noise component is ARMA(2,2), the same as in the previous experiment with SD=0.3. The regions of activation are clearly captured (Figure 6) when the contrast ratio is high. The effect of the contrast ratio on the detection of region of activation is depicted in Figures 7. It is evident that the level of detection depends on the contrast ratio.

### 5. A Real Data Analysis

# 5.1. Experiment Paradigm and Data Description

In this study, an fMRI data set was obtained from one human subject performing a predefined event sequence as visually instructed. The stimulus sequence includes two different events: right-hand and left-hand finger tapping. Each finger tapping 

З



FIG 2. Plot of the stimulus series (16 trials) with on-off pattern stimuli (every 18 seconds). The whole duration lasts about  $16 \times 18 = 288$  seconds. (b) The estimated power spectrum provides frequency of the occurred events. The frequency associated with the first peak is about 18/288 = 0.0625. The second peak gives frequency of every other event, etc.

movement lasted around 1 second. The order of the sequence was predefined in a random way. To avoid the overlapping of consecutive events, the time interval between two successive events was randomly selected from Uniform[18, 22]. A typical sequence of stimuli is  $\{R, L, R, R, L, L\}$ .

During the experiment, 47 MR scans were acquired on a modified 3T Siemens MAGNETOM Vision system. Each acquisition consisted of 49 contiguous slices. Each slice contained  $64 \times 64$  voxels. Hence there were  $64 \times 64 \times 49$  voxels from each scan. The size of each voxel is  $3\text{mm} \times 3\text{mm} \times 3\text{mm}$ . Each acquisition took 2.9388 seconds, with the scan to scan repetition time (TR) set to be 3 seconds.

# 5.2. Analysis and Results

The data set was preprocessed using SPM5 [21]. The preprocessing included realignment, slice timing correction, coregistration and spatial smoothing. We then analyzed the processed data set using both our proposed method and SPM5. When using SPM5, we used a canonical HRF with time and dispersion derivatives to model the hemodynamic response [25] and its functional form is shown in Figure 5. A *t*-statistic map was generated to show the activations triggered by the stimuli and part of them is shown on the first row of Figure 8.

When using the proposed method to detect which regions of the brain were activated by the finger tapping movements, we generated a spatial color map of the *p*-value for each voxel. The *p*-values were calculated based on the test defined by (2.9). Thus the activation regions are identified by the F statistics that are significant. The *p*-map generated this way is shown on the second row in Figure 8. The four image slices represent the spatial maps of the right-hand activation. The red areas illustrate activated brain regions. Brighter color indicates higher intensity. Our *p*-maps demonstrate the classic brain activation patterns during hand movement as described above. However, the *t*-maps of the same four slices generated using SPM5 do not show any activation, as seen from the first row of Figure 8. 



FIG 3. The fMRI responses. In these experiments, the responses are generated from  $Y(t) = (0.5) \int h(t-u)X(u) \, du + \epsilon(t)$  with  $\epsilon(t) = 0.7\epsilon(t-1) + z(t), z \sim N(0, .3^2)$  in (a) and (b);  $\epsilon(t) - 0.8897\epsilon(t-1) + 0.4858\epsilon(t-2) = z(t) - 0.2279z(t-1) + 0.2488z(t-2), z \sim N(0, .3^2)$  in (c) and (d). The stimuli are the same as in Fig. 2.

Next we plot the estimated HRFs at voxels which are shown to be activated according to Figure 8. Figure 9 displays the estimated HRFs for voxels (with F > 20) selected from primary motor cortex (PMC). Figure 10 displays the estimated HRFs in cerebellum, and Figure 11 shows the estimated HRFs in the supplementary motor area (SMA). These figures were obtained by first computing the *F*-statistics (2.9) followed with the selection of those voxels with the *F* values greater than 20. This threshold was chosen to adjust for the multiple comparison effect and was carried out by computing the *F* statistics over those voxels that are known to be not activated by the stimuli. For example, we used the WFU PickAtlas software [36] to generate region of interest (ROI) mask in the cerebro spinal fluid area of the brain. Then computed the *F* statistics over this area, followed with a density estimate (eg, kernel method, or simply histogram) to select the thresholding value. There are about 20,000 voxels in the cerebro spinal fluid area which can be used to calibrate the null distribution for detecting fMRI activity.

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009



FIG 4. (a) and (c): The estimated coherency function with pointwise confidence intervals. Large coherence values at the event frequencies indicate perfect linear time invariant system used in this simulation. (b) and (d): The F-test for coherency with the dashed-line showing the significant level. (a) and (b) have the same conditions as (a) and (b) of Fig. 3; similarly for (c) and (d).



FIG 5. (a) 95% variation bands (VB) are obtained from estimates of the HRF using 100 random samples simulated from the model by trimming off the 2.5% both ends. (b) Same as in (a) except the noise SD has increased to 0.5. This clearly indicates the variance of the proposed estimate depends on the variance of the noise. (c) Same as in (b) by doubling the number of runs. This illustrates the variance of the estimate is inversely proportional to the duration of the experiment. In these experiments, the responses are generated from  $Y(t) = \int h(t-u)X(u)\,du + \epsilon(t)$  with  $\epsilon \sim {\rm ARMA}(2,2).$  The parameters are the same as in Fig. 3.



imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009



FIG 6. Estimated region of activation. In these experiments, the responses are generated from  $Y(t) = a \int h(t-u)X(u) du + \epsilon(t)$  with a = 0.1 outside and a = 1.0 inside the sub-region, respectively. There are two slices with different region sizes and locations. Each slice is  $8 \times 8$ . (a) and (c) are the true regions, (b) and (d) are the estimated regions. The noise  $\epsilon$  is generated from an ARMA(2,2). The stimuli are the same as in Fig. 2.

All of them have the form that agrees with empirical experience. It is well established that the contralateral cerebral hemisphere motor areas such as primary motor cortex (PMC), and ipsilateral cerebellar areas play dominant role in motor functions in normal human subjects [1, 3, 17, 19]. Our new methods validate unequivocally this known motor activation pattern with single finger movement in a single subject, whereas traditional SPM5 failed to do. Adequate imaging analysis techniques to demonstrate the involvement of those structures during motor function is superior important. PMC is the primary brain region directly control human movement, basal ganglia and cerebellar modulate its functions through a number of cortical motor associated areas of the brain (such as SMA). Dysfunctions of these structures have known to cause a variety of movement disorders such as Parkinson's disease and cerebellar ataxia [1]. Our methods might provide "higher resolution" statistical analysis methods for clinical and neuroscientist to define the roles of these structures in disease stages using fMRI.

# 5.3. Implications

(1) We demonstrated that our method handles statistical issues related to eventrelated experiments well. (2) It is nonparametric in the sense that the functional form of HRF is not specified a priori. Hence it is an useful diagnostic tool for other approaches that may be biased because of misspecification of HRF. (3) Variation of HRF in the brain has been under active research [18], and the nonparametric

#### Bai, Truong and Huang



FIG 7. Estimated region of activation with varying contrast ratios. In these experiments, the responses are generated from  $Y(t) = a \int h(t-u)X(u) du + \epsilon(t)$  with a = 0.2, 0.3, 0.4, 0.5 outside the sub-region and a = 1.0 inside the sub-region. The noise  $\epsilon$  is generated from an ARMA(2,2). These results illustrate that accuracy of the estimates depends on the signal-to-noise (or contrast) ratio: The contrast ratio is proportional to 1/a. (a) Here a = 0.2 implies that the signal is weaker than that in Fig. 6, but the contrast is still high and so the estimated region can still be clearly identified. (b) The contrast here is weaker with a = 0.3. (c) Weaker contrast with a = 0.4, and (d) fuzzy region due to the weakest contrast used in this experiment.

approach offers a systematic way to study the variation without requiring HRF to have the same shape over all voxels. (4) The linear relationship specified through the BOLD signal can be examined statistically by carrying out a formal test of hypothesis. This is important in verifying the linearity assumption employed in SPM [20, 23, 24, 45] in the process of constructing the *T*-map. (5) It is relatively easy to interpret the results using our approach as no prior specification of HRF is required (as is done in SPM [21]/FSL [43]/AFNI [15]).

### 5.4. Discussions

There are many ways to detect fMRI activity. The critical problem is to estimate the statistical significance, which depends on the estimation of both the magnitude of the response to the stimulus and the serial dependence of the time series and especially on the assumptions made in that estimation. Nonparametric spectral density estimation is shown to be self-calibrating and accurate when compared to several other time-domain approaches [12, 13], SPM: [20, 23, 24, 45, 46]. In particular, spectral technique to detect periodic and event-related activations has a distribution theory with significance levels down to 1 in 100,000, levels which are needed when a whole brain image is under consideration. The technique is especially resistant to high frequency artifacts that are found in some datasets and

З



FIG 8. The four related slices that contain the areas activated by right-hand finger tapping. The first row consists of the t-maps generated by SPM5 and they do not show any activation. The second row contains the p-maps generated by the proposed method. The first slice indicates the activated areas in cerebellum. The second slice contains basal ganglia. The third slice contains supplementary motor area (SMA) and the fourth slice shows primary motor cortex (PMC).

it was demonstrated that time-domain approaches may be sufficiently susceptible to these effects to give misleading results. Also, these techniques are capable for detecting activations in clumps of a few (even one) voxel in periodic designs, yet produce essentially no false positive detections at any voxels in null datasets [37].

### 6. Concluding Remarks

It is now widely accepted that fMRI modeling requires flexible HRF modeling, with the HRF varying spatially and between subjects. Flexibility in linear modeling has been introduced with the use of basis functions [22]. However, basis functions suffer from a number of limitations. They impose a hard constraint on the allowed HRF shape and often the extent of the constraint is difficult to control and/or interpret. To overcome these problems, we formulated a procedure based on model (1.1) and FFT. The usefulness has been demonstrated empirically.

We remark that time-domain methods such as SPM [21], FSL [43], FIR [34, 35] and local adaptive spline estimate [47] in modeling the HRF are generally very sensi-tive to the linearity assumption and the error structures they employ. Any approach proposed within the time-domain may have difficulty providing resistant estimates. There is also no guarantee that the parametric noise model chosen will be sufficiently flexible to capture the true form of the correlation structure even if artifacts are removed and a model selection procedure is employed [18, 37]. Therefore significant loss in statistical efficiency can occur if these assumptions are invalidated. In contrast, if these assumptions are valid then the use of a frequency approach will result in a comparatively small loss in efficiency [10]. When considering voxel time series from fMRI datasets there can be no guarantees that the correct time domain 

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009



FIG 9. The histogram for the F-stat and the HRF estimates in the primary motor cortex (PMC) area with F > 20. Each finger-tapping task lasted around 1 second. The order of the sequence was predefined in a random way. The time interval between two successive events was randomly selected from Uniform (18,22). Each acquisition took 2.9388 seconds, with the scan to scan repetition time (TR) set to 3 seconds.



FIG 10. The histogram for the F-stat and the HRF estimates in the cerebellum area with F > 15.

approach has been chosen and a frequency approach seems the most prudent in this context. It is often demonstrated that the assumptions of commonly proposed time-domain models are not resistant to high frequency artifacts.

It is generally believed that the direct analysis of nonperiodic designs will not be as simple as that of the periodic designs, since the response due to the stimulus will be spread over a range of frequencies. Marchini and Ripley [37] suggested that this may be addressed by combining their method with the iteratively reweighted least squares [33] in the spectral domain and the basis functions [22]. However, this method will not be easily extended to model the HRF discussed in this paper. By formulating the problem using point processes, the frequency method advocated by [37] can be easily generalized to handle event-related designs. We also observe that our method is applicable to block designs since the stimuli can be put next to each other to form a block. Thus this unified approach significantly improves the estimation of the HRF described in [33, 37].

The flexible frequency approach proposed here acts as an insurance policy against the results being badly affected by artifacts, and is guaranteed to be near-optimal under all realistic operational conditions. It also offers a quick and accurate way to



FIG 11. The histogram for the F-stat and the HRF estimates in the cerebellum area with F > 15.

check the calibration of the procedure. Further investigations will be carried out for an extensive comparative study on these maps. A concern about our procedure is the choice of weight function  $W(\cdot)$  and bandwidth b given in (2.7). The former is less crucial and it can be addressed by choosing one of commonly used weight functions described in Newton [39]. Bandwidth selection appears to be more serious and it would seem to require adaptive methods such as cross-validation. Based on our experience and the fact that the HRF (blood flow) is relatively smooth, the choice of bandwidth therefore plays a less significant role. Nevertheless, we do observe that the spectral properties of the stimuli can be closely examined by the designs of the experimental protocols, which to some extent can help determine the smoothness of the estimate of HRF. This project is currently underway along with the use of splines for estimating the spectra.

### 7. Sampling Properties of the Estimates

The sampling properties of the HRF estimate will rely on the spectral properties of the stimulus X(t), which is a point process. They also depend on the spectral properties of the response Y(t) and the noise series  $\epsilon(t)$ , which are real-valued ordinary stationary time series. Thus this section starts with a brief summary of the spectral properties of stationary time series which will be denoted by X(t). This is followed by a discussion of the cumulants, which are essential for establishing asymptotic distributions of the estimates. Subsequent sections describe sampling properties of various statistics involved in establishing the properties of HRF estimate.

### 7.1. Point Process

Consider a point process X(t) with points occurring at times  $0 \le \tau_1 \le \tau_2 \le \cdots$ with X(t) denoting the number of points in the interval (0, t]. When it exists, the rate of the process at time t is given by

$$p_X(t) = \lim_{v \downarrow 0} \frac{1}{v} \mathsf{E} \big( X(t+v) - X(t) \big).$$

The expected number of points in the small interval (t, t+v] is given by  $p_X(t)v + o(v)$ . Suppose orderliness, that is, the points of its realizations are isolated, multiple

50

51

| 1                    | points do not occur. Then $dX(t) = 0$ or 1 and one must have   | 1                    |
|----------------------|--|----------------------|
| 2<br>3               | $P\{dX(t)=1\}=p_X(t)dt.$   | 2<br>3               |
| 4<br>5<br>6          | The rate function $p_X(t)$ is seen to have an interpretation as a probability.<br>In the second-order case one defines the second-order product density as   | 4<br>5<br>6          |
| 7<br>8               | $p_{XX}(t_1, t_2) = \lim_{v_1, v_2 \downarrow 0} \frac{1}{v_1 v_2} E \big( X(t_1 + v_1) - X(t_1) \big) \big( X(t_2 + v_2) - X(t_2) \big),  t_1 \neq t_2.$  | 7<br>8               |
| 9<br>10<br>11        | In view of the orderliness of the process, $P\{dX(t) = 1 \text{ and } dX(t) = 1\} = P\{dX(t) = 1\}$ , the case $t_1 = t_2$ can be included via   | 9<br>10<br>11        |
| 12                   | $P\{dX(t_1) = 1 \text{ and } dX(t_2) = 1\} = E(dX(t_1)dX(t_2))$  | 12                   |
| 13<br>14             | $= \left( p_{XX}(t_1, t_2) + \delta(t_1 - t_2) p_X(t_2) \right) dt_1 dt_2,$  | 13<br>14             |
| 15<br>16<br>17<br>18 | where $\delta(\cdot)$ is the Dirac delta function: $\delta(\cdot) \ge 0$ and $\int \delta(t)\varphi(t) dt = \varphi(0)$ for infinitely differentiable real-valued function $\varphi$ with compact support. It is useful to recall here that $\delta$ is the (generalized) derivative of the Heaviside function: $H(\cdot) = 1_{(0,\infty)}(\cdot)$ .<br>See [42] | 15<br>16<br>17<br>18 |
| 19                   | The covariance density of the process is defined by  | 19                   |
| 20<br>21             | $q_{XX}(t_1, t_2) = p_{XX}(t_1, t_2) - p_X(t_1)p_X(t_2),$  | 20<br>21             |
| 22                   | with the interpretation  | 22                   |
| 23<br>24             | $cov \{ dX(t_1) \ dX(t_2) \} = (a_{XX}(t_1, t_2) + \delta(t_2 - t_1)a_{X}(t_2)) dt_1 dt_2$   | 23<br>24             |
| 25<br>26<br>27       | The conditional intensity of the process is defined by $p_{XX}(t)/p_X(t)$ with the inter-<br>pretation   | 25<br>26<br>27       |
| 28                   | $P\{dX(t_2) = 1 \mid dX(t_1) = 1\} = \left(p_{XX}(t_1, t_2)/p_X(t_1)\right) dt_2.$   | 28                   |
| 29<br>30<br>31       | A point process is said to be stationary when its probability properties are un-<br>affected by simple shifts of time. In this case one has  | 29<br>30<br>31       |
| 32<br>33             | $P\{dX(t)=1\} = p_X dt,$   | 32<br>33             |
| 34                   | $P\{dX(t_1) = 1 \text{ and } dX(t_2) = 1\} = (p_{XX}(t_2 - t_1) + \delta(t_2 - t_1)p_X) dt_1 dt_2$   | 34                   |
| 35<br>36             | $\cos\{dX(t_1), dX(t_2)\} = (a_{XX}(t_2 - t_1) + \delta(t_2 - t_1)p_X) dt_1 dt_2.$   | 35<br>36             |
| 37<br>38             | By analogy with what is done in the ordinary time series case one may define the power spectrum at frequency $f$ by  | 37<br>38             |
| 39<br>40<br>41       | $s_{XX}(f) = \frac{1}{2\pi} \int e^{-ifu} \left( \cos\{dX(t+u), dX(t)\}/dt \right) du.$  | 39<br>40<br>41       |
| 42<br>43<br>44       | For multivariate process $X(t) = \{X_1(t), \ldots, X_m(t)\}$ , it may be convenient to consider  | 42<br>43<br>44       |
| 45                   | $P\{dX_j(t) = 1\} = C_j dt, \qquad j = 1, \dots, m,$   | 45                   |
| 46                   | and $\cos\left[dY_{1}(t+u), dY_{2}(t)\right] = C_{1}(du)dt$ $i, h = 1$ m   | 46                   |
| 47<br>48             | $\operatorname{cov}\{a\boldsymbol{\Lambda}_{j}(t+u), a\boldsymbol{\Lambda}_{k}(t)\} = \mathbb{C}_{jk}(au)  at, \qquad j, \kappa = 1, \dots, m.$  | 47                   |
| 49                   | I ne power spectrum at frequency $f$ is defined by   | 49                   |

 $s_{jk}(f) = \frac{1}{2\pi} \int e^{-ifu} C_{jk}(du), \qquad j,k = 1,\dots,m.$ 

7.2. Stationary Time Series

Let  $X(t) = (X_1(t), \ldots, X_r(t)), t = 0, 1, 2, \ldots$ , denote a vector-valued stationary time series. Set

$$C_{jk}(u) = \operatorname{cov}\{X_j(t+u), X_k(t)\}, \qquad j, k = 1, \dots, r.$$

The power spectrum at frequency f is defined by

$$s_{jk}(f) = \frac{1}{2\pi} \sum_{u} e^{-ifu} C_{jk}(u), \qquad f \in \mathbb{R}, \quad j, k = 1, \dots, r.$$

# 7.3. Cumulants and Spectra

**Definition 1.** Let  $X_1, X_2, \ldots, X_r$  denote random variables with finite rth moment. The rth order joint cumulant of  $X_1, X_2, \ldots, X_r$  is defined by

$$\operatorname{cum}(X_1, X_2, \dots, X_r) = \sum (-1)^{p-1} (p-1)! \left(\prod_{j \in \nu_1} X_j\right) \dots \left(\prod_{j \in \nu_p} X_j\right),$$

where the summation extends over all partitions  $\nu_1, \ldots, \nu_p$ ,  $p = 1, \ldots, r$  of  $\{1, 2, \ldots, r\}$ .

### Remarks

- 1. When  $X_1 = X_2 = \cdots = X_r$ , the definition gives the cumulant of order r of a univariate random variable.
- 2.  $\operatorname{cum}(X_1, X_2, \ldots, X_r)$  is also given by the coefficient of  $(i)^r t_1 \ldots t_r$  in the Taylor series expansion of  $\log(E \exp i \sum_{j=1}^r X_j t_j)$ .

Given r time series  $X_1(t), X_2(t), \ldots, X_r(t)$  with each having finite rth moment, we define

$$C_{1,\ldots,r}(t_1,t_2,\ldots,t_r) = \operatorname{cum}(X_1(t_1),X_2(t_2),\ldots,X_r(t_r)).$$

For stationary time series,

$$C_{1,\dots,r}(t_1,t_2,\dots,t_r) = C_{1,\dots,r}(t_1-t_r,t_2-t_r,\dots,t_{r-1}-t_r,0),$$

which is a function of r-1 variables. In this case, the *r*th order cumulant spectrum,  $s_{1,\ldots,r}(f_1, f_2, \ldots, f_{r-1})$ , is defined by

$$s_{1,\ldots,r}(f_1, f_2, \ldots, f_{r-1}) =$$

$$(2\pi)^{-k+1} \sum_{u_1, u_2, \dots, u_{r-1}} C_{1, \dots, r}(u_1, u_2, \dots, u_{r-1}) \exp\left(-i \sum_{j=1}^{r-1} u_j f_j\right),$$

 $u_{1}, u_{2}, \dots, u_{r-1}$   $f_{1}, f_{2}, \dots, f_{r-1} \in \mathbb{R}, \quad r \ge 2.$ 

For a more detailed discussion of cumulants and their spectra, see [10].

#### imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

.

# 7.4. Fast Fourier Transforms

Let  $a_j(\cdot) : \mathbb{R} \to \mathbb{R}, j = 1, 2$ , denote **tapering functions**. The discrete Fourier transform for the univariate series  $X_j$  is defined by

$$\varphi_j^T(f) \equiv \varphi_{X_j}^T(f) = \sum_t a_j(t/T) X_j(t) \exp(-ift), \qquad f \in \mathbb{R}, \quad j = 1, 2.$$

For vector-valued series  $\mathbf{X}$ , it is given by

$$\boldsymbol{\varphi}^{T}(f) \equiv \boldsymbol{\varphi}_{\mathbf{X}}^{T}(f) = \begin{pmatrix} \varphi_{1}^{T}(f) \\ \varphi_{2}^{T}(f) \end{pmatrix}, \qquad d \in \mathbb{R}.$$
<sup>9</sup>
<sup>10</sup>
<sup>11</sup>
<sup>12</sup>

Set  $a_j^T(t) = a_j(t/T), j = 1, 2$ . For  $j_m \in \{1, 2\}, m = 1, \dots, M$ ,

$$A_{j_1,\dots,j_M}^T(f) = \sum_t \left(\prod_{m=1}^M a_{j_m}^T(t)\right) \exp(-ift), \quad f \in \mathbb{R}.$$

**Condition 1.** The tapering function  $a(\cdot) : \mathbb{R} \to \mathbb{R}$  has a compact support with bounded first derivative. Furthermore,

$$\int a(u) \, du = 1$$
 and  $\int |a(u)| \, du < \infty$ .

**Condition 2.** The covariance function satisfies

$$\sum_{u} C_{jk}(u) < \infty,$$

and

$$\sum_{u_1,\dots,u_{M-1}} C_{j_1\dots j_M}(u_1,\dots,u_{M-1}) < \infty, \qquad j_1,\dots,j_M = 1,2.$$

The second part of the above condition is necessary for establishing the asymptotic properties of the estimates to be considered in this section.

Lemma 1. Suppose Conditions 1 and 2 hold. Then

$$\sup_{f_1,\dots,f_M} \left| \operatorname{cum}(\varphi_{j_1}^T(f_1),\dots,\varphi_{j_M}^T(f_M)) - (2\pi)^{M-1} A_{j_1,\dots,j_M}^T(f_1+\dots+f_M) s_{j_1,\dots,j_M}(f_1,\dots,f_M) \right| = o(T).$$

**Condition 3.** The covariance function satisfies

$$\sum_{u} |u|c_{jk}(u) < \infty,$$

and

$$\sum_{u_1,\dots,u_{M-1}} |u_1\cdots u_{M-1}| C_{j_1\dots j_M}(u_1,\dots,u_{M-1}) < \infty, \qquad j_1,\dots,j_M = 1,2.$$

Lemma 2. Under Conditions 1 and 3,

$$\sup_{f_1,...,f_M} |\operatorname{cum}(\varphi_{j_1}^T(f_1),\ldots,\varphi_{j_M}^T(f_M))$$
48
49
50

$$-(2\pi)^{M-1}A_{j_1,\ldots,j_M}^T(f_1+\cdots+f_M)s_{j_1,\ldots,j_M}(f_1,\ldots,f_M)\Big|=O(1).$$

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

*Proof.* We now prove Lemmas 1 and 2. If follows from  $|a_{i}(t+u)a_{k}(t+v) - a_{i}(t)a_{k}(t)| \leq |a_{i}(t+u)a_{k}(t+v) - a_{i}(t+u)a_{k}(t)|$  $+ |a_{i}(t+u)a_{k}(t) - a_{i}(t)a_{k}(t)|$ and Condition 1 that there is a constant  $K_1$  such that  $\left|\sum_{t} a_{j_{1}}^{T}(t+u_{1})\cdots a_{j_{M-1}}^{T}(t+u_{M-1})a_{j_{M}}^{T}(t)\exp(-ift) - A_{j_{1}\dots j_{M}}^{T}(f)\right|$  $< K_1(|u_1| + \dots + |u_{M-1}|).$ By the cumulant property,  $\operatorname{cum}(\varphi_{i_1}^T(f_1),\ldots,\varphi_{i_M}^T(f_M))$  $=\sum_{i}\cdots\sum_{j=1}a_{j_{1}}^{T}(t_{1})\cdots a_{j_{M}}^{T}(t_{M})\exp\left(-i\sum_{j=1}^{M}f_{m}t_{m}\right)$  $\times C_{j_1,\ldots,j_M}(t_1-t_M,\ldots,t_{M-1}-t_M)$  $=\sum_{u_1=-2(T-1)}^{2(T-1)}\cdots\sum_{u_{M-1}=-2(T-1)}^{2(T-1)}\exp\left(-i\sum_{m=1}^{M-1}f_mt_m\right)C_{j_1,\dots,j_M}(u_1,\dots,u_{M-1})$  $\times \sum_{t} a_{j_1}^T (t+u_1) \cdots a_{j_{M-1}}^T (t+u_{M-1}) a_{j_M}^T (t) \exp\left(-i \sum_{m=1}^M f_m t\right)$  $= \epsilon_T + \sum_{n_1=-2(T-1)}^{2(T-1)} \cdots \sum_{n_M}^{2(T-1)} \exp\left(-i\sum_{m=1}^{M-1} f_m t_m\right)$  $C_{j_1,\ldots,j_M}(u_1,\ldots,u_{M-1})A_{j_1\ldots,j_M}^T(f_1+\cdots+f_M),$ where  $|\epsilon_T| \le K_2 \sum_{u_1=-2(T-1)}^{2(T-1)} \cdots \sum_{u_M} \sum_{i=-2(T-1)}^{2(T-1)} (|u_1| + \dots + |u_{M-1}|) C_{j_1,\dots,j_M}(u_1,\dots,u_{M-1}).$ It now follows from Condition 3,  $|T^{-1}|\epsilon_T| \le K_2$   $\sum_{\alpha(T-1)}^{2(T-1)} \cdots \sum_{\alpha(T-1)}^{2(T-1)} T^{-1}(|u_1| + \dots + |u_{M-1}|)$  $C_{i_1,\ldots,i_M}(u_1,\ldots,u_{M-1}),$  $T^{-1}(|u_1| + \cdots + |u_{M-1}|) \to 0$  and the dominated convergence theorem that  $|\epsilon_T| = o(T).$ (7.1)Lemmas 1 and 2 follow from this and  $s_{i_1,\ldots,i_M}(f_1,\ldots,f_{M-1})$  $= (2\pi)^{M-1} \sum \cdots \sum \exp\left(-i \sum_{i=1}^{M-1} f_m u_m\right) C_{j_1,\dots,j_M}(u_1,\dots,u_{M-1}) + o(1).$ imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

# 7.5. Complex Normal

Let X denote an k-dimensional random vector whose components are complexvalued random variables. If, for some  $\mu \in \mathbb{C}^k$  and  $k \times k$  Hermitian non-negative definite matrix  $\Sigma$  (that is,  $\Sigma = \overline{\Sigma}^{\top}$ ),

$$\begin{pmatrix} \operatorname{Re} X \\ \operatorname{Im} X \end{pmatrix} \sim N_{2k} \left( \begin{pmatrix} \operatorname{Re} \mu \\ \operatorname{Im} \mu \end{pmatrix}, \frac{1}{2} \begin{pmatrix} \operatorname{Re} \Sigma & -\operatorname{Im} \Sigma \\ \operatorname{Im} \Sigma & \operatorname{Re} \Sigma \end{pmatrix} \right),$$

we say X has a complex normal distribution with mean  $\mu$  and covariance matrix  $\Sigma$ , and is abbreviated by  $X \sim N_k^c(\mu, \Sigma)$ .

The FFT is asymptotically normal with mean specified according to the frequency f as described below.

**Theorem 7.1.** Under Conditions 1 and 2,  $\varphi_i^T(f)$  is asymptotically

1.  $N_1^c(0, 2\pi T s_{jj}(f) A_{jj}(0))$  if  $f \neq 0 \mod \pi$ ,

2.  $N_1(Tc_jA_{jj}(0), 2\pi Ts_{jj}(0)A_{jj}(0))$  if  $f = 0, \pm 2\pi, \dots,$ 

3.  $N_1(0, 2\pi T s_{jj}(\pi) A_{jj}(0))$  if  $f = \pm \pi, \dots$ 

Note that  $A_{jj}(0) = \int a_j^2$ . The above result implies that the real and the imaginary part of  $\varphi_j^T(f)$  are approximately independent. Each is approximately normal with mean and variance  $\pi T s_{jj}(f) \int a_j^2$ .

*Proof.* To prove Theorem 7.1, we note that by Condition 1,  $A_{j_1,\ldots,j_M}^T(f) = O(T)$ . Recall that the Gaussian distribution has cumulants of order greater than 2 vanishes. The desired result now follows from Lemmas 1, 2 and that fact that

$$T^{-M/2} \operatorname{cum}(\varphi_{j_1}^T(f_1), \dots, \varphi_{j_M}^T(f_M))$$
  
=  $T^{-M/2} (2\pi)^{M-1} A_{j_1, \dots, j_M}^T (f_1 + \dots + f_M) s_{j_1, \dots, j_M} (f_1, \dots, f_M) + o(T^{1-M/2})$   
 $\to 0 \quad \text{for } M > 2 \text{ as } T \to \infty.$ 

# 7.6. Asymptotics for Periodograms

The distributions of the FFT suggests the following statistic:

$$\tilde{s}_{jj}(f) = |\varphi_j^T(f)|^2 / \left(2\pi \sum_t [a_j(t/T)]^2\right), \qquad f \in \mathbb{R}.$$

This is called **periodogram** and is an estimate of the spectral density function  $s_{jj}$ . For more historical remarks, see [10]. Note that if there is no tapering function, the periodogram is given by

$$\tilde{s}_{jj}(f) = (2\pi T)^{-1} |\varphi_j^T(f)|^2, \qquad f \in \mathbb{R}.$$

Let  $f_m = 2\pi m/T$ ,  $m = 0, \pm 1, \pm 2, \dots, \pm T/2$  denote the Fourier frequencies. The result below describes the asymptotic distribution of the periodograms.

Theorem 7.2. Under Conditions 1–3,  $\tilde{s}_{jj}(f_m)$ , m = 1, ..., M = T/2, are asymptotically independent  $s_{jj}(f_m)\chi_2^2/2$ . Also  $\tilde{s}_{jj}(f)$  is asymptotically  $s_{jj}(f)\chi_1^2$  for  $f = 50 \pm \pi, \pm 3\pi, ...$ , independent of the  $\tilde{s}_{jj}(f_m)$ , m = 1, ..., T/2.

*Proof.* The proof follows from Theorem 7.1 and the definition of the chi-square distribution.

The above result shows that the asymptotic variance of the periodogram is approximately  $s_{jj}(f)^2$ , which is usually positive. Thus the periodogram is not a consistent estimate of the spectral density function. The following section will present a class of consistent estimates obtained by smoothing the periodograms.

# 7.7. Window Estimates — The Smoothed Periodograms

A class of consistent estimates can be obtained by using a running mean or local average of the periodograms. Specifically, set

$$\hat{s}(f_m) = (2K+1)^{-1} \sum_{k=-K}^{K} \tilde{s}_{jj} \left(\frac{2\pi}{T}(m+k)\right).$$
<sup>15</sup>
<sup>16</sup>
<sup>17</sup>
<sup>16</sup>
<sup>17</sup>

It follows from the asymptotic distributional properties of the periodograms (Theorem 7.2) that  $\hat{s}(f_m)$ ,  $m = 1, \ldots, T/2$ , are asymptotically independent with  $\hat{s}(f) \sim s(f)\chi_{4K+2}^2/(4K+2)$  if  $f \neq 0$ , and  $\hat{s}(0) \sim s(0)\chi_{2K}^2/(2K)$ . An important implication of the above result is that consistency can be achieved by letting  $K \to \infty$  and  $K/T \to 0$  as  $T \to \infty$ .

More generally, let  $W(\cdot)$  denote a weight function. Set

(7.2) 
$$\hat{s}_{jj'}(f) = \sum_{k \neq 0} b_T^{-1} W\left(b_T^{-1}\left(f - \frac{2\pi k}{T}\right)\right) \tilde{s}_{jj'}\left(\frac{2\pi k}{T}\right),$$

where

$$\tilde{s}_{jj'}(f) = (2\pi T)^{-1} \varphi_j^T(f) \overline{\varphi_{j'}^T(f)}, \qquad f \in \mathbb{R},$$

and  $b_T$  is referred to as the **bandwidth** or **window width** that will be specified more clearly later. Certain properties of the weight function  $W(\cdot)$  will be required in order to assure that the above estimate is consistent.

**Condition 4.** The weight function  $W(\cdot) : \mathbb{R} \to \mathbb{R}$  is a symmetric probability density function with a compact support  $[-\pi, \pi]$ .

Under this condition, the bias of the window estimate is given by

$$E(\hat{s}_{jj}(f)) = \int W(\lambda) s_{jj}(f - b_T \lambda) \, d\lambda + O(T^{-1}b_T^{-1}).$$

In fact, more properties can be obtained and are stated in the following result.

**Theorem 7.3.** Under Conditions 1–3 and suppose that the spectral density function  $s_{jj}$  does not vanish. Let  $b_T \to 0$  and  $b_T T \to \infty$  as  $T \to \infty$ . Then,  $\hat{s}(f_m)$ ,  $m = 1, \ldots, M$ , are asymptotically normal with mean zero and covariance structure given by

49  
50 (7.3) 
$$\lim_{T \to \infty} b_T T \operatorname{cov}(\hat{s}(f_1), \hat{s}(f_2)) = \begin{cases} 0 & \text{if } f_1 \neq f_2, \\ 2\pi \sigma(f_1)^2 \in W^2 & \text{otherwise} \end{cases}$$

(7.3) 
$$\lim_{T \to \infty} b_T T \operatorname{cov}(\hat{s}(f_1), \hat{s}(f_2)) = \begin{cases} 3 & \text{or } f(f_1) \\ 2\pi s(f)^2 \int W^2 & \text{otherwise.} \end{cases}$$

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

*Proof.* Direct computation shows that

$$\operatorname{cov}(\tilde{s}_{jj}(f_1), \tilde{s}_{jj}(f_2)) = \left( (\sin T(f_1 + f_2)/2)^2 - (\sin T(f_1 - f_2)/2)^2 \right)$$

$$= s_{jj}(f_1) \left\{ \left( \frac{\sin T(f_1 + f_2)/2}{T \sin (f_1 + f_2)/2} \right) + \left( \frac{\sin T(f_1 - f_2)/2}{T \sin (f_1 - f_2)/2} \right) \right\} + O(1/T).$$

Moreover,

$$\operatorname{cov}(\hat{s}_{jj}(f_1), \hat{s}_{jj}(f_2)) = 2\pi T^{-1} \int W^T(f_1 - \lambda) W^T(f_2 - \lambda) s_{jj}(\lambda)^2 d\lambda$$

$$+ 2\pi T^{-1} \int W^T (f_1 - \lambda) W^T (f_2 + \lambda) s_{jj}(\lambda)^2 d\lambda$$

$$+ O(b_T^{-2}T^{-2}) + O(T^{-1}), 13$$

where

$$W^T(\lambda) = b_T^{-1} \sum_{k=-\infty}^{\infty} W(b_T^{-1}(\lambda + 2\pi k)).$$

The indicated covariance structure (7.3) is an easy consequence of these results.

To obtain the asymptotic normality, we need to show that all cumulants of order higher than 2 tend to zero as  $T \to \infty$ . This is carried out by directly computing the cumulants of the window estimates in a manner similar to the proof of Lemma 1.

### 7.8. Estimating the Transfer Function

**Theorem 7.4.** Suppose that

1.  $\epsilon(t), t = 0, 1, \dots$  satisfy Condition 2 and have mean zero,

2. X(t) is uniformly bounded and  $s_{11} \neq 0$ ,

- 3.  $\sum_{u} |u|h(u) < \infty$ ,
  - 4.  $\overline{W}$  in Condition 4 is a uniform kernel.

Let  $b_T \to 0$ ,  $b_T T \to \infty$ ,  $b_T^5 T \to 0$  as  $T \to \infty$ . Then  $\hat{H}(f_1), \ldots, \hat{H}(f_M)$  is complex normal with mean  $(E\hat{H}(f_1), \ldots, E\hat{H}(f_M))$  and covariance matrix whose entries are given by

$$\operatorname{cov}(\hat{H}(f_1), \hat{H}(f_2)) = \eta(f_1 - f_2) \frac{2\pi s_{\epsilon\epsilon}(f_1)}{b_T T s_{11}(f_1)} \int W^2,$$

where  $\eta(0) = 1$  and  $\eta(f) = 0$  for  $f \neq 0$ .

The weight function W is assumed to be uniform on  $[-\pi, \pi]$  in order to simplify the presentation of the above asymptotic properties of the estimate of the transfer function. A more general approach can be found in [10].

*Proof.* We begin the proof of Theorem 7.4 with two lemmas.

Lemma 3. Let  $(\mathbf{V}_n)$  denote a sequence of random vectors converging in distribution to **V**. Then there exists a probability space such that  $\mathbf{V}_n$  converges to **V** almost surely.

Proof. The proof can be found in [5].

**Lemma 4.** Let  $(\mathbf{V}_n)$  denote a sequence of random vectors in  $\mathbb{R}^p$  converging in 50 distribution to  $N_p^c(\mathbf{0}, \mathbf{I}_p)$  and  $(\mathbf{U}_n)$  a sequence of  $p \times p$  unitary matrices. Then 51  $\mathbf{U}_n \mathbf{V}_n$  converges to  $N_p^c(\mathbf{0}, \mathbf{I}_p)$  as  $n \to \infty$ .

-±3 

*Proof.* This follows from Lemma 3.

З

З

Before proceeding to the proof, we remark that the following argument is simplified by assuming the series X to be non-random. The result nevertheless holds for general random X. Let  $\varphi_i^T$ , j = 0, 1, be the Fourier transform of Y and X, respectively. Let  $2\pi k/T$  denote the Fourier frequency that is nearest to  $\lambda$ . Then  $\varphi_0^T (2\pi (k+l)/T)$  $= H(2\pi(k+l)/T)\varphi_1^T(2\pi(k+l)/T) + \varphi_c^T(2\pi(k+l)/T) + O(1)$  $= H(\lambda)\varphi_1^T(2\pi(k+l)/T) + \varphi_{\epsilon}^T(2\pi(k+l)/T) + O(1), \quad l = 0, \pm 1, \dots, \pm m,$ where O(1) is uniformly in l. Now let  $\mathbf{D}_0$  denote the  $1 \times (2m+1)$  matrix given by  $\mathbf{D}_{0} = \frac{1}{\sqrt{2\pi T}} \left[ \varphi_{0}^{T}(2\pi(k-m)/T) \cdots \varphi_{0}^{T}(2\pi k/T) \cdots \varphi_{0}^{T}(2\pi(k+m)/T) \right].$ Define  $\mathbf{D}_1$  and  $\mathbf{D}_{\epsilon}$  similarly. Then  $\mathbf{D}_0 = H(f)\mathbf{D}_1 + \mathbf{D}_{\epsilon} + O(T^{-1/2}).$ Let  $\mathbf{U} \equiv \mathbf{U}^T = [\mathbf{U}_1, \mathbf{U}_0]$  be a  $(2m+1) \times (2m+1)$  unitary matrix whose first column is  $\mathbf{U}_1 = \mathbf{D}_1^H (\mathbf{D}_1 \mathbf{D}_1^H)^{-1/2}$ , where  $\mathbf{D}^H = \overline{\mathbf{D}}^\top$  is the conjugate transpose of **D**. Then  $\mathbf{D}_0 \mathbf{U} = H(f) \mathbf{D}_1 \mathbf{U} + \mathbf{D}_{\epsilon} \mathbf{U} + O(T^{-1/2}).$ The first and the remaining columns of these matrices yield  $[\hat{H}(f) - H(f)]\hat{s}_{11}(f)^{1/2}(2m+1)^{1/2} = \mathbf{D}_{\epsilon}\mathbf{U}_1 + O(T^{-1/2}),$ (7.4) $\mathbf{D}_0 \mathbf{U}_0 = \mathbf{D}_{\boldsymbol{\epsilon}} \mathbf{U}_0 + O(T^{-1/2}).$ (7.5)By the property of the unitary matrix,  $(2m+1)\hat{s}_{00} = \mathbf{D}_0\mathbf{D}_0^H = \mathbf{D}_0\mathbf{U}_1\mathbf{U}_1^H\mathbf{D}_0^H + \mathbf{D}_0\mathbf{U}_0\mathbf{U}_0^H\mathbf{D}_0^H$  $= \mathbf{D}_0 \mathbf{D}_1^H (\mathbf{D}_1 \mathbf{D}_1^H)^{-1} \mathbf{D}_1 \mathbf{D}_0^H + \mathbf{D}_0 \mathbf{U}_0 \mathbf{U}_0^H \mathbf{D}_0^H.$ Thus  $\hat{s}_{\epsilon\epsilon} = \mathbf{D}_0 \mathbf{U}_0 \mathbf{U}_0^H \mathbf{D}_0^H = \mathbf{D}_{\epsilon} \mathbf{U}_0 \mathbf{U}_0^H \mathbf{D}_{\epsilon}^H + O_n (T^{-1/2}).$ (7.6)Now, according to Theorem 7.1,  $\mathbf{D}_{\epsilon} \rightarrow_d N_{2m+1}^c N(\mathbf{0}, s_{\epsilon,\epsilon}(f)\mathbf{I})$  and therefore  $s_{\epsilon}(f)^{-1/2}\mathbf{D}_{\epsilon} \rightarrow_d N_{2m+1}^c N(\mathbf{0}, \mathbf{I}).$  By Lemma 4,  $s_{\epsilon}(f)^{-1/2}\mathbf{D}_{\epsilon}\mathbf{U} \rightarrow_d N_{2m+1}^c N(\mathbf{0}, \mathbf{I}),$ or  $\mathbf{D}_{\epsilon}\mathbf{U} \to_d N_{2m+1}^c N(\mathbf{0}, s_{\epsilon}(f)\mathbf{I})$ . This, together with (7.4) and (7.6) yield the desired result. This completes the proof of the theorem. 7.9. Estimating the Hemodynamic Response Function From  $H(f) = \sum h(u) \exp(-iuf),$ we see that the hemodynamic response function is given by  $h(u) = \frac{1}{2\pi} \int_{0}^{2\pi} H(f) \exp(iuf) df, \qquad u = 0, \pm 1, \pm 2, \dots$ 

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

Bai, Truong and Huang

Let  $\hat{H}(f)$  denote an estimate of H(f) given by the last section, and let  $Q \equiv Q_T$  denote a sequence of positive integers tending to  $\infty$  with T. As an estimate of h(u) by approximating the integral using finite sums, we define

$$\hat{h}(u) = \frac{1}{Q} \sum_{q=0}^{Q} \hat{H}\left(\frac{2\pi q}{Q}\right) \exp\left(i\frac{2\pi u q}{Q}\right), \qquad u = 0, \pm 1, \pm 2, \dots$$

# **Theorem 7.5.** Suppose that

1.  $\epsilon(t), t = 0, 1, \dots$  satisfy Condition 2 and have mean zero,

2. X(t) is uniformly bounded and  $s_{11} \neq 0$ ,

3.  $\sum_{u} |u|h(u) < \infty$ ,

4.  $\overline{W}$  in Condition 4 is a uniform kernel.

Let  $Qb \to 0$  as  $T \to \infty$ . Then

$$E\hat{h}(u) = h(u) + \sum_{q \neq 0} h(u + qQ) + O(b) + O(T^{-1/2}).$$

In particular,  $\hat{h}(u)$  is asymptotically unbiased. Furthermore,  $\hat{h}(u_1), \ldots, \hat{h}(u_M)$  are asymptotically normal with mean  $h(u_1), \ldots, h(u_M)$  and covariance structure

$$\operatorname{cov}(\hat{h}(u), \hat{h}(v)) = \frac{2\pi}{QbT} \Lambda^{T}(u, v) \int W^{2} + O(T^{-1}), \quad u, v = 0, \pm 1, \pm 2, \dots,$$

where

$$\Lambda^{T}(u,v) = \frac{1}{Q} \sum_{q=0}^{Q} \exp\left(i\frac{2\pi(u-v)q}{Q}\right) s_{\epsilon\epsilon}(2\pi q/Q) / s_{11}(2\pi q/Q).$$

*Proof.* The proof of Theorem 7.5 is tedious and very computational. We outline the argument here. The proof starts by assuming the X series to be non-random. The asymptotic normality then follows from the computation of the joint cumulants, it is shown that cumulants of order greater than 2 of  $\hat{h}(u_1), \ldots, \hat{h}(u_M)$  tend to zero as  $T \to \infty$ . The desired result for the random X follows by invoking a standard technique in nonparametric regression for handling ratio of two random variates.

### Acknowledgements

We like to thank Mechelle Lewis, Andrew Smith, Suman Sen and Roxanne Poole for their comments and helps to make the data available to us. We are grateful to three referees for their critical and constructive comments that improve the readability of the paper significantly. Finally, we thank Allen Song for the physics behind MRI.

# References

| 45 | [1] | AFIFI, A. K. and BERGMAN, R. A. (1998). Functional Neuroanatomy: Text and Altas. McGraw-Hill,       | 45 |
|----|-----|---|----|
| 46 |     | USA.  | 46 |
| 17 | [2] | Aguirre, G. K., Zarahn, E. and D'Esposito, M. (1998). The Variability of Human, BOLD Hemo-          | 17 |
| 47 |     | dynamic Responses. NeuroImage, 8, 360–369.  | 41 |
| 48 | [3] | ALEXANDER, G. E., DELONG, M. R. and STRICK, P. L. (1986). Parallel organization of functionally     | 48 |
| 49 |     | segregated circuits linking basal ganglia and cortex. Annu. Rev. Neurosci., 9, 357–381.             | 49 |
|    | [4] | BÉNAR, C. G., GROSS, D. W., WANG, Y., PETRE, V., PIKE, B., DUBEAU, F. and GOTMAN, J. (2002).        |    |
| 50 |     | The BOLD Response to Interictal Epileptiform Discharges. <i>NeuroImage</i> , <b>17</b> , 1182–1192. | 50 |
| 51 | [5] | BILLINGSLEY, P. (1995). Probability and Measures, 3rd ed., Wiley, New York, USA.                    | 51 |

imsart-coll ver. 2008/08/29 file: Truong.tex date: March 25, 2009

З

| 1<br>2 | [6]           | BOYNTON, G. M., ENGEL, S. A., GLOVER, G. H. and HEEGER, D. J. (1996). Linear Systems Analysis of<br>Functional Magnetic Resonance Imaging in Human V1. <i>The Journal of Neuroscience</i> , <b>16</b> , 4207– | 1<br>2 |
|--------|---------------|---|--------|
| 3      | [7]           | BRILLINGER, D. R. (1974). Cross-Spectral Analysis of Processes with Stationary Invrements Including   | 3      |
| 4      | [8]           | the Stationary $G/G/\infty$ Queue. $AOP$ , 2, 815–827.  | 4      |
| 5      | [9]           | BRILLINGER, D. R. (1978). A note on the estimation of evoked response. <i>Biological Cubernetics</i> .  | 5      |
| 6      | r.1           | <b>31:3</b> , 141–144.  | 6      |
| 7      | [10]          | BRILLINGER, D. R. (1981). Time Series: Data Analysis and Theory. Holden-Day, San Francisco,   | 7      |
| 8      | [11]          | CA.<br>BUGUERD D. L. AND KOUTGERLAL W. AND COLLECTED D. L. AND DALE A. M. AND BOTTED M. and   | 8      |
| 9      | [11]          | ROSEN, B. R. (1998). Functional-Anatomic Study of Episodic Retrieval: II. Selective Averaging of  | 9      |
| 10     |               | Event-Related fMRI Trials to Test the Retrieval Success Hypothesis. NeuroImage, 7, 163–175.   | 10     |
| 11     | [12]          | BULLMORE, E., BRAMMER, M., WILLIAMS, S. C. R., RABE-HESKETH, S., JANOT, N., DAVID, A., MELLERS,   | 11     |
| 12     |               | J., HOWARD, R. and SHAM, P. (1996). Statistical methods of estimation and inference for functional MR image analysis. <i>IEEE Trans. Med. Ima.</i> <b>35</b> , 261–277  | 12     |
| 13     | [13]          | Bullmore, E. T., Suckling, J., Overmeyer, S., Rabe-Hesketh, S., Taylor, E. and Brammer, M.  | 13     |
| 14     | . ,           | J. (1999). Statistical methods of estimation and inference for functional MR image analysis. <i>IEEE</i>  | 14     |
| 15     | <b>6</b> - 11 | Trans. Med. Img., 18, 32–42.  | 15     |
| 16     | [14]          | CAULHOUN, V. D., STEVENS, M. C., PEARLSON, G. D. and KIEHL, K. A. (2004). IMRI analysis with the  | 16     |
| 17     |               | derivative terms. NeuroImage, 22, 252–257.  | 17     |
| 10     | [15]          | Cox, R. W. (1996). AFNI: Software for analysis and visualization of functional magnetic resonance   | 10     |
| 10     | [1.0]         | neuroimages. Computers and Biomedical Research, 29, 162–173. http://afni.ninh.nih.gov/afni.   | 10     |
| 19     | [16]          | DALE, A. M. and BUCKNER, R. L. (1997). Selective Averaging of Rapidly Presented Individual Trials<br>Using fMRL Human Brain Manning 5, 329–340  | 19     |
| 20     | [17]          | DELONG, M. R., ALEXANDER, G. E., GEORGOPOULOS, A. P., CRUTCHER, M. D., MITCHELL, S. J. and  | 20     |
| 21     | . ,           | RICHARDSON, R. T. (1984). Role of basal ganglia in limb movements. Hum. Neurobiol., 2, 235–244.   | 21     |
| 22     | [18]          | DUANN, J. R., JUNG, T. P., KUO, W. J., YEH, T. C., MAKEIG, S., HSIEH, J. C. and SEJNOWSKI, T. J. (1999). Single Third Visit Visit in Front Political POL D Single Magnetization 202, 225                      | 22     |
| 23     | [19]          | (1998). Single-Irial variability in Event-Related BOLD Signals. <i>Neuroimage</i> , 823–835.<br>Fried I Katz A McCarthy G Sass K J Williamson P and Spencer S S and Spencer                                   | 23     |
| 24     | [10]          | D. D. (1991). Functional organization of human supplementary motor cortex studied by electrical   | 24     |
| 25     |               | stimulation. J. Neurosci., 11, 3656-3666.   | 25     |
| 26     | [20]          | FRISTON, K., HOLMES, A., WORSLEY, K., POLINE, J. B., GRASBY, P., WILLIAMS, S., FRACKOWIAK, R.   | 26     |
| 27     | [21]          | and IURNER, R. (1995). Analysis of IMRI time series revisited. <i>NeuroImage</i> , 2, 45–53.<br>FRISTON K J ASHRIENER J KIEBEL S J NICHOLS T E and PENNY W D (2007) Statistical                               | 27     |
| 28     | []            | Parametric Mapping: The Analysis of Functional Brain Images, Academic Press. http://www.fil.  | 28     |
| 29     |               | ion.ucl.ac.uk/spm/.   | 29     |
| 30     | [22]          | FRISTON, K. J., FRITH, C. D., TURNER, R. and FRACKOWIAK, R. S. J. (1995). Characterizing Evoked   | 30     |
| 31     | [23]          | FRISTON, K. J., HOLMES, A. P., WORSLEY, K. J., POLINE, J. P., FRITH, C. D. and FRACKOWIAK, R.   | 31     |
| 32     |               | S. J. (1995). Statistical parametric maps in functional imaging: a general linear approach. Human   | 32     |
| 33     | [0,1]         | Brain Mapping, <b>2</b> , 189–210.  | 33     |
| 34     | [24]          | FRISTON, K. J., JEZZARD, P. and TURNER, R. (1994). Analysis of Functional MRI Time-Series. Human<br>Brain Manning 1 153–171   | 34     |
| 35     | [25]          | FRISTON, K. J., JOSEPHS, O., REES, G. and TURNER, R. (1998). Nonlinear event-related responses in   | 35     |
| 36     |               | fMRI. Magnetic Resonance in Medicine, <b>39</b> , 41–52.  | 36     |
| 27     | [26]          | GENOVESE, C. R. (2000). A Bayesian Time-Course Model for Functional Magnetic Resonance Imag-  | 27     |
| 20     | [27]          | GLOVER, G. H. (1999). Deconvolution of impulse response in event-related fmri. NeuroImage, 9.   | 20     |
| 30     | [=.]          | 416-429.  | 30     |
| 39     | [28]          | Gössl, C., FAHRMEIR, L. and AUER, D. P. (2001). Bayesian Modeling of the hemodynamic response   | 39     |
| 40     | [20]          | function in BOLD fMRI. NeuroImage, 14, 140–148.   | 40     |
| 41     | [29]          | fMRI using smooth FIR filters. <i>IEEE Transactions on Medical Imagina</i> , <b>19</b> , 1188–1201.   | 41     |
| 42     | [30]          | HUETTEL, S. A., SONG, A. W. and MCCARTHY, G. (2004). Functional Magnetic Resonance Imaging.   | 42     |
| 43     |               | Sinauer Associates, Inc   | 43     |
| 44     | [31]          | JOSEPHS, O., TURNER, R. and FRISTON, K. (1995). Event-Related fMRI. Human Brain Mapping, 5, 243–248   | 44     |
| 45     | [32]          | KANG, J. K., BÉNAR, C. G., AL-ASMI, A., KHANI, Y. A., PIKE, G. B., DUBEAU, F. and GOTMAN, J.  | 45     |
| 46     |               | (2003). Using patient-specific hemodynamic response functions in combined EEG-fMRI studies in   | 46     |
| 47     | [0.0]         | epilepsy. NeuroImage, <b>20</b> , 1162-1170.  | 47     |
| 48     | [33]          | LANGE, N. and ZEGER, S. L. (1997). Non-linear Fourier Time Series Analysis for Human Brain<br>Mapping by Europian Magnetic Resonance Imaging (with discussion). <i>IRSS</i> 46:1 1, 20                        | 48     |
| 49     | [34]          | LINDQUIST, M. and WAGER, T. D. (2007). Validity and power in hemodynamic response modeling:   | 49     |
| 50     | . 1           | A comparison study and a new approach. Human Brain Mapping, Available online.   | 50     |
| 51     | [35]          | LINDQUIST, M. and WAGER, T. D. (2006). Validity and Power in Hemodynamic Response Modeling:   | 51     |
|        |               |   |        |

|    | 216   | Bai, Truong and Huang   |    |
|----|-------|---|----|
|    |       |   |    |
| 1  |       | A comparison study and a new approach. Human Brain Mapping, DOI: 10.1002/hbm.20310.   | 1  |
| 2  | [36]  | MALDJIAN, J. A., LAURIENTI, P. J. and BURDETTE, J. H. (2004). Precentral Gyrus Discrepancy in   | 2  |
| 3  | [37]  | MARCHINI, J. L. and RIPLEY, B. D. (2000). A new statistical approach to detecting significant acti-                                       | 3  |
| 4  |       | vation in functional MRI. NeuroImage, 12, 366–380.  | 4  |
| 5  | [38]  | MIEZIN, F. M., MACCOTTA, L., OLLINGER, J. M., PETERSEN, S. E. and BUCKNER, R. L. (2000). Char-  | 5  |
| 6  |       | Possibility of Ordering Brain Activity Based on Relative Timing. <i>NeuroImage</i> , <b>11</b> , 735-759.                                 | 6  |
| 7  | [39]  | NEWTON, H. J. (1988). Timeslab: A Time Series Analysis Laboratory. Wadsworth & Brooks/Cole.   | 7  |
| 8  | [40]  | S. OGAWA, D. W. TANK, R. MENON, J. M. ELLERMANN, S. G. KIM, H. MERKLE and K. UĞURBIL (1992).  | 8  |
| 9  |       | Intrinsic signal changes accompanying sensory stimulation: functional brain mapping with magnetic resonance imaging. PNAS, 89, 5951–5955. | 9  |
| 10 | [41]  | ROSEN, B. R., BUCKNER, R. L. and DALE, A. M. (1998). Event-related functional MRI: past, present,   | 10 |
| 11 |       | and future. Proceedings of the National Academy of Sciences of the United States of America,  | 11 |
| 12 | [42]  | SCHWARTZ, L. (1966). Mathematics for the Physical Sciences. Hermann.  | 12 |
| 13 | [43]  | Smith, S. M., Jenkinson, M., Woolrich, M. W., Beckmann, C. F., Behrens, T. E. J., Johansen-   | 13 |
| 14 |       | BERG, H., BANNISTER, P. R., DE LUCA, M., DROBNJAK, I. AND FLITNEY, D. E., NIAZY, R., SAUNDERS,  | 14 |
| 15 |       | in functional and structural MR image analysis and implementation as FSL. NeuroImage, 23(S1).   | 15 |
| 16 |       | 208-219. http://www.fmrib.ox.ac.uk/fsl/.  | 16 |
| 17 | [44]  | WOOLRICH, M. W., BEHRENS, T. E. and SMITH, S. M. (2004). Constrained linear basis sets for HRF  | 17 |
| 18 | [45]  | WORSLEY, K. and FRISTON, K. (1995). Analysis of fMRI time series revisited-again. NeuroImage, 2.  | 18 |
| 19 | 1 - 1 | 173–181.  | 19 |
| 20 | [46]  | ZARAHN, E., AGUIRRE, G. K. and D'ESPOSITO, M. (1997). Empirical analyses of BOLD fMRI statistics.   | 20 |
| 21 | [47]  | ZHANG, C. M., JIANG, Y. and YU, T. (2007). A comparative study of one-level and two-level semi-   | 21 |
| 22 | L . J | parametric estimation of hemodynamic response function for fMRI data. Statistics in Medicine,   | 22 |
| 23 |       | <b>26</b> , 3845–3861.  | 23 |
| 24 |       |   | 24 |
| 25 |       |   | 25 |
| 26 |       |   | 26 |
| 27 |       |   | 27 |
| 28 |       |   | 28 |
| 29 |       |   | 29 |
| 30 |       |   | 30 |
| 31 |       |   | 31 |
| 32 |       |   | 32 |
| 33 |       |   | 33 |
| 34 |       |   | 34 |
| 35 |       |   | 35 |
| 36 |       |   | 36 |
| 37 |       |   | 37 |
| 38 |       |   | 38 |
| 39 |       |   | 39 |
| 40 |       |   | 40 |
| 41 |       |   | 41 |
| 42 |       |   | 42 |
| 43 |       |   | 43 |
| 44 |       |   | 44 |
| 45 |       |   | 45 |
| 46 |       |   | 46 |
| 47 |       |   | 47 |
| 48 |       |   | 48 |
| 49 |       |   | 49 |
| 50 |       |   | 50 |
| 51 |       |   | 51 |

| Penalized Least-Squares         Jolanda Muñoz Maldonado¹         Michigan Technological University         Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalence for smoothing spline estimators and best linear unbiased prediction and/or posterior analysis of certain Gaussian signal-plus-noise models is estamined in a more general setting. These connections allow for the application of an efficient, linear time algorithm, to estimate parameters, compute random effects predictions and evaluate likelihoods in a large class of model scenarios. We also show that the methods of generalized cross-validation, restricted maximum likelihood and unbiased risk prediction can be used to estimate the variance components or adaptively select the smoothing parameters in any of the three settings.         Contents       2         Introduction       2         Equivalence Theorem       2         3.1       Varying Coefficient Models       2         3.2       Ridge Regression and Penalized Spline Regression       2         3.3       Mixed-Effects Model       2         3.4       State-Space Forms       2         4.5       state-Space Forms       2         6.       mmodel methodology, penalized least-squares and Bayesian randor flects models are widely used statistical tools. However, due to the dissimilar rule of the settings in which they are typicaly formulated, connections betwe hese three   | Penalized Least-Squares         Yolanda Muñoz Maldonado¹         Michigan Technological University         Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalence for smoothing spline estimators and best linear unbiased prediction and/or posterior analysis of certain Gaussian signal-plus-noise models is examined in a more general setting. These connections allow for the application of an efficient, linear time algorithm, to estimate parameters, compute random effects predictions and evaluate likelihoods in a large class of model scenarios. We also show that the methodos of generalized cross-validation, restricted maximum likelihood and unbiased risk prediction can be used to estimate the variance components or adaptively select the smoothing parameters in any of the three settings.         Contents         1       Introduction         2       Equivalence Theorem         3.1       Varying Coefficient Models         3.2       Ridge Regression and Penalized Spline Regression         3.3       Mixed-Effects Model         4       Summary         A State-Space Forms         References         Mixed-effects model methodology, penalized least-squares and Bayesian re effects models are widely used statistical tools. However, due to the dissimit ture of the settings in which they are typically formulated, connections by these three techniques as well as the fundamental reasons of the connection of the set words has the fundamen |        |
|--|---|--------|
| <b>Yolanda Muñoz Maldonado</b> <sup>1</sup> Michigan Technological University         Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalance for smoothing spline estimators and best linear unbiased prediction of and/or posterior analysis of certain Gaussian signal-plus-noise models is examined in a more general setting. These connections allow for the application of afficient, linear time algorithm, to estimate parameters, compute random effects predictions and evaluate likelihoods in a large class of model scenarios. We also show that the methods of generalized cross-validation, restricted maximum likelihood and unbiased risk prediction can be used to estimate the variance components or adaptively select the smoothing parameters in any of the three settings. <b>Contents</b> <u>Introduction</u> 2 <u>Equivalence Theorem</u> 2 <u>2</u> . Ridge Regression and Penalized Spline Regression       2 <u>3</u> . Mixed-Effects Model       2 <u>3</u> . Mixed-Effects Model       2 <u>Summary</u> 2 <u>A</u> State-Space Forms       2 <u>C</u> Acterication       2 <u>Mixed-effects</u> model methodology, penalized least-squares and Bayesian randor fiects models are widely used statistical tools. However, due to the dissimilar rise in other set tings in which they are typically formulated, connections betwee hese three techniques as well as the fundamental reasons for the connections, he fien been overloo  | <b>Yolanda Muñoz Maldonado</b> <sup>1</sup> Michigan Technological University         Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalence for smoothing spline estimators and best linear unbiased prediction and evaluate likelihoods in a large class of model scenarios. We also show that the methods of generalized cross-validation, restricted maximum likelihood and unbiased risk prediction can be used to estimate the maximum likelihood and unbiased risk prediction can be used to estimate the three settings. <b>Contents</b> 1       Introduction         2       Equivalence Theorem         3       Examples         3.1       Varying Coefficient Models         3.3       Mixed-Effects Model         4       Summary         A State-Space Forms         References   |        |
| Michigan Technological University         Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalence for smoothing spline estimators and best linear unbiased prediction and/or posterior analysis of certain Gaussian signal-plus-noise models is examined in a more general setting. These connections allow for the application of an efficient, linear time algorithm, to estimate parameters, compute random effects predictions and evaluate likelihoods in a large class of model scenario. We also show that the methods of generalized cross-validation, restricted maximum likelihood and unbiased risk prediction can be used to estimate the variance components or adaptively select the smoothing parameters in any of the three settings.         Contents       2         Introduction       2         Equivalence Theorem       2         3.1 Varying Coefficient Models       2         3.2 Ridge Regression and Penalized Spline Regression       2         3.3 Mixed-Effects Model       2         Summary       2         Astae-Space Forms       2         Astae-Space Forms       2         Active-effects model methodology, penalized least-squares and Bayesian random fiects models are widely used statistical tools. However, due to the dissimilar prime of the settings in which they are typically formulated, connections betwee hese three techniques as well as the fundamental reasons for the connections, has the noverlooked. In this paper, we review some of the well known resu hat | Michigan Technological University         Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalence for smoothing spline estimators and best linear unbiased prediction and/or posterior analysis of certain Gaussian signal-plus-noise models is examined in a more general setting. These connections allow for the application of effects predictions and evaluate likelihoods in a large class of model scenarios. We also show that the methods of generalized cross-validation, restricted wariance components or adaptively select the smoothing parameters in any of the three settings.         Contents         1       Introduction         2       Equivalence Theorem         3       Examples         3.1       Varying Coefficient Models         3.2       Ridge Regression and Penalized Spline Regression         3.3       Mixed-Effects Model         4       Summary         A State-Space Forms         References         1       Introduction   |        |
| Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalence for smoothing spline estimators and best linear unbiased prediction and/or posterior analysis of certain Gaussian signal-plus-noise models is examined in a more general setting. These connections allow for the application of an efficient, linear time algorithm, to estimate parameters, compute random effects predictions and evaluate likelihoods in a large class of model scenarios. We also show that the methods of generalized cross-validation, restricted maximum likelihood and unbiased risk prediction can be used to estimate the variance components or adaptively select the smoothing parameters in any of the three settings.         Contents         Introduction       2         Equivalence Theorem       2         2.1< Varying Coefficient Models   | Abstract: This paper reviews the connections between estimators that derive from three different modeling methodologies: Mixed-effects models, Bayesian models and Penalized Least-squares. Extension of classical results on the equivalence for smoothing spline estimators and best linear unbiased prediction and/or posterior analysis of certain Gaussian signal-plus-noise models is examined in a more general setting. These connections allow for the application of an efficient, linear time algorithm, to estimate parameters, compute random effects predictions and evaluate likelihoods in a large class of model scenarios. We also show that the methods of generalized cross-validation, restricted maximum likelihood and unbiased risk prediction can be used to estimate the variance components or adaptively select the smoothing parameters in any of the three settings.         Contents         1       Introduction  |        |
| Contents         Introduction       2         Equivalence Theorem       2         Examples       2         3.1 Varying Coefficient Models       2         3.2 Ridge Regression and Penalized Spline Regression       2         3.3 Mixed-Effects Model       2         Summary       2         A State-Space Forms       2         References       2         Introduction       2         Mixed-effects model methodology, penalized least-squares and Bayesian random flects models are widely used statistical tools. However, due to the dissimilar rule of the settings in which they are typically formulated, connections betwee hese three techniques as well as the fundamental reasons for the connections, ha ften been overlooked. In this paper, we review some of the well known resu hat connect smoothing spline estimators, Gaussian signal-plus-noise models an est linear unbiased prediction of mixed effects models and show that they are hese three techniques and spline estimators.   | Contents         1 Introduction         2 Equivalence Theorem         3 Examples         3.1 Varying Coefficient Models         3.2 Ridge Regression and Penalized Spline Regression         3.3 Mixed-Effects Model         4 Summary         A State-Space Forms         References         I. Introduction         Mixed-effects model methodology, penalized least-squares and Bayesian rate         effects models are widely used statistical tools. However, due to the dissimit ture of the settings in which they are typically formulated, connections b         these three techniques as well as the fundamental reasons for the connection often been overlooked. In this paper, we review some of the well known that connect smoothing spline estimators, Gaussian signal-plus-noise models and short they are typically and short they are the theorem of the setting spline estimators.  |        |
| Introduction       2         Equivalence Theorem       2         Examples       2         3.1 Varying Coefficient Models       2         3.2 Ridge Regression and Penalized Spline Regression       2         3.3 Mixed-Effects Model       2         Summary       2         A State-Space Forms       2         References       2         Introduction       2         Mixed-effects model methodology, penalized least-squares and Bayesian random ffects models are widely used statistical tools. However, due to the dissimilar rule of the settings in which they are typically formulated, connections betwee hese three techniques as well as the fundamental reasons for the connections, ha ften been overlooked. In this paper, we review some of the well known result hat connect smoothing spline estimators, Gaussian signal-plus-noise models are widely and affecta models and show that they are here the theoremeters.  | 1       Introduction         2       Equivalence Theorem         3       Examples         3       Introduction         3.1       Varying Coefficient Models         3.2       Ridge Regression and Penalized Spline Regression         3.3       Mixed-Effects Model         4       Summary         5       Examples         6       State-Space Forms         7       References         8       References         9       References         9       Mixed-effects model methodology, penalized least-squares and Bayesian ra         9       effects models are widely used statistical tools. However, due to the dissimit         9       ture of the settings in which they are typically formulated, connections b         9       these three techniques as well as the fundamental reasons for the connection         9       paper, we review some of the well known         9       that connect smoothing spline estimators, Gaussian signal-plus-noise mode  |        |
| Equivalence Theorem       2         3.1 Varying Coefficient Models       2         3.1 Varying Coefficient Models       2         3.2 Ridge Regression and Penalized Spline Regression       2         3.3 Mixed-Effects Model       2         3.3 Mixed-Effects Model       2         Summary       2         A State-Space Forms       2         References       2         References       2         Introduction       2         Mixed-effects model methodology, penalized least-squares and Bayesian random ffects models are widely used statistical tools. However, due to the dissimilar rule of the settings in which they are typically formulated, connections betwee hese three techniques as well as the fundamental reasons for the connections, has ften been overlooked. In this paper, we review some of the well known resu hat connect smoothing spline estimators, Gaussian signal-plus-noise models a west linear unbiased prediction of mixed effects models and show that they are here the test they are here the set of the                   | <ul> <li>2 Equivalence Theorem</li></ul>  | . 21'  |
| 3.1       Varying Coefficient Models       2         3.1       Varying Coefficient Models       2         3.2       Ridge Regression and Penalized Spline Regression       2         3.3       Mixed-Effects Model       2         3.3       Mixed-Effects Model       2         4       Summary       2         5       Summary       2         4       State-Space Forms       2         2       References       2         2       References       2         2       References       2         2       References       2         3       Introduction       2         4       Introduction       2         5       Introduction       2         6       Introduction       2         7       Introduction       2         7       Introduction       2         8       Introduction       2         9       Inthey are typically formulated, connections between these t   | <ul> <li>3 Examples</li></ul>   | . 22   |
| 3.1       Varying Coefficient Models       2         3.2       Ridge Regression and Penalized Spline Regression       2         3.3       Mixed-Effects Model       2         3.3       Mixed-Effects Model       2         4       Summary       2         4       Summary       2         4       State-Space Forms       2         2       References       2         2       References       2         2       References       2         3       Introduction       2         4       Introduction       2         5       Introduction       2         6       Introduction       2         6       State-Space Forms       2         7       Introduction       2         7       Introduction       2         7       Introduction       2         7       Introduction       2         8       Introduction       2         9       Introduction       2         9       Introduction       2         9       Introduction       2         9       Inthey are typically formulated, connections betwee   | 3.1       Varying Coefficient Models         3.2       Ridge Regression and Penalized Spline Regression         3.3       Mixed-Effects Model         4       Summary         4       Summary         5       A State-Space Forms         6       References         7       References         8       References         9       Mixed-effects model methodology, penalized least-squares and Bayesian rate         9       effects models are widely used statistical tools. However, due to the dissimination of the settings in which they are typically formulated, connections b         10       these three techniques as well as the fundamental reasons for the connection often been overlooked. In this paper, we review some of the well known that connect smoothing spline estimators, Gaussian signal-plus-noise models are the production of the setting of mixed effects models and show that the theorem.   | . 22   |
| 3.2       Ridge Regression and Penalized Spline Regression       2         3.3       Mixed-Effects Model       2         4       Summary       2         4       Summary       2         A       State-Space Forms       2         References       2         References       2         .       Introduction         Mixed-effects model methodology, penalized least-squares and Bayesian random         effects models are widely used statistical tools. However, due to the dissimilar r         ure of the settings in which they are typically formulated, connections betwee         hese three techniques as well as the fundamental reasons for the connections, has         ften been overlooked. In this paper, we review some of the well known resu         hat connect smoothing spline estimators, Gaussian signal-plus-noise models a         act linear unbiased prediction of mixed effects models and show that they are here  | 3.2       Ridge Regression and Penalized Spline Regression         3.3       Mixed-Effects Model         4       Summary         4       Summary         A       State-Space Forms         A       State-Space Forms         References       References         1.       Introduction         Mixed-effects model methodology, penalized least-squares and Bayesian rateffects models are widely used statistical tools. However, due to the dissimitation of the settings in which they are typically formulated, connections between the techniques as well as the fundamental reasons for the connection often been overlooked. In this paper, we review some of the well known that connect smoothing spline estimators, Gaussian signal-plus-noise models are the production of the setting of mixed effects models are unbiased production of the set model and show that the production of the set models are spline estimators.  | . 22   |
| 3.3 Mixed-Effects Model       2         Summary       2         A State-Space Forms       2         References       2         References       2         . Introduction       2         Mixed-effects model methodology, penalized least-squares and Bayesian random         affects models are widely used statistical tools. However, due to the dissimilar r         ure of the settings in which they are typically formulated, connections betwee         hese three techniques as well as the fundamental reasons for the connections, has         aften been overlooked. In this paper, we review some of the well known resu         hat connect smoothing spline estimators, Gaussian signal-plus-noise models a         aget linear unbiased prediction of mixed affects medels and show that they are here   | 3.3 Mixed-Effects Model         4 Summary         5 State-Space Forms         A State-Space Forms         References         1. Introduction         Mixed-effects model methodology, penalized least-squares and Bayesian rate         effects models are widely used statistical tools. However, due to the dissiming         curve of the settings in which they are typically formulated, connections be         chese three techniques as well as the fundamental reasons for the connection         often been overlooked. In this paper, we review some of the well known         chat connect smoothing spline estimators, Gaussian signal-plus-noise models  | . 22   |
| A State-Space Forms  | 4 Summary   | . 23   |
| A State-Space Forms  | A State-Space Forms   | . 23   |
| References   | <b>I. Introduction</b><br>Mixed-effects model methodology, penalized least-squares and Bayesian ra<br>effects models are widely used statistical tools. However, due to the dissimi-<br>cure of the settings in which they are typically formulated, connections b<br>chese three techniques as well as the fundamental reasons for the connection<br>often been overlooked. In this paper, we review some of the well known<br>that connect smoothing spline estimators, Gaussian signal-plus-noise models<br>were the production of mixed effects models and show that they are   | . 23   |
| <b>1. Introduction</b><br>Mixed-effects model methodology, penalized least-squares and Bayesian random<br>effects models are widely used statistical tools. However, due to the dissimilar run<br>ure of the settings in which they are typically formulated, connections betwee<br>hese three techniques as well as the fundamental reasons for the connections, has<br>fiten been overlooked. In this paper, we review some of the well known resu<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise models a<br>west linear unbiased prediction of mixed effects models and show that they are here.   | 1. Introduction<br>Mixed-effects model methodology, penalized least-squares and Bayesian ra<br>effects models are widely used statistical tools. However, due to the dissimi-<br>cure of the settings in which they are typically formulated, connections b<br>chese three techniques as well as the fundamental reasons for the connection<br>often been overlooked. In this paper, we review some of the well known<br>that connect smoothing spline estimators, Gaussian signal-plus-noise models<br>over the production of mixed effects models and show that they are  | . 23   |
| Mixed-effects model methodology, penalized least-squares and Bayesian random<br>effects models are widely used statistical tools. However, due to the dissimilar run<br>ure of the settings in which they are typically formulated, connections betwee<br>hese three techniques as well as the fundamental reasons for the connections, has<br>often been overlooked. In this paper, we review some of the well known resu<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise models a<br>pert linear unbiased prediction of mixed effects models and show that there are h  | Mixed-effects model methodology, penalized least-squares and Bayesian ra<br>effects models are widely used statistical tools. However, due to the dissimi-<br>cure of the settings in which they are typically formulated, connections b<br>chese three techniques as well as the fundamental reasons for the connection<br>often been overlooked. In this paper, we review some of the well known<br>that connect smoothing spline estimators, Gaussian signal-plus-noise models<br>can be an endiption of mixed effects models and show that they are   |        |
| Mixed-effects model methodology, penalized least-squares and Bayesian random<br>effects models are widely used statistical tools. However, due to the dissimilar run<br>run of the settings in which they are typically formulated, connections betwee<br>hese three techniques as well as the fundamental reasons for the connections, has<br>often been overlooked. In this paper, we review some of the well known result<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise models a<br>pert linear unbiased prediction of mixed effects models and show that they are here.   | Mixed-effects model methodology, penalized least-squares and Bayesian ra<br>effects models are widely used statistical tools. However, due to the dissimi-<br>ture of the settings in which they are typically formulated, connections b<br>these three techniques as well as the fundamental reasons for the connection<br>often been overlooked. In this paper, we review some of the well known<br>that connect smoothing spline estimators, Gaussian signal-plus-noise models<br>est linear unbiased prediction of mixed effects models and show that the   |        |
| effects models are widely used statistical tools. However, due to the dissimilar r<br>ure of the settings in which they are typically formulated, connections betwee<br>hese three techniques as well as the fundamental reasons for the connections, has<br>often been overlooked. In this paper, we review some of the well known resu<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise models a<br>post linear unbiased prediction of mixed effects models and show that they are h   | effects models are widely used statistical tools. However, due to the dissimi-<br>cure of the settings in which they are typically formulated, connections b<br>chese three techniques as well as the fundamental reasons for the connection<br>often been overlooked. In this paper, we review some of the well known<br>that connect smoothing spline estimators, Gaussian signal-plus-noise mode<br>page t linear unbiased prediction of mixed effects models and show that they are   | ndom   |
| ure of the settings in which they are typically formulated, connections betwee<br>hese three techniques as well as the fundamental reasons for the connections, has<br>often been overlooked. In this paper, we review some of the well known resu<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise models a<br>post linear unbiased prediction of mixed effects models and show that they are h   | ure of the settings in which they are typically formulated, connections b<br>hese three techniques as well as the fundamental reasons for the connection<br>often been overlooked. In this paper, we review some of the well known<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise mode<br>part linear unbiased prediction of mixed effects models and show that they  | lar na |
| hese three techniques as well as the fundamental reasons for the connections, ha<br>often been overlooked. In this paper, we review some of the well known resu<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise models a<br>post linear unbiased prediction of mixed effects models and show that they are h  | hese three techniques as well as the fundamental reasons for the connection<br>often been overlooked. In this paper, we review some of the well known<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise mode   | etwee  |
| often been overlooked. In this paper, we review some of the well known resu<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise models a<br>post linear unbiased prediction of mixed effects models and show that there are h   | often been overlooked. In this paper, we review some of the well known<br>hat connect smoothing spline estimators, Gaussian signal-plus-noise mode<br>part linear unbiased prediction of mixed effects models and show that they  | s, hav |
| hat connect smoothing spline estimators, Gaussian signal-plus-noise models a   | hat connect smoothing spline estimators, Gaussian signal-plus-noise mode  | result |
| uset linear unbiaged prediction of mixed effects models and show that they are h   | agt linear unbiaged prediction of mired offects models and show that they   | els an |
| rest mean unbrased prediction of mixed-effects models and show that they are p   | Jest mean unblased prediction of mixed-enects models and show that they a   | re bu  |
| one aspect of a general framework that allows for "cross-platform" development   | one aspect of a general framework that allows for "cross-platform" developm   | nent i |
| nixed-effects models, using frequentist or Bayesian approaches, and/or penaliz   | nixed-effects models, using frequentist or Bayesian approaches, and/or per  | nalize |
|  |   |        |
| IZEDALINENE UL MALIEINALUAL OLENCES, MICHPALL LECHNOLOPICAL ULIVEISILV HONORON   | ISA, email: vmunoz@mtu.edu  |        |

Keywords and phrases: Smoothing splines, p-splines, ridge regression, varying coefficient mod els, Bayesian prediction, Kalman fiter, adaptive selection

least-squares (PLS) criteria.

The relationship between particular cases of frequentist and Bayesian mixed-З effects models and PLS has been exposed before. For example, Lindley and Smith [29] proposed the use of prior information on the parameters of a fixed effects linear model under the assumption of the parameters having exchangeable distributions. In the early development of the Bayesian theory for smoothing splines, Wahba [43] noticed the intimate connection between estimators resulting from spline smooth-ing and a Gaussian model with diffuse initial conditions. Robinson [35] remarked on applications of Best Linear Unbiased Predictors (BLUP's) for estimation of variance parameters, randomized block designs and their link to empirical Bayes methods and Kriging. Speed [40] pointed out, in a comment to Robinson's article. that smoothing spline estimators were in fact BLUP's of a certain mixed effects model. In the PLS framework, it is well known that smoothing splines estimators are a special case of penalized splines estimators (P-splines) [37]. Wahba [47] and Cressie [7] discussed the links between splines and kriging estimates and Nychka used the representation of smoothing splines as a type of ridge estimator to further relate smoothing spline estimation and kriging [32].

More recently, researchers have been using the connection between smoothing spline estimators and particular mixed-effects models to compute smoothing spline estimators [see 4, 19, 48]. Ruppert et al. [36] mentioned the correspondence be-tween penalized spline smoothers and prediction in the mixed-effects model and remarked on the advantages of using existing mixed-effects model techniques and software in a semi-parametric regression setting. Eubank et al. [11] took advantage of the relationship between smoothing splines and the Gaussian model of [43] to provide a general development that includes the efficient computation of estimators in a varying coefficient model context. 

Using connections that have been established for various special cases, we syn-thesize them and present a formal result that details precisely when penalized least-squares estimation, BLUP for a mixed-effects model and posterior mean analysis of a mixed-effects model with diffuse priors on some of the random effects (hereafter referred to as simply the Bayesian model) produce identical estimators. We then describe how this can be exploited in many cases of interest to provide a computationally efficient algorithm for evaluation of estimators and likelihoods, computation of predictions, and construction of Bayesian prediction intervals. The implemented algorithm reduces the computational effort of calculating the aforementioned quan-tities by two orders of magnitude over what would normally be the case for a direct mixed-effects model approach. We also establish a result showing that the methods of Generalized Cross-validation (GCV), Restricted Maximum Likelihood (REML) or the equivalent technique of Generalized Maximum Likelihood (GML) and Unbiased Risk Prediction (UBR) can be used in any of the three settings to adaptively estimate the smoothing parameters or variance components.

The following three examples will be used throughout the paper to illustrate the utility of our approach.

**Example 1: Varying Coefficient Models.** Varying coefficient models generalize ordinary linear regression models by allowing for regression coefficients that change dynamically as a function of independent variables. The simplest example of this are the so-called time varying coefficient models where there is only one effect modifying covariate. In that setting, we have response variables  $y_{ij}$ , i = 1, ..., n,  $j = 1, ..., n_i$ , that depend on some predictor variables  $x_{1ij}, ..., x_{Kij}$  through a relationship of the

form

(1)

 $y_{ij} = \sum_{k=1}^{K} \beta_k(t_{ij}) x_{kij} + e_{ij},$ 

where the  $\beta_k(\cdot)$ 's are unknown coefficient functions of a covariate t and the  $e_{ij}$  represent random error terms. Models like (1) were first introduced by [21] who proposed obtaining estimators through minimization of the PLS criterion

(2) 
$$\sum_{i=1}^{n} \sum_{j=1}^{n_i} \left\{ y_{ij} - \sum_{k=1}^{K} f_k(t_{ij}) x_{kij}(t_{ij}) \right\}^2 + \sum_{k=1}^{K} \lambda_k \int_0^1 [f_k^{(r)}(t)]^2 dt$$

over functions  $f_1, \ldots, f_K$  having r square integrable derivatives, and  $g^{(s)}(t)$  being the s<sup>th</sup> derivative of the function g. The parameters  $\lambda_k \geq 0$  control the smoothness of the coefficient functions and the minimizers can be shown to be natural splines of degree 2r - 1 with knots at the unique elements of the set  $\{t_{ij}\}$ .

**Example 2: Ridge Regression.** Consider the linear regression model

where  $\boldsymbol{y}$  is a  $n \times 1$  vector of responses,  $\mathbf{X}$  is a known  $n \times p$  matrix of predictor variables of rank p,  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of unknown coefficients and  $\boldsymbol{e}$  is a normally distributed vector of errors with  $\mathbf{E}(\boldsymbol{e}) = \mathbf{0}$  and  $\mathbf{E}(\boldsymbol{e}\boldsymbol{e}^{\mathrm{T}}) = \sigma_{\boldsymbol{e}}^{2}\mathbf{I}$ , with "T" denoting the transpose of a matrix and  $\mathbf{I}$  an identity matrix of suitable dimension. The generalized ridge regression estimator of  $\boldsymbol{\beta}$  is then given by  $\hat{\boldsymbol{\beta}} = [\mathbf{X}^{\mathrm{T}}\mathbf{X} + \mathbf{K}]^{-1}\mathbf{X}^{\mathrm{T}}\boldsymbol{y}$ . This estimator can be obtained by minimizing the PLS criterion

(4) 
$$(\boldsymbol{y} - \mathbf{X}\boldsymbol{a})^{\mathrm{T}}(\boldsymbol{y} - \mathbf{X}\boldsymbol{a}) + \boldsymbol{a}^{\mathrm{T}}\mathbf{K}\boldsymbol{a}$$

over  $\{a : a \in \mathbb{R}^p\}$ , with **K** a diagonal matrix having elements  $\lambda_i \geq 0$ , for  $i = 1, \ldots, p$ . A special instance of (4) is given by ordinary ridge regression in which case the predictor variables are usually standardized and **K** has the form  $\lambda \mathbf{I}$ , for  $\lambda > 0$ . Other variations of generalized ridge regression are the P-splines estimators of [9] and of [36]. We will now describe the later approach in more detail.

Suppose that we have a collection of points on the plane,  $(t_i, y_i)$ , i = 1, ..., n, and want to fit them using scatter-plot smoothing methodology. P-splines provide one popular approach for accomplishing this that arise from using a spline basis to construct the **X** matrix in (3). That is, for some integer  $m \ge 0$  and a fixed set of knots  $\xi_1 < \xi_2 < \cdots < \xi_p$ , we take  $\mathbf{X} = [\mathbf{x}_1, \ldots, \mathbf{x}_{m+p}]$  with  $\mathbf{x}_1$  a *n*-vector of all ones,  $\mathbf{x}_j = [t_1^{j-1}, \ldots, t_n^{j-1}]^{\mathrm{T}}$ ,  $j = 2, \ldots, m$  and  $\mathbf{x}_{m+j} = [(t_1 - \xi_j)_+^{m-1}, \ldots, (t_n - \xi_j)_+^{m-1}]^{\mathrm{T}}$ , for  $j = 1, \ldots, p$  with  $(x)_+^r$  being  $x^r$  for  $x \ge 0$  and zero otherwise. A P-spline smoother is then found by minimizing (4), with the matrix **K** having the form

(5) 
$$\mathbf{K} = \begin{bmatrix} \mathbf{0}_{m \times m} & \mathbf{0}_{m \times p} \\ \mathbf{0}_{p \times m} & \lambda \mathbf{I} \end{bmatrix},$$

with  $\mathbf{0}_{r \times s}$  being an r by s matrix of all zeros.

Example 3: Randomized Block Design. Linear mixed-effects models have been
 applied for analysis of data arising from situations involving repeated measures and
 experimental designs with factors that can be seen as a combination of fixed and

З

random effects. Some types of randomized block designs fall in the last category, for example, when the experimental units are randomly selected and each has repeated measurements. For this particular type of design, the experimental units are assumed to be the factor (or blocking criterion), that makes them relatively homogeneous with respect to a measured response. One way of modeling this type of problems is

8 (6) 
$$\boldsymbol{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{b} + \boldsymbol{e}$$

where **X** is the design matrix for the fixed-effects,  $\boldsymbol{\theta}$  is the parameter vector for the fixed-effects and  $\boldsymbol{b}$  is a random vector of blocking factors. This is not the only model that can be used with this type of design, but it will serve the purpose of this paper.

The remainder of the paper is organized as follows. In section 2 we present a result that connects estimators/predictions that are obtained from mixed-effects models, penalized least-squares estimation and Bayesian formulations. We also address the issue of estimation of the variance components and smoothing parameters that arise from their respective contexts. In this latter respect, we establish that GCV, REML/GML and UBR can all be used to obtain the above mentioned estimators. Section 3 illustrates the implementation of our main result using the three examples mentioned in this section. Section 4 concludes with some comments about the use of the theorems in section 2 and the employment of the Kalman filter algorithm.

2. Equivalence Theorem

To begin, we will give a detailed description of the three modeling scenarios that are the focus of this section.

• Mixed-effects model : Consider first the linear mixed-effects model

(7) 
$$y = \mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b} + \boldsymbol{e},$$

where  $\boldsymbol{y}$  is a  $n \times 1$  vector of responses and  $\mathbf{T}$  and  $\mathbf{U}$  are design matrices for the fixed and random effects of dimensions  $n \times m$  and  $n \times q$ , respectively. Here, we take  $\boldsymbol{\theta}$  to be a  $m \times 1$  vector of fixed effects and  $\boldsymbol{b}$  to be a  $q \times 1$  normally distributed random vector with zero mean and variance-covariance matrix  $\operatorname{Var}(\boldsymbol{b}) = \sigma_b^2 \mathbf{R}$ . The random effects  $\boldsymbol{b}$  are assumed to be independent of the  $n \times 1$  vector of random errors,  $\boldsymbol{e}$ , which in turn, is assumed to be normally distributed with zero mean and variance-covariance matrix  $\sigma_e^2 \mathbf{I}$ . For this model, as well as for the Bayesian model below, the parameters  $\sigma_e^2$  and  $\sigma_b^2$  are the so called variance components. It is often convenient to reparameterized the variance components as  $\lambda = \sigma_b^2/\sigma_e^2$  so that  $\operatorname{Var}(\boldsymbol{y}) = \sigma_e^2(\lambda \mathbf{URU}^{\mathrm{T}} + \mathbf{I})$ . The value of  $\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}$  can be predicted using its BLUP.

• Bayesian Model: Similar to the previous case, in this setting we assume that

(8)  $\boldsymbol{y} = \mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b} + \boldsymbol{e},$ 

46 with **T** and **U** fixed matrices. However, we now also take  $\boldsymbol{\theta}$  to be random and model 47 it as being independent of  $\boldsymbol{b}$  and  $\boldsymbol{e}$ , with a zero mean, normal prior distribution 48 having variance-covariance matrix  $\operatorname{Var}(\boldsymbol{\theta}) = \nu \mathbf{I}$ . The vector of random effects,  $\boldsymbol{b}$ , 49 is also assumed to be normally distributed with zero mean and  $\operatorname{Var}(\boldsymbol{b}) = \sigma_b^2 \mathbf{R}$ . 50 Prediction of  $\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}$  can be accomplished via the use of its posterior mean. In 51 the absence of an informative prior for  $\boldsymbol{\theta}$  a diffuse formulation can be employed

wherein  $\nu$  is allowed to diverge. Note: notice that this is not truly a Bayesian model since there are no priors on the variance components. It is named Bayesian model for the sake of identification. • Penalized Least-Squares: In this case we have  $y = T\theta + Ub + e$  with  $\theta$  and b being non random and e is a vector of zero mean, normally distributed random errors with variance-covariance matrix  $\operatorname{Var}(e) = \sigma_e^2 \mathbf{I}$ . The parameters are to be estimated by minimizing the PLS criterion PLS $(\boldsymbol{a}, \boldsymbol{c}) = (\boldsymbol{y} - \mathbf{T}\boldsymbol{a} - \mathbf{U}\boldsymbol{c})^T (\boldsymbol{y} - \mathbf{T}\boldsymbol{a} - \mathbf{U}\boldsymbol{c}) + \lambda \boldsymbol{c}^T \mathbf{R}^{-1} \boldsymbol{c},$ (9)with respect to a and c. Here,  $\mathbf{R}^{-1}$  is a penalty matrix and  $\lambda$  is the parameter that controls how heavily we penalize the coefficients c. Having these three scenarios in mind, we now state the following theorem. **Theorem 2.1.** The Best Linear Unbiased Predictor (BLUP) of  $\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}$  in ( $\boldsymbol{\gamma}$ ) is given explicitly by  $\hat{y} = \mathbf{A}_{\lambda} y$ (10)where  $\mathbf{A}_{\lambda} = \{\mathbf{I} - \mathbf{Q}^{-1} [\mathbf{I} - \mathbf{T} (\mathbf{T}^{T} \mathbf{Q}^{-1} \mathbf{T})^{-1} \mathbf{T}^{T} \mathbf{Q}^{-1}]\},\$ (11)and  $\mathbf{Q} = (\lambda \mathbf{U} \mathbf{R} \mathbf{U}^T + \mathbf{I}).$ (12)This result is numerically the same as the limiting value (as  $\nu \to \infty$ ) of  $E[\mathbf{T}\boldsymbol{\theta} +$  $\mathbf{Ub}[\mathbf{y}]$  in (8) and the minimizer of (9). *Proof.* To simplify the proof let us assume that the design matrices U and T, as well as **R**, are all full rank matrices (we will later relax this assumption). Under model (7), the first two moments of y are given by  $E(\boldsymbol{y}) = \mathbf{T}\boldsymbol{\theta}$  and  $Var(\boldsymbol{y}) = \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^T + \sigma_e^2 \mathbf{I}.$ Using the distribution of y given b and the distribution of b, we can then find the joint density of  $\boldsymbol{u}$  and  $\boldsymbol{b}$  and obtain the normal equations of [23]:  $\mathbf{T}^{\mathrm{T}}\mathbf{T}\boldsymbol{\theta} + \mathbf{T}^{\mathrm{T}}\mathbf{U}\boldsymbol{b} = \mathbf{T}^{\mathrm{T}}\boldsymbol{u}$  $\mathbf{U}^{\mathrm{T}}\mathbf{T}\boldsymbol{\theta} + (\mathbf{U}^{\mathrm{T}}\mathbf{U} + \mathbf{R}_{\lambda}^{-1})\boldsymbol{b} = \mathbf{U}^{\mathrm{T}}\boldsymbol{y},$ for  $\mathbf{R}_{\lambda} = \lambda \mathbf{R}$ . After some algebra and using the Sherman-Morrison-Woodbury formula in [24] we have  $\mathbf{Q}^{-1} = \mathbf{I} - \mathbf{U}(\mathbf{U}^{\mathrm{T}}\mathbf{U} + \mathbf{R}_{\lambda})^{-1}\mathbf{U}^{\mathrm{T}}$ (13) $\hat{\boldsymbol{\theta}} = (\mathbf{T}^{\mathrm{T}}\mathbf{Q}^{-1}\mathbf{T})^{-1}\mathbf{T}^{\mathrm{T}}\mathbf{Q}^{-1}\boldsymbol{y}$ and  $\hat{\boldsymbol{b}} = (\mathbf{U}^{\mathrm{T}}\mathbf{U} + \mathbf{R}_{\lambda}^{-1})^{-1}\mathbf{U}^{\mathrm{T}}[\mathbf{I} - \mathbf{T}(\mathbf{T}^{\mathrm{T}}\mathbf{Q}^{-1}\mathbf{T})^{-1}\mathbf{T}^{\mathrm{T}}\mathbf{Q}^{-1}]\boldsymbol{y}.$ In this way, the predicted values of  $\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}$  are given by  $\hat{\boldsymbol{y}} = \{ \mathbf{I} - \mathbf{Q}^{-1} [\mathbf{I} - \mathbf{T} (\mathbf{T}^{\mathrm{T}} \mathbf{Q}^{-1} \mathbf{T})^{-1} \mathbf{T}^{\mathrm{T}} \mathbf{Q}^{-1}] \} \boldsymbol{y}.$ (14)imsart-coll ver. 2008/08/29 file: Munoz.tex date: March 25, 2009

З

To show that minimization of the PLS criterion produces the same numerical answer as the BLUP of (7), we differentiate PLS(a, c) with respect to a and c to obtain normal equations which together with (13) give us the same answer as in (14).

It remains to show that under the Bayesian model with diffuse prior,  $\lim_{n\to\infty} E(\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}|\boldsymbol{y})$  also agrees with (14). In this case, the joint distribution of  $\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}$  and  $\boldsymbol{y}$  is found to be normal with zero mean vector and variance-covariance matrix given by

$$\left( \begin{array}{cc} \nu \mathbf{T} \mathbf{T}^{\mathrm{T}} + n^{-1} \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^{\mathrm{T}} & \nu \mathbf{T} \mathbf{T}^{\mathrm{T}} + n^{-1} \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^{\mathrm{T}} \\ (\nu \mathbf{T} \mathbf{T}^{\mathrm{T}} + n^{-1} \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^{\mathrm{T}})^{\mathrm{T}} & \nu \mathbf{T} \mathbf{T}^{\mathrm{T}} + n^{-1} \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^{\mathrm{T}} + \sigma_e^2 \mathbf{I} \end{array} \right).$$

Standard multivariate analysis results then produce

$$E(\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}|\boldsymbol{y}) = Cov(\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}, \boldsymbol{y})[Var(\boldsymbol{y})]^{-1}\boldsymbol{y}$$
14

$$= (\nu \mathbf{T} \mathbf{T}^{\mathrm{T}} + n^{-1} \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^{\mathrm{T}})$$

$$\times (\nu \mathbf{T} \mathbf{T}^{\mathrm{T}} + n^{-1} \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^{\mathrm{T}} + \sigma_e^2 \mathbf{I})^{-1} \boldsymbol{y}.$$
<sup>17</sup>

$$\langle (\nu \mathbf{T} \mathbf{T}^{\mathrm{T}} + n^{-1} \sigma_b^2 \mathbf{U} \mathbf{R} \mathbf{U}^{\mathrm{T}} + \sigma_e^2 \mathbf{I})^{-1} \boldsymbol{y}.$$

Letting  $\lambda$  be as in (12),  $\eta = \nu / \sigma_e^2$  and recalling equation (12) we obtain

(15) 
$$E(\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}|\boldsymbol{y}) = (\eta \mathbf{T}\mathbf{T}^{\mathrm{T}} + \mathbf{U}\mathbf{R}_{\lambda}\mathbf{U}^{\mathrm{T}})(\eta \mathbf{T}\mathbf{T}^{\mathrm{T}} + \mathbf{Q})^{-1}\boldsymbol{y}.$$

Applying the Sherman-Morrison-Woodbury formula [24] on  $(\eta \mathbf{T}\mathbf{T}^{\mathrm{T}} + \mathbf{Q})^{-1}$  and using a little algebra we get

$$\begin{aligned} (\eta \mathbf{T} \mathbf{T}^{\mathrm{T}} + \mathbf{Q})^{-1} &= \mathbf{Q}^{-1} - \\ & \mathbf{Q}^{-1} \mathbf{T} (\mathbf{T}^{\mathrm{T}} \mathbf{Q}^{-1} \mathbf{T})^{-1} [\eta^{-1} (\mathbf{T}^{\mathrm{T}} \mathbf{Q}^{-1} \mathbf{T})^{-1} + \mathbf{I}]^{-1} \mathbf{T}^{\mathrm{T}} \mathbf{Q}^{-1}. \end{aligned}$$

For  $\eta$  sufficiently large, the eigenvalues of  $(\eta^{-1}(\mathbf{T}^{\mathrm{T}}\mathbf{Q}^{-1}\mathbf{T})^{-1})$  are all less than one. So, applying a power series expansion on  $(\eta \mathbf{T}\mathbf{T}^{\mathrm{T}} + \mathbf{Q})^{-1}$  [16], substituting this expansion in (15), and with the aid of some straight forward calculus we have that  $\lim_{n\to\infty} \mathbb{E}(\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}|\boldsymbol{y})$  is exactly the same expression as in (14).

Now, let us go back to our assumption of U, T and R being full rank matrices. This may not be always the case. For example, if we approach estimation from the PLS criterion perspective, there are cases (such as spline smoothing), where  $\mathbf{R}$  has less than full rank. To deal with this instance, suppose that the matrix  $\mathbf{U}\mathbf{R}\mathbf{U}^{T}$  is not invertible. In this situation, the matrix  $\mathbf{Q} = (\lambda \mathbf{U} \mathbf{R} \mathbf{U}^T + \mathbf{I})$  will still be invertible and our only concern is that the matrix  $\mathbf{T}$  is less than full rank. In that case, we can employ conditional inverses [e.g., 18, pp. 31] and the theorem will still hold.

A result such as theorem 2.1 is important because, as pointed out by [4, 19, 48]and [36], one can take advantage of existing methodology and software to facilitate and enhance our analyses. The difference here is that theorem 2.1 is not restricted to the smoothing spline case of [43]; the BLUP result by [40] and referenced by [4, 19, 48]; or to the Bayesian mixed model of [29]. Instead we see that, quite generally, methodology from any particular one of the three frameworks can be potentially applied to obtain useful developments for the other two.

In each of the scenarios described by theorem 2.1, it will generally be necessary to estimate the parameter  $\lambda$ . The following result is a generalization of theorem 5.6 in [12] that allows us to apply three standard methods to the problem of adaptively selecting this parameter. The methods considered here are GCV, UBR and GML 

З

which respectively produce estimators of  $\lambda$  via minimization of

<sup>2</sup>  
<sup>3</sup>  
<sub>4</sub> (16) 
$$\operatorname{GCV}(\lambda) = \frac{n^{-1} \operatorname{RSS}(\lambda)}{[n^{-1} \operatorname{tr}(\mathbf{I} - \mathbf{A}_{\lambda})]^2},$$

(17) 
$$UBR(\lambda) = n^{-1}RSS(\lambda) + 2n^{-1}\sigma_e^2 tr(\mathbf{A}_{\lambda}),$$

and

(18) 
$$\operatorname{GML}(\lambda) = \frac{\boldsymbol{y}^T (\mathbf{I} - \mathbf{A}_{\lambda}) \boldsymbol{y}}{|\mathbf{I} - \mathbf{A}_{\lambda}|_+^{1/(n-m)}}.$$

Here, tr denotes the trace of a matrix,  $\text{RSS}(\lambda) = (\boldsymbol{y} - \hat{\boldsymbol{y}})^T (\boldsymbol{y} - \hat{\boldsymbol{y}})$  and  $|\mathbf{I} - \mathbf{A}_{\lambda}|_+$  is the product of the nonzero eigenvalues of  $\mathbf{I} - \mathbf{A}_{\lambda}$ .

We note in passing that GML is equivalent to the method of REMS that is a popular approach to variance component estimation. [See, e.g. 40]. In terms of the relationship between criteria (16)-(18) we can establish the following result.

**Theorem 2.2.** E[GCV(
$$\lambda$$
)], E[UBR( $\lambda$ )] and E[REMS/GML( $\lambda$ )] are all minimized at  $\lambda = \sigma_b^2 / \sigma_e^2$ .

*Proof.* To establish theorem 2.2 first note that the arguments in [12, pp. 244–247] can be easily modified to account for either the GCV or UBR part of the theorem. The main difference is that here we are not working with diffuse priors. Thus, we will concentrate on sketching the part of the proof that pertains to equation (18). Let  $\lambda_o = \sigma_b^2/\sigma_e^2$  and write  $\mathbf{I} - \mathbf{A}_\lambda = \mathbf{B}(\mathbf{B}^T\mathbf{Q}\mathbf{B})^{-1}\mathbf{B}^T$ , for a **B** such that  $\mathbf{B}^T\mathbf{B} = \mathbf{I}$ ,  $\mathbf{B}\mathbf{B}^T = \mathbf{I} - \mathbf{T}(\mathbf{T}^T\mathbf{Q}^{-1}\mathbf{T})^{-1}\mathbf{T}^T\mathbf{Q}^{-1}$  and  $\mathbf{B}^T\mathbf{T} = 0$ . Then,

Define the matrix of eigenvalues for 
$$\mathbf{B}^{\mathrm{T}}\mathbf{U}\mathbf{R}\mathbf{U}^{\mathrm{T}}\mathbf{B}$$
 with corresponding matrix of eigenvectors  $\mathbf{V}$  as  $\mathbf{\Lambda} = \text{diag}\{d_1, \ldots, d_{n-m}\}$ . Then, we can write

$$\mathbf{B}^{\mathrm{T}}\mathbf{Q}\mathbf{B} = \mathbf{V}(\lambda\mathbf{\Lambda} + \mathbf{I})\mathbf{V}^{\mathrm{T}}.$$

Now, taking expectation with respect to e and b we can show that

$$\mathbf{E}[\mathbf{REMS/GML}(\lambda)] = \frac{\sigma_e^2 \mathrm{tr}[(\mathbf{I} - \mathbf{A}_{\lambda})] + \lambda_o \mathrm{tr}[(\mathbf{I} - \mathbf{A}_{\lambda})(\mathbf{Q} - \mathbf{I})]}{[\prod_{i=1}^{n-m} (\lambda d_i + 1)^{-1/(n-m)}]},$$

$$2$$
  $n-m$  ()  $l + 1$ )

$$= \frac{\sigma_e^2}{\prod_{i=1}^{n-m} (\lambda d_i + 1)^{-1/(n-m)}} \sum_{i=1}^{n-m} \frac{(\lambda_o d_i + 1)}{(\lambda d_i + 1)}.$$

Now, take the difference of the logarithms of the expectations  $E[\text{REMS/GML}(\lambda)]$ and  $E[\text{REMS/GML}(\lambda_o)]$ . A sufficient condition for minimization of the REMS/GML criterion at  $\lambda_o$  is then seen to be

$$\log\left[\frac{1}{n-m}\sum_{i=1}^{n-m}\frac{(\lambda_o d_i+1)}{(\lambda d_i+1)}\right] - \frac{1}{(n-m)}\sum_{i=1}^{n-m}\log\left[\frac{(\lambda_o d_i+1)}{\lambda d_i+1}\right] \ge 0.$$

However, this is an immediate consequence of Jensen's inequality.

Criteria (16)-(18) have long been used for the selection of smoothing or penalty parameters. Golub et al. [17] proposed (16) as a method to choose the ridge regression parameter in a standard regression model like (3) and Craven and Wahba

З

[6] introduced GCV as a method for choosing the smoothing parameter in nonparametric regression. Wahba [46], Kohn et al. [27] and Stein [41] compared the performance of GCV and REML/GML for the smoothing spline case.

Unlike the methods of REML/GML in the PLS framework, GCV and UBR have not been applied in the context of mixed-effects models. Theorem 2.2 suggests that GCV may be another suitable method for estimation of variance components in this context. The fact that the GCV estimator of the variance components shares the REML/GML estimator attribute of minimizing the expectation of the risk, seems to indicate that both estimators will have similar properties and behavior under the mixed-effects model (as it has been shown for the PLS and the Bayesian models [see 27, 46] ). However, this needs to be confirmed by studying the distributional and consistency properties of the GCV estimator of  $\sigma_e^2$  and  $\sigma_b^2$  under the mixed-effects model and this is a topic for future research.

# 3. Examples

З

In this section we focus on the examples introduced in section 1 and exemplify the advantages of using existing methodology for one particular framework (the Bayesian model) to the other two. In particular, we will use a Kalman filter al-gorithm to compute estimators and predictions that arise in the three scenarios considered in theorem 2.1. Perhaps the most common application of the Kalman filter has been in a Bayesian context [see 3, 28]. Specifically, Kohn and Ansley [25], using Wahba's Gaussian model (a particular case of our Bayesian model), reformu-lated the model into a state-space representation and thereby obtained an efficient O(n) algorithm for computing smoothing spline estimators. Theorem 2.1 allows us to extend this approach to non spline smoothing situations and obtain an efficient, Kalman filter based, computational algorithm provided that the random compo-nents in theorem 2.2 admit a state-space representation. This algorithm also permits the evaluation of likelihood functions, making it possible to obtain REMS/GML es-timators for variance components or smoothing parameters.

Description of the Kalman filter is beyond the scope of this paper. Instead, we will focus on establishing a state-space representation for the three examples and refer the reader to [11, 13] and [14] for a more complete development. To accomplish this, it suffices to give only a brief discussion concerning the form of a state-space model.

Any response  $y_i$  can be represented using a state-space model if the observation at time *i* can be expressed as a function of the observation at time i - 1. More formally, a state-space model is composed of a set of response equations

(19) 
$$y_i = \mathbf{h}^{\mathrm{T}}(t_i)\boldsymbol{x}(t_i) + e_i,$$

and a system of state equations

$$x(t_{i+1}) = \mathbf{F}(t_i) \mathbf{x}(t_i) + \mathbf{u}(t_i).$$

with  $t_i \in [0, 1]$  and  $0 = t_0 \leq t_1 < \cdots < t_n$ . The  $y_i$  are observed quantities and the  $e_i, u(t_i), x(t_i)$ , are all unobservable with  $u(t_0) \ldots, u(t_{n-1}), e_1, \ldots, e_n$  and the initial state,  $x(t_0)$ , all being zero mean, uncorrelated normal random variables. In general, the  $x(t_i)$  and  $u(t_i)$  may be vector valued with  $u(t_i)$  having variancecovariance matrix  $\mathbf{R}_{u(t_i)}$ . For our purposes we will treat the vectors  $\mathbf{h}(t_i)$  and the transition matrix  $\mathbf{F}(t_i)$  in (19)– (20) as being known.

We will proceed now to demonstrate the application of the equivalence theoremin the context of our three examples.

# 3.1. Varying Coefficient Models

To illustrate the varying coefficient case, we will examine the progesterone profiles data (figure 1) of [4]. The data consists of metabolite progesterone profiles, measured daily in urine over the course of a menstrual cycle in a group of 51 women.



FIG 1. Observed progesterone measurements for subject 11 in the non-conceptive group. The plots correspond to three of the four cycles for subject 11 and show the log progesterone concentration versus day in the cycle. All cycles have missing observations. Days corresponding to the menses were excluded.

The women in the study were divided into two groups: 29 in the non-conceptive group and 22 in the conceptive group. Each woman contributed a different number of cycles, ranging from 1 to 5 cycles and some of the cycles have missing values.

The goal of the analysis is to detect differences between the conceptive and nonconceptive group profiles. To do this we will express the varying coefficient model (1)with the formulation in (9), apply theorem 2.1 and find the equivalent formulation (8) in the Bayesian framework in order to use the efficient Kalman filter algorithm of [13].

For simplicity, assume that we have complete data and the same number of cycles per woman (later we will relax these assumptions). Let the log progesterone level of the  $c^{th}$  cycle for the  $w^{th}$  woman at time  $t_i$  be denoted by  $y_{wci}$  and model this response as

$$y_{wci} = \beta_1(t_i)X_{1wci} + \beta_2(t_i)X_{2wci} + e_{wci}$$

where i = 1, ..., 24, and  $t_1 = -8, t_2 = -7, ..., t_{24} = 15$  are the days in a menstrual cycle. The cycles c range from 1 to 5 and w = 1, ..., 29 correspond to women in the non-conceptive group and the rest belong to the conceptive group. 

Assume that the  $\beta_k(\cdot)$ 's, k = 1, 2, are smooth functions of t. Usually, this translates into assuming that the functions belong to a Hilbert space of order m [see 

З

22]. To find the estimated profiles we minimize a particular PLS criterion, where the penalty is applied to the integral of the square of the second derivative of the  $\beta_k(\cdot)$ 's. The minimizers of  $\beta_1(\cdot)$ ,  $\beta_2(\cdot)$  are natural splines of order m, with m = 3, that can be represented by a linear combinations of basis functions

$$\sum_{q=0}^{m-1} \theta_{kq} t_i^q + \sum_{r=1}^{24} b_{kr} \xi_r(t_i),$$

$$\int_{0}^{0} b_{kq} c_{i} + \sum_{r=1}^{0} b_{kr} \varsigma_{r} (c_{i}),$$

with knots at each of the design points  $t_i$ , and

(21) 
$$\xi_r(t_i) = \int_0^{\min\{t_r, t_i\}} \frac{(t_i - u)^{m-1}(t_r - u)^{m-1} du}{[(m-1)!]^2}.$$

Equation (21) is one of the usual reproducing kernels of a Hilbert space of order m [22].

Let  $\boldsymbol{y}_{wc} = [y_{wc}(t_1), \ldots, y_{wc}(t_{24})]^{\mathrm{T}}$  be the vector of responses for women w that contains all the daily observations in the c cycle, and  $\boldsymbol{\xi}_i = [\xi_1(t_i), \xi_2(t_i), \ldots, \xi_{24}(t_i)]^{\mathrm{T}}$  and  $\boldsymbol{t}_i = [t_i^0, t_i^1, \ldots, t_i^{m-1}]^{\mathrm{T}}$  be the vectors of basis functions evaluated at the times  $t_i$ 's.

Denote the vector of coefficients for  $\beta_1$  and  $\beta_2$  as  $\boldsymbol{\theta}_1 = [\theta_{10}, \theta_{11}, \dots, \theta_{1,(m-1)}]^{\mathrm{T}}$ ,  $\boldsymbol{\theta}_2 = [\theta_{20}, \theta_{21}, \dots, \theta_{2,(m-1)}]^{\mathrm{T}}$ ,  $\boldsymbol{b}_1 = [b_{10}, b_{11}, \dots, b_{1,24}]^{\mathrm{T}}$ ,  $\boldsymbol{b}_2 = [b_{21}, b_{22}, \dots, b_{2,24}]^{\mathrm{T}}$ , respectively. Construct  $\boldsymbol{t}, \boldsymbol{\xi}$  and  $\boldsymbol{X}$  such that

$$\boldsymbol{t} = \begin{bmatrix} \boldsymbol{t}_{1}^{\mathrm{T}} \\ \boldsymbol{t}_{2}^{\mathrm{T}} \\ \vdots \\ \boldsymbol{t}_{24}^{\mathrm{T}} \end{bmatrix}, \, \boldsymbol{\xi} = \begin{bmatrix} \boldsymbol{\xi}_{1}^{\mathrm{T}} \\ \boldsymbol{\xi}_{2}^{\mathrm{T}} \\ \vdots \\ \boldsymbol{\xi}_{24}^{\mathrm{T}} \end{bmatrix} \text{ and } \boldsymbol{X} = \begin{bmatrix} \boldsymbol{X}_{1wc1} & \boldsymbol{X}_{2wc1} \\ \boldsymbol{X}_{1wc2} & \boldsymbol{X}_{2wc2} \\ \vdots & \vdots \\ \boldsymbol{X}_{1wc24} & \boldsymbol{X}_{2wc24} \end{bmatrix}.$$

Let  $\mathbf{T}_{wc} = \mathbf{t} \bigotimes \mathbf{X}$  and  $\mathbf{U}_{wc} = \mathbf{\xi} \bigotimes \mathbf{X}$ , where  $\mathbf{A} \bigotimes \mathbf{B}$  denotes the Kronecker product of the matrices  $\mathbf{A}$  and  $\mathbf{B}$  and it is equal to  $a_{ij}\mathbf{B}$ .

For each woman's cycle we have the model  $\mathbf{T}_{wc}\boldsymbol{\theta}^* + \mathbf{U}_{wc}\boldsymbol{b}^* + \boldsymbol{e}_{ew}$ , where  $\boldsymbol{\theta}^* = [\boldsymbol{\theta}_1^{\mathrm{T}}, \boldsymbol{\theta}_2^{\mathrm{T}}]^{\mathrm{T}}$ ,  $\boldsymbol{b}^* = [\boldsymbol{b}_1^{\mathrm{T}}, \boldsymbol{b}_2^{\mathrm{T}}]^{\mathrm{T}}$  and  $\boldsymbol{e}_{wc}$  is the corresponding vector of errors. Denote by  $\boldsymbol{y}$  and  $\boldsymbol{e}$  the vectors resulting from stacking the vectors  $\boldsymbol{y}_{wc}$  and  $\boldsymbol{e}_{wc}$ , (i.e.,  $\boldsymbol{y} = [\boldsymbol{y}_{1,1}^{\mathrm{T}}, \boldsymbol{y}_{1,2}^{\mathrm{T}}, \dots, \boldsymbol{y}_{1,5}^{\mathrm{T}}, \boldsymbol{y}_{2,1}^{\mathrm{T}}, \dots, \boldsymbol{y}_{51,5}^{\mathrm{T}}]^{\mathrm{T}}$ ), and let  $\mathbf{T} = \text{diag}\{\mathbf{T}_{wc}\}_{\substack{w=1,51\\c=1,5}}^{w=1,51}$  and  $\mathbf{U} = \text{diag}\{\mathbf{U}_{wc}\}_{\substack{w=1,51\\c=1,5}}^{w=1,51}$ . Then, we can construct the model  $\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b} + \boldsymbol{e}$ , where  $\boldsymbol{\theta} = \mathbf{1} \otimes \boldsymbol{\theta}^*$ ,  $\boldsymbol{b} = \mathbf{1} \otimes \boldsymbol{b}^*$ , and minimize criteria (9), where  $\mathbf{R}^{-1} = \mathbf{U}$ .

By theorem 2.1, this is equivalent to find  $\lim_{\substack{\substack{\nu \to \infty \\ lim \nu \to \infty \\ e \ }} \mathbb{E}[\mathbf{T}\theta + \mathbf{U}b|\mathbf{y}]$ , where  $\theta$ , band e are independent of each other and normally distributed with zero mean and variance-covariance matrices  $\nu \mathbf{I}$ ,  $\sigma_b^2 \mathbf{U}^{-1}$ , and  $\sigma_e^2 \mathbf{I}$ , respectively. In this case, the smoothing parameter  $\lambda$  in the PLS model can be found using the variance components,  $\sigma_e^2$  and  $\sigma_b^2$  since  $n\lambda = \sigma_b^2/\sigma_e^2$ , where n is the total number of observations in the data.

The equivalent Bayesian representation of the varying coefficient model will allow us to make use of the Bayesian theory and apply it to our PLS setting. Specifically, we can follow [13] and transform the Bayesian model into a state-space model, as they indicate, and apply their efficient algorithm to compute the varying coefficients and respective confidence bands. Their approach also shows how to reformulate the matrices in the Bayesian model so the unbalanced design does not represent a prob-lem in the computation of the estimators. For details on how to find the state-space model form, or on how to apply this efficient algorithm, we refer the readers to the 

appendix and to the above mentioned authors, respectively. To see what are the advantages of using this equivalence representation of the PLS, let us first explore the extent that the Kalman filter can speed up computations. To investigate this issue we carried out a run time comparison between our Kalman filter approach, the "standard way" of estimation assuming a mixed-effects model approach (both in SAS), and, only as a reference, we provide the time used in the method developed by Brumback and Rice [4]. We need to point out that these are the reported times in their 1998 paper and that there has been great improvement in computational speed since the publication of this paper. Table 1 shows the required times for computing the estimated conceptive and non-conceptive functions (see figure 2).





#### Log Hormone Profiles with respective "Confidence" Interva



FIG 2. Smooth estimates for non conceptive and conceptive mean groups with respective 95% pointwise confidence intervals. The corresponding smoothing parameters were computed using the GML method implemented through a Kalman filter algorithm.

The first time in table 1 corresponds to the time employed by the Kalman filter algorithm of [11] implemented in SAS and using a computer with a 3.2GHz processor and 1G RAM. This algorithm used 2004 observations (missing values were

imsart-coll ver. 2008/08/29 file: Munoz.tex date: March 25, 2009

З

З

omitted) and calculated the estimated coefficient functions and corresponding 95% confidence intervals. The second time is the result of using a mixed-effects model representation and taking advantage of SAS proc mixed (the same equipment was used). The last time is the one reported by [4]. They implemented an eigenvalueeigenvector decomposition on a mixed-effects model representation of the profiles, separately for each group, and combined the times for both groups and the estimation of the variance components. We calculated the smoothing parameters via REML/GML using the Kalman filter algorithm and it took approximately 10.5 seconds in SAS (this time is included in the computation of the Kalman filter in table 1). These parameters were used in both the SAS and Kalman filter calculation of the varying coefficient functions (we didn't want to calculate the smoothing parameters with SAS Proc Mixed given that it already took 4 hrs. to calculate the profiles without estimating the variance components). The reason why SAS takes so long to estimate the functions is due to the complex covariance structure of the model and the number of observations. The convenient SAS built-in covariance structures were not an option [see comment by 4], and the inversion of a general  $n \times n$  matrix requires  $O(n^3)$  operations versus the O(n) used by the Kalman filter.

Another advantage of using the Bayesian interpretation in our PLS model is that the Kalman filter allows us to easily obtain confidence intervals as well as point estimators. In this respect, we use the relationship between PLS and the Bayesian model to provide Bayesian  $100(1 - \alpha)\%$  confidence (or prediction) intervals which parallel those developed by [44] and [31]. Specifically, we "estimate" the  $i^{th}$  component of  $\beta_k(t_i)$  via the interval  $\beta_k(t_i) \pm z_{1-\alpha/2}\sqrt{\hat{\sigma}_e^2 \times a_{ii}}$ , where  $\hat{\sigma}_e^2 = [(\boldsymbol{y} - \mathbf{A}_\lambda \boldsymbol{y})^T (\boldsymbol{y} - \mathbf{A}_\lambda \boldsymbol{y})]/(n - m)$ ,  $a_{ii}$  is the  $i^{th}$  diagonal element of the corresponding hat matrix  $\mathbf{A}_\lambda$  for  $\beta_k$  and  $z_{1-\alpha/2}$  is the  $100(1 - \alpha/2)$  standard normal percentile.

Wahba's "Bayesian Confidence Intervals" have been often used in the nonparametric community. Wahba [44] showed that the average of the coverage probability across points of these pointwise intervals is very close to nominal level for large n. She also commented that even if the confidence intervals are derived from a Bayesian perspective, they perform well in the frequentist realm. Nychka [32] offers an excellent discussion on why this is true.

In their paper, Brumback and Rice [4] utilized a hierarchical bootstrap method to assess the variability of the fitted functions instead of using the variance components estimators (it is well know that these estimators often underestimate the true parameters). For each bootstrapped sample 1.5 hours was required to obtain the estimated sample profiles (as reported by [4]). As a result, a partially parametric version of the method was implemented [see 4, for more details]). They computed 35 bootstrap samples and this took approximately 45 mins. In contrast, the confidence intervals computed in this paper for the progesterone profiles were obtained with the same computational effort involved in the estimation of the profiles. 

Our estimated function profiles seem to agree with the ones obtained by Brum-back and Rice. In addition, the "confidence" intervals also allow us to see that, on average, the production of progesterone in the conceptive group drops significantly from day 4 until around day 8 (when ovulation occurs) as compared to the hormone production of the non conceptive group. This result differs from the findings by [4]. Their bootstrap sample suggested that the decrease in progesterone for the con-ceptive group was not significant. The discrepancy between our findings and those of Brumback and Rice may be due to the small bootstrap sample they employed in their analysis or with our interpretation of the confidence intervals. Nychka [32] 

imsart-coll ver. 2008/08/29 file: Munoz.tex date: March 25, 2009

pointed out that these intervals may not be reliable at specific points, even more if those points are part of a sharp peak or deep valley in the function. However, he also mentioned that it provides "a reasonable measure of the spline estimate's accuracy provided that the point for evaluation is chosen independently of the shape" of the function. It is known that the women are more fertile around day 3 to day 9, making it an interval targeted for studying before the start of the analysis. Also, we do not consider that the bump that forms in the interval of interest is that sharp. Hence, we believe that the confidence intervals provide reasonable evidence that the profiles are different.

### 3.2. Ridge Regression and Penalized Spline Regression

To exemplify the use of theorem 2.1 in the ridge regression setting, we have selected a data set that is freely distributed by *StatLib* at http://lib.stat.cmu.edu. The data set consists of 1150 heights measured at 1 micron intervals along the drum of a roller (i.e. parallel to the axis of the roller). The units of height are not given and the zero reference height is arbitrary.

To fit this data we used a model of the form (3) with  $\mathbf{X} = [\mathbf{T}, \mathbf{U}]$  and corresponding vector of coefficients  $\boldsymbol{\beta} = [\boldsymbol{\theta}^{\mathrm{T}}, \boldsymbol{b}^{\mathrm{T}}]^{\mathrm{T}}$ , where

(22) 
$$\mathbf{T} = \begin{bmatrix} 1 & t_1 \\ 1 & t_2 \\ \vdots & \vdots \\ 1 & t_{150} \end{bmatrix} \text{ and } \mathbf{U} = \begin{bmatrix} (t_1 - \xi_1)_+ & \cdots & (t_1 - \xi_k)_+ \\ (t_2 - \xi_1)_+ & \cdots & (t_2 - \xi_k)_+ \\ \vdots & \cdots & \vdots \\ (t_{1150} - \xi_1)_+ & \cdots & (t_{1150} - \xi_k)_+ \end{bmatrix}$$

The generalized ridge regression estimator of  $\beta$  is then obtained by minimizing the PLS criterion (4), with **K** as in (5) and m = 2.

Applying the results of theorem 2.1 we can write a parallel mixed-effects model representation for this ridge regression problem. This particular framework was considered by [36] who describe in section 4.9 of their book how to represent p-splines as BLUP's and illustrated how to use available software packages, like SAS proc mixed or the S-PLUS function lme, to obtain a fitted curve for the data. In view of the equivalence theorem, an alternative approach would be to use the connection between PLS and the Bayesian model so the Kalman filter can be implemented for purposes of computing estimators and "confidence" intervals. Another compre-hensive description of the use of P-splines in the semi-parametric regression setting using Bayesian techniques is given in [5]. In this paper, we will use the Bayesian connection.

Assume that the vectors  $\boldsymbol{\theta}$ ,  $\boldsymbol{b}$  and  $\boldsymbol{e}$  are independently normally distributed with zero mean and respective variance-covariance matrices  $\nu \mathbf{I}$ ,  $\sigma_b^2 \mathbf{I}$  and  $\sigma_e^2 \mathbf{I}$ . Then, by the equivalence theorem, the minimizer of (4) is the same as the limit, when  $\nu$  is allowed to go to infinity, of the posterior mean of  $\mathbf{T}\boldsymbol{\theta} + \mathbf{U}\boldsymbol{b}|\boldsymbol{y}$ .

Again, this Bayesian model representation of the ridge regression example will
permit the use of the Kalman filter algorithm for the computation of the estimated
function and its respective "confidence intervals". For the explicit form of the statespace model see the appendix.

For this particular example, we considered two different model versions, one using k = 150 knots and the other with k = 1150 knots. This because we wanted to contrast the computational performance of the P-splines versus the computational effort required using smoothing splines and the Kalman filter. We should remark

З
here that, when k = 1150, basically we have a smoothing spline estimator which basis functions are the polynomials and the truncated power functions  $(t_i - \xi_k)_+$ , for i = 1, ..., 1150.



FIG 3. Figure (a) shows the P-spline estimator with 150 knots for the roller height with its respective 95% confidence bands. The corresponding figure for the P-spline estimator with 1150 knots looks exactly the same and has been omitted. Figure (b) shows the comparison between the p-spline estimator with 150 knots and the P-spline with 1150 knots. There is no visual difference and both procedures yielded an estimated error variance of 0.36.

Figure 3 shows the smooth estimated curve for the roller height and corresponding 95% "confidence intervals". The smoothing parameters were, respectively,  $\lambda_{150} = 0.043$  and  $\lambda_{1150} = 0.095$ . They were selected via GCV and as we can see the GCV methods adjusts the smoothing parameters according to the number of knots used.

One of the main arguments in favor of using P-splines in lieu of smoothing splines is that, by reducing the number of knots involved in the model, we increase the computational efficiency involved in calculating the spline estimator. This is true when using brute force methods, i.e., direct inversion of the matrix (12). However, when using the proposed Kalman filter algorithm, the computational advantage of the P-splines over the smoothing splines disappear as we can see in table 2.

| functions the polynom<br>and the tim | ials and equation (21)). Bother the second second second terms of the second seco | h estimators were computed using coo<br>tion of the smoothing parameter. | le in R |
|--------------------------------------|--|--|---------|
|                                      | Knots  | Real Time  |         |
|                                      | P-spline 150   | 48.34 secs.  |         |
|                                      | P-spline 1150  | 54.55 secs.  |         |
|                                      | Smoothing Spline   | 48.26 secs.  |         |
|                                      |  |  |         |
|                                      |  |  |         |
|                                      |  |  |         |
| 3.3. Mixed-Effect                    | ts Model   |  |         |

imsart-coll ver. 2008/08/29 file: Munoz.tex date: March 25, 2009

of a mixed-effects model representation we will demonstrate the use of Kalman filtering for estimating parameters in a setting that it has seldom being used and that it can benefit from the reduced computational burden of estimating parameters and variance components. It is true that, if we have a "reasonable" number of observations and a specific covariance structure, like the ones provided by existing software, it will be advisable to use these procedures in lieu of the Kalman filter. However, there are occasions where the number of observations is really large. Then we can take advantage of the computational efficiency of the Kalman filter.

Our example deals with a randomized block design, where the data consists of 37 patients, which represent the random blocks, and a set of consecutive Hamilton depression scores measured over the course of 6 weeks (see figure 4). The data set is part of a study conducted by [34] and it is available at http://tigger.uic.edu/ hedeker/. We model the data as

$$y_{ij} = \beta_0 + \beta_1 \text{week} + b_i + e_{ij},$$

where the  $y_{ij}$ 's are the depression scores, for i = 1, ..., 37,  $\beta_0$  and  $\beta_1$  are fixed parameters and week = 0, 1..., 5, is the week number where the score was measured. The random effects due to each patient are denoted by the  $b_i$ 's and they are independent of the errors  $e_{ij}$ 's which are generated by an autoregressive process of order 1, i.e.,  $e_{ij} = \phi e_{i,j-1} + a_j(t_i)$ , with  $\phi$  a constant and  $a_j(t_i)$  independent, identically distributed zero mean errors with variance  $\sigma_e^2$ .



FIG 4. (a) Hamilton depression scores for 37 patients measured over the period of 6 weeks. (b) Estimated regression line, y = 24.41 - 2.36 week, with respective 95% confidence bands.

Let  $\boldsymbol{y}$  be the vector of depression scores such that  $\boldsymbol{y} = [\boldsymbol{y}_1^{\mathrm{T}}, \dots, \boldsymbol{y}_{37}^{\mathrm{T}}]^{\mathrm{T}}$ , for  $\boldsymbol{y}_i = [y_{i0}, y_{i1}, \dots, y_{i5}]^{\mathrm{T}}$ . Denote by  $\mathbf{1}_n$ , the vector of all ones of dimension  $n \times 1$  and  $\mathbf{week} = [0, 1, \dots, 5]^{\mathrm{T}}$ . In matrix form our model becomes

$$oldsymbol{y} \hspace{.1in} = \hspace{.1in} \mathbf{T}oldsymbol{ heta} + oldsymbol{b} + oldsymbol{e},$$

with 
$$\boldsymbol{\theta} = [\beta_0, \beta_1]^{\mathrm{T}}, \mathbf{T} = [\mathbf{1}_{37} \bigotimes \mathbf{1}_5, \mathbf{1}_{37} \bigotimes \mathbf{week}], \text{ and } \boldsymbol{b} = \mathbf{1}_5 \bigotimes [b_1, b_2, \dots, b_{37}]^{\mathrm{T}} \text{ and}$$
  
 $\boldsymbol{e} = [\boldsymbol{e}_1^{\mathrm{T}}, \boldsymbol{e}_2^{\mathrm{T}}, \dots, \boldsymbol{e}_{37}]^{\mathrm{T}}, \text{ for } \boldsymbol{e}_i = [e_{i0}, e_{i1}, \dots, e_{i5}]^{\mathrm{T}}. \text{ Here, we model the } \boldsymbol{b} \text{ as normally}$ 

distributed with zero mean and variance-covariance matrix  $\mathbf{R} = \{\xi_r(t_1)\}_{\substack{r=1,5\\i=1,5}}$  and

$$\frac{\sigma_e^2}{1-\phi^2} \quad \frac{\phi}{1-\phi^2} \quad \frac{\phi^2}{1-\phi^2} \quad \cdots \quad \frac{\phi^n}{1-\phi^2}$$

$$\frac{\overline{\psi}}{1-\phi^2} \quad \frac{\overline{\psi}}{1-\phi^2} \quad \frac{\overline{\psi}}{1-\phi^2} \quad \cdots \quad \frac{\overline{\psi}}{1-\phi^2} \quad ,$$

$$\mathbf{W} = \begin{bmatrix} \frac{\sigma_e^2}{1-\phi^2} & \frac{\phi}{1-\phi^2} & \frac{\phi^2}{1-\phi^2} & \cdots & \frac{\phi^n}{1-\phi^2} \\ \frac{\phi}{1-\phi^2} & \frac{\sigma_e^2}{1-\phi^2} & \frac{\phi}{1-\phi^2} & \cdots & \frac{\phi^{n-1}}{1-\phi^2} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \frac{\phi^n}{1-\phi^2} & \frac{\phi^{n-1}}{1-\phi^2} & \cdots & \frac{\phi}{1-\phi^2} & \frac{\sigma_e^2}{1-\phi^2} \end{bmatrix},$$

З

where W is the variance-covariance matrix of the AR(1) errors.

To find the corresponding Bayesian model, let b and e keep their distributions and assume that  $\theta$  is normally distributed with zero mean and variance-covariance matrix  $\nu \mathbf{I}$ . Once in the Bayesian form, we check that our observations  $Y_{ij}$  can be represented using the state-space equations (19-20). The equivalence theorem hold regardless of the state-space structure but, if we have that structure, then we can apply the efficient Kalman filter algorithm of [13] and estimate all our parameters with linear computational efficiency.

Figure 4 shows the estimated regression line for the Hamilton depression scores over the 6 week period. The variance components for this example are estimated via REML/GML and are  $\hat{\phi} = 0.97$ ,  $\hat{\sigma}_e^2 = 1.214$  and  $\hat{\sigma}_b^2 = 0.00132$ . The corresponding estimated values for the regression coefficients are  $\hat{\theta}_0 = 24.41$  and  $\hat{\theta}_1 = -2.36$ .

#### 4. Summary

In this paper, we have reviewed known results concerning the numerical equivalence of 1) a smoothing spline estimator and a particular mixed-effects model and 2) a smoothing spline estimator and the posterior mean of Wahba's Gaussian model and focus on the more general framework of frequentist and Bayesian mixed-effects models and penalized least-squares estimation as seen in Theorem 2.1. This result broadens the number of methodological resources available for computing BLUPs, posterior means, likelihoods and minimizers of penalized least squares criteria and facilitates the use of existing methodological tools, as exemplified by theorem 2.2and our examples.

The link between the Bayesian mixed-effects model and the two other model settings allowed us to obtain Bayesian "confidence" intervals for the profile groups (instead of the computationally demanding bootstrap method of Brumback and Rice) and facilitated the analysis of the profile differences during the fertile days. Example 2 showed us that the Kalman filter implementation is not restricted to Wahba's Bayesian model. More generally, the idea carries over to settings involving p-splines, Kernel estimators, differences, etc. Lastly, this link allows for the imple-mentation of a computationally efficient Kalman filter algorithm in many cases of interest. Kalman filter algorithms have been used to compute smoothing splines type estimators [19, 25, 26, 48]. But, they have been sparsely used in mixed-effects model settings. To this author knowledge, only [38] and, more recently, [30] have ap-plied the Kalman filter to mixed-effects models. In the mixed-effects framework, the techniques employed for the analysis of large data sets require the use of computer intensive methods like the EM or MCMC algorithms [1, 39], conjugate gradient iterative methods [42], or the use of high performance computing environments. Some of the methods mentioned in these references assume that observations are gener-ated by Brownian motion or ARMA processes and, whenever we have this type of processes, we have a state-space structure that can be exploited, as demonstrated 

in our examples, to reduce the computational burden. Observations generated by longitudinal analysis (as in example 3), repeated measurements or any process that depends on an ordering variable can also frequently be assumed to have a statespace representation and can, as a result, benefit from the computational efficiency of the Kalman filter.

# Appendix A: State-Space Forms

In this section we will explicitly describe the state-space forms used for the application of the Kalman filter in each of our examples. Since the form of the errors  $u(t_i)$  in equation (20) is assumed to be the same for the varying-coefficient case and the mixed-effects model, we will show the derivation for the varying coefficient case and detail the small changes needed for the mixed-effects case. We will leave for the last the ridge regression example.

To employ the Kalman filter for computation of the varying coefficient example we need to show that the varying coefficients have a state-space representation. That is, we need to be able to write equation (1) using equations (19)–(20). Since the  $\beta_k(\cdot)$  are assumed to be smooth functions of t, we model them as

(A.1) 
$$\beta_k(t_i) = \sum_{q=0}^{m-1} \theta_{kq} t^q + \sigma_b^2 Z_k(t_i),$$

for k = 1, 2 and m = 2, where (without loss of generality) we can take t in [0, 1] and

$$Z_k(t) = \int_0^1 \frac{(t-u)_+^{m-1}}{(m-1)!} dW_k(u),$$
26  
27  
28  
28  
28  
28

with  $W_k(\cdot)$  standard Wiener processes. To simplify matters, first assume that  $\beta_k(t_i) = \sigma_b^2 Z_k(t_i)$ . Then,  $\beta_k(t_{i+1})$  can be written as  $\sigma_b^2$  times

$$\int_0^{t_i} \frac{(t_{i+1}-s)_+^{m-1}}{(m-1)!} dW_k(s) + \int_{t_i}^{t_{i+1}} \frac{(t_{i+1}-s)_+^{m-1}}{(m-1)!} dW_k(s).$$

Taking

$$u_k(t_i) = \int_{t_i}^{t_{i+1}} \frac{(t_{i+1} - u)_+^{m-1}}{(m-1)!} dW_k(u),$$
35
36
36
37

for  $t_i < t_j$ , the covariance between  $u_k(t_i)$  and  $u_k(t_j)$  is found to be equal to

$$\int_0^{t_i} \frac{(t_i - u)^{m-1} (t_j - u)^{m-1}}{[(m-1)!]^2} du.$$

For the remaining integral, add and subtract  $t_i$  inside  $(t_{i+1} - u)^{m-1}$  and apply the Binomial theorem. Upon doing this, a state-space representation results with  $\mathbf{F}(t_i)$ equal to

$$\begin{bmatrix} 1 & (t_{i+1} - t_i) & \frac{(t_{i+1} - t_i)^2}{2!} & \dots & \frac{(t_{i+1} - t_i)^{m-1}}{(m-1)!} \end{bmatrix}$$

imsart-coll ver. 2008/08/29 file: Munoz.tex date: March 25, 2009

 $\mathbf{Z}_{k}(t_{i}) = [Z_{k}(t_{i}), Z_{k}^{(1)}(t_{i}), \dots, Z_{k}^{(m-1)}(t_{i})]^{T}, \mathbf{u}_{k}(t_{i}) = [u_{k}(t_{i}), u_{k}^{(1)}(t_{i}), \dots, u_{k}^{(m-1)}(t_{i})]^{T}$ and  $\mathbf{Z}_{k}(t_{i+1}) = \mathbf{F}(t_{i})\mathbf{Z}_{k}(t_{i}) + \mathbf{u}_{k}(t_{i}).$ Now, rearranging the observations with respect to the time  $t_i$  define  $\boldsymbol{y}_{iw}^{\mathrm{T}} = [y_{iw1}, \dots, y_{iwc_w}]^{\mathrm{T}}$ with  $\boldsymbol{y}_{iw}^{\mathrm{T}}$  the responses for woman w at time  $t_i$  observed at cycles  $1, \ldots, c_w$ , with corresponding vector of random errors  $\boldsymbol{e}_{iw}$ . Let  $\boldsymbol{x}(t_i) = [\boldsymbol{Z}_1(t_i), \boldsymbol{Z}_2(t_i)]^T, \boldsymbol{u}(t_i) =$  $[\boldsymbol{u}_1(t_i), \boldsymbol{u}_2(t_i)]^T$  and  $\boldsymbol{X}_{kwi} = [X_{kw1i}, \dots, X_{kwc_wi}]^T$ . Then, taking  $h(t_i) = [X_{1_{avi}}^T, 0, ..., 0, X_{2_{avi}}^T, 0, ..., 0]^T,$ we arrive at the state-space model  $\boldsymbol{y}_{iw} = \boldsymbol{h}(t_i)\boldsymbol{x}(t_i) + \boldsymbol{e}_{iw},$  $\boldsymbol{x}(t_{i+1}) = \mathbf{F}^{\star}(t_i)\boldsymbol{x}(t_i) + \boldsymbol{u}(t_i),$ where  $\mathbf{F}^{\star}(t_i)$  is the block diagonal matrix of size  $2m \times 2m$  with diagonal blocks  $\mathbf{F}(t_i), i = 1, \ldots, n.$ Application of the standard Kalman filter to the vector of observations  $\boldsymbol{y}_{iw}$  will yield coefficient functions estimates that disregard the polynomial term in (A.1). To account for that, we must employ the diffuse Kalman filter as in [13]. This entails a slight modification of our approach wherein the Kalman filter is applied to the vector of observations  $\boldsymbol{y}_{iw}$  and each of the vectors  $\boldsymbol{1}_{\boldsymbol{n}}, \boldsymbol{t}, \boldsymbol{t}^2 \dots, \boldsymbol{t}^{(m-1)}$ , where  $\mathbf{t}^r = [t_1^r, t_2^r, \dots, t_n^r]^{\mathrm{T}}$  [see 11, for a detailed derivation]. For our mixed-effects example we need to show that e can be represented in a state-space form and stack the respective state vectors, errors and matrices. We will proceed as follows: since the errors  $e(t_i)$  are generated by an AR(1) process, they can be written as  $e(t_{i+1}) = \phi e(t_i) + a_j(t_i)$ , with  $\phi$  a non-random coefficient. This entails that the transition matrix  $\mathbf{F}^{\star}(t_i) = diag\{\mathbf{F}(t_i), \phi\}$ , with  $\mathbf{F}(t_i)$  as in (A.2)

and  $\boldsymbol{h}(t_i) = [1, 0, 1]$ . Take the state vector,  $\boldsymbol{x}(t_i)$ , to be equal to  $[\boldsymbol{Z}_k(t_i)^{\mathrm{T}}, e(t_i)]^{\mathrm{T}}$ ,  $\boldsymbol{u}(t_i) = [\boldsymbol{u}_k(t_i)^{\mathrm{T}}, a_j(t_i)]^{\mathrm{T}}$  with m = 2, where  $\boldsymbol{Z}_k(t_i)$  and  $\boldsymbol{u}_k(t_i)$  are as in the varying coefficient case. Specific details about the form of the state vector and the vector  $\boldsymbol{u}(t_i)$  of the state equation (20), as well as a more general form for an ARMA model, can be found in [14].

Lastly, the state-space representation for the ridge regression example is found by taking the state vector to be  $\boldsymbol{x}(t_i) = [x(t_i), x^{(1)}(t_i), \dots, x^{(m-1)}(t_i)]^T$ , with  $x(t_i) = \sum_{k=1}^{j} \beta_k (t_i - \xi_k)^{m-1}$  for  $t_i \in [\xi_j, \xi_{j+1})$  (using the definition of the truncated power function), and  $x^{(r)}(t_i)$  the  $r^{th}$  derivative of  $x(t_i), r = 1, \dots, (m-1)$ . Then,

$$\boldsymbol{x}(t_{i+1}) = \mathbf{F}(t_i)\boldsymbol{x}(t_i) + \boldsymbol{u}(t_i),$$

with 
$$\mathbf{F}(t_i)$$
 as in (A.2) and  $\boldsymbol{u}(t_i) = [u(t_i), u^{(1)}(t_i), \dots, u^{(m-1)}(t_i)]^{\mathrm{T}}$ , where

$$u(t_i) = \begin{cases} 0 & \text{if } t_{i+1} \in [\zeta_j, \zeta_{j+1}) \\ \beta_{j+1}(t_{i+1} - \xi_{j+1})^{m-1} & \text{if } t_{i+1} \in [\xi_{j+1}, \xi_{j+2}) \end{cases}.$$

To complete the state-space formulation, take the vector  $\mathbf{h}(t_i)$  to have dimension  $m \times 1$  with one in the first position and the rest of its elements equal to zero.

З

| 1   | Re   | eferences   | 1  |
|-----|------|---|----|
| 2   |      |   | 2  |
| 3   | [1]  | AITKIN, M. (1999). A general maximum likelihood analysis of variance components in generalized  | 3  |
| 4   | [0]  | linear models. Statistical Methods in Medical Research, 55, 117–128.  | 4  |
| 5   | [2]  | ANDERSON, B. D. O. and MOORE, J. B. (1979). <i>Optimal Filtering</i> . Prentice Hall, Englewood Cliffs,   | 5  |
| 6   | [3]  | ANSLEY, C. F. and KOHN, R. (1985). Estimation, filtering, and smoothing in state space models with  | 6  |
| 7   | . ,  | incompletely specified initial conditions. Annals of Statistics, 13, 1286–1316.   | 7  |
| 8   | [4]  | BRUMBACK, B. A. and RICE, J. A. (1998). Smoothing spline models for the analysis of nested and  | 8  |
| 9   | [5]  | CRAINICEANIL C. M. RUDDERT, D. and WAND, M. P. (2005). Bayesian analysis for populized spline   | 9  |
| 10  | [0]  | regression using WinBugs. Journal of Statistical Software, 14, 1–24.  | 10 |
| 11  | [6]  | CRAVEN, P. and WAHBA, G. (1979). Smoothing noisy data with spline functions. Numerical Mathe-   | 10 |
| 10  | [=1  | matic, <b>31</b> , 377–403.   | 10 |
| 12  | [7]  | CRESSIE, N. (1990). Reply: letters to the editor. <i>The American Statistician</i> , <b>44</b> , 250–258.<br>DEMPSTER A P. LAIRD N.M. and RUBIN D.B. (1977). Maximum likelihood from incomplete data. | 12 |
| 13  | [0]  | via the EM algorithm. Journal of the Royal Statistical Society, Series B, <b>39</b> , 1–22.   | 13 |
| 14  | [9]  | EILERS, P. H. C. and MARX, B. D. (1996). Flexible smoothing with <i>B</i> -splines and penalties. <i>Statistical</i>  | 14 |
| 15  | [10] | Science, 11, 89–102.  | 15 |
| 16  | [10] | and preschool children in the North Central Region of the United States of America. World Review  | 16 |
| 17  |      | of Nutrition and Dietetics, 14, 269–332.  | 17 |
| 18  | [11] | EUBANK, R. L., HUANG, C., MUÑOZ MALDONADO, Y., WANG, N., WANG, S. and BUCHANAN, R. J. (2004).   | 18 |
| 19  |      | Smoothing spline estimation in varying coefficient models. Journal of the Royal Statistical Society,  | 19 |
| 20  | [12] | EUBANK, R. L. (1988). Spline Smoothing and Nonparametric Regression. 1st ed. Marcel Dekker  | 20 |
| 21  | []   | Inc., New York.   | 21 |
| 22  | [13] | EUBANK, R. L., HUANG, C. and WANG, S. (2003). Adaptive order selection for spline smoothing.  | 22 |
| 23  | [14] | Journal of Computational and Graphical Statistics, <b>12</b> , 2546–2559.   | 23 |
| 24  | [14] | DE FINETTI, B. (1964). Foresight: its logical laws, its subjective sources in <i>Studies in Subjective</i>  | 24 |
| 25  |      | Probability. Wiley, New York.   | 25 |
| 26  | [16] | GANTMAKHER, F. R. (1959). The Theory of Matrices. Chelsea Pub. Co., New York.   | 26 |
| 27  | [17] | GOLUB, G. H., HEATH, M. and WAHBA, G. (1979). Generalized-cross validation as a method for choosing a good ridge parameter. <i>Technometrics</i> <b>58</b> , 215–223                                  | 27 |
| 28  | [18] | GRAYBILL, F. A. (1976). The theory and application of the linear model. Duxbury, North Scituate,  | 28 |
| 20  |      | Massachusetts.  | 20 |
| 20  | [19] | Guo, W. (2002). Functional mixed effects models. <i>Biometrics</i> , <b>58</b> , 121–128.   | 23 |
| 30  | [20] | tical Methods in Medical Research 13 1–24   | 30 |
| 31  | [21] | HASTIE, T. and TIBSHIRANI, R. (1993). Varying-coefficient models. Journal of the Royal Statistical  | 31 |
| 32  |      | Society, Series B, 4, 757–796.  | 32 |
| 33  | [22] | HECKMAN, N. (1997). The theory and application of penalized least squares methods or reproducing  | 33 |
| 34  | [23] | HENDERSON, C. R., KEMPTHORNE, O., SEARLE, S. R. and KROSIGK, C. M. (1959). The estimation of  | 34 |
| 35  | [-]  | environmental and genetic trends from records subject to culling. <i>Biometrics</i> , <b>15</b> , 192–218.  | 35 |
| 36  | [24] | HOUSEHOLDER, A. (1964). The Theory aof Matrices in Numerical Analysis. Dover, New York.   | 36 |
| 37  | [25] | KOHN, R. and ANSLEY, C. F. (1987). A new algorithm for spline smoothing based on smoothing a stochastic process. SIAM Journal of Scientific and Statistical Computing 8, 33–48                        | 37 |
| 38  | [26] | KOHN, R. and ANSLEY, C. F. (1989). A fast algorithm for signal extraction, influence and cross-   | 38 |
| 39  |      | validation in state-space models. Biometrika, <b>76</b> , 65–79.  | 39 |
| 40  | [27] | KOHN, R., ANSLEY, C. F. and THARM, D. (1993). Performance of cross-validation and maximum like-   | 40 |
| 41  |      | lihood estimators of spline smoothing parameters. <i>Journal of the American Statistical Association</i> ,<br><b>86</b> , 1042–1050   | 41 |
| 42  | [28] | KOOPMAN, S. J. and DURBIN, J. (1998). Fast filtering and smoothing for multivariate state space   | 42 |
| 43  | . ,  | models. Journal of Times Series Analysis, 21, 281–296.  | 43 |
| 44  | [29] | LINDLEY, D. V. and SMITH, A. F. M. (1972). Bayes estimates for the linear model. <i>Journal of the</i>  | 44 |
| 45  | [30] | PIEHPO, H. P. and OGUTU, J. O. (2007). Simple state-space models in a mixed-model framework   | 45 |
| 46  | [00] | The American Statistician, <b>61</b> , 224–232.   | 46 |
| 47  | [31] | NYCHKA, D. (1988). Bayesian confidence intervals for smoothing splines. Journal of the American   | 47 |
| 48  | [20] | Statistical Association, 83, 1134-1143.   | 49 |
| 10  | [34] | proaches, Computation and Application. Wiley, New York.   | 40 |
| -19 | [33] | PATTERSON, H. D. and THOMPSON, R. (1971). Recovery of interblock information when cell sizes are  | 49 |
| 50  | Ic t | unequal. Biometrika, <b>58</b> , 545–554.   | 50 |
| 51  | [34] | RIESBY, N., GRAM, L.F., BECH, P., NAGY, A., PETERSEN, G. O., ORTMANN, J., IBSEN, I., DENCKER, S.  | 51 |

| 1        |            | J., JACOBSEN, O., KRAUTWALD, O., SØNDERGAARD, I. and CHIRSTIANSEN, J. (1977). Imipramine: clinical effects and pharmacokinetic variability. <i>Psychopharmacology</i> , <b>54</b> , 263–272. | 1      |
|----------|------------|--|--------|
| 3        | [35]       | ROBINSON, G. K. (1991). That BLUP is a good thing: the estimation of random effects. <i>Statistical Science</i> <b>6</b> 15–32   | 3      |
| 4        | [36]       | RUPPERT, D., WAND, M. P. and CARROLL, R. J. (2003). Semiparametric Regression. Cambridge   | 4      |
| 5        | [37]       | University Press, New York.<br>RUPPERT, D. (2002). Selecting the number of knots for penalized splines. <i>Journal of Computational</i>  | 5      |
| 6        |            | and Graphical Statistics, 11, 735–757.   | 6      |
| 7        | [38]       | SALLAS, W. M. and HARVILLE, D. A. (1981). Best linear recursive estimation for mixed linear models.<br>Journal of the American Statistical Association, <b>76</b> , 860–869.                 | 7      |
| 0<br>0   | [39]       | SCHAFER, J. L. and RECAI, M. Y. (2002). Computational strategies for multivariate linear mixed-  | ہ<br>م |
| 10       |            | effects models with missing values. Journal of Computational and Graphical Statistics, 11, 437–  | 10     |
| 11       | [40]       | SPEED, T. (1991). That BLUP is a good Thing: the estimation of random effects: comment. <i>Statis</i> -  | 11     |
| 12       | [41]       | tical Science, $6$ , 42–44.  | 12     |
| 13       | [41]       | for estimating the parameters of a stochastic process. Annals of Statistics, <b>18</b> , 1139–1157.  | 13     |
| 14       | [42]       | STRANDÉN, I. and LIDAUER, M. (1999). Solving large mixed linear models using preconditioned con-   | 14     |
| 15       | [43]       | WAHBA, G. (1978). Improper priors, spline smoothing, and the problem of guarding against model   | 15     |
| 16       | <br>[ 4 4] | errors in regression. Journal of the Royal Statistical Society, Series B, 40, 364–372.   | 16     |
| 17       | [44]       | of the Royal Statistical Society, Series B, 40, 364–372.   | 17     |
| 18       | [45]       | WAHBA, G. (1990). Spline models for observational data. <i>SIAM</i> , <b>59</b> . Philadelphia, Pennsylvania.  | 18     |
| 19       | [46]       | WAHBA, G. (1985). A comparison of GCV and GML for choosing the smoothing parameter in the generalized spline smoothing problem. <i>Annals of Statistics</i> , <b>13</b> , 1378–1402.         | 19     |
| 20       | [47]       | WAHBA, G. (1990). Comment on Cressie: letters to the editor. The American Statistician, 44,  | 20     |
| 22       | [48]       | 255–256.<br>WANG Y (1998) Smoothing spline models with correlated random errors. <i>Journal of the American</i>  | 22     |
| 23       | [10]       | Statistical Association, 93, 341–348.  | 23     |
| 24       | [49]       | WANG, Y. (1998). Mixed effects smoothing spline analysis of variance. Journal of the Royal Sta-<br>tistical Society. Series B 60, 159–174  | 24     |
| 25       |            | <i>usucai becciy, beries D</i> , <b>00</b> , 100 114.  | 25     |
| 26       |            |  | 26     |
| 27       |            |  | 27     |
| 28       |            |  | 28     |
| 29       |            |  | 29     |
| 30       |            |  | 30     |
| 30       |            |  | 32     |
| 33       |            |  | 33     |
| 34       |            |  | 34     |
| 35       |            |  | 35     |
| 36       |            |  | 36     |
| 37       |            |  | 37     |
| 38       |            |  | 38     |
| 39       |            |  | 39     |
| 40       |            |  | 40     |
| 41       |            |  | 41     |
| 42       |            |  | 42     |
| 43<br>44 |            |  | 43     |
| 45       |            |  | 45     |
| 46       |            |  | 46     |
| 47       |            |  | 47     |
| 48       |            |  | 48     |
| 49       |            |  | 49     |
| 50       |            |  | 50     |
| 51       |            |  | 51     |

# Yolanda Muñoz Maldonado

imsart-coll ver. 2008/08/29 file: Munoz.tex date: March 25, 2009

| Xia Chen <sup>1,*</sup> and Davar Khoshnevisan <sup>2,†</sup><br>University of Tennessee and University of Utah Abstract: We prove that two seemingly-different models of random walk in<br>random environment are generically quite close to one another. One model<br>comes from statistical physics, and describes the behavior of a randomly-<br>charged random polymer. The other model comes from probability theory, and<br>was originally designed to describe a large family of asymptotically self-similar<br>processes that have stationary increments. Contents 1 Introduction and the Main Results   |
|---|
| University of Tennessee and University of Utah         Abstract: We prove that two seemingly-different models of random walk in random environment are generically quite close to one another. One model comes from statistical physics, and describes the behavior of a randomly-charged random polymer. The other model comes from probability theory, and was originally designed to describe a large family of asymptotically self-similar processes that have stationary increments.         Contents         1       Introduction and the Main Results         2       Preliminary Estimates         3       Proofs of the Main Results         4       Acknowledgement.         5       Acknowledgement.         6       References         7       Introduction and the Main Results         8       Introduction and the Main Results         9       Preliminary Estimates         1       Introduction and the Main Results         2       Preliminary Estimates         3       Proofs of the Main Results         4       Acknowledgement.         6       Acknowledgement.         7       References         1       Introduction and the Main Results         1       Introduction and the Main Results         1       Introduction and the Main Results         1       Introduction and the Main Results      < |
| <ul> <li>Abstract: We prove that two seemingly-different models of random walk in random environment are generically quite close to one another. One model comes from statistical physics, and describes the behavior of a randomly-charged random polymer. The other model comes from probability theory, and was originally designed to describe a large family of asymptotically self-similar processes that have stationary increments.</li> <li>Contents</li> <li>1 Introduction and the Main Results</li></ul>  |
| Contents          1       Introduction and the Main Results       2         2       Preliminary Estimates       2         3       Proofs of the Main Results       2         3       Proofs of the Main Results       2         Acknowledgement.       2         References       2         I. Introduction and the Main Results       2         The principal goal of this article is to show that two apparently-disparate models one from statistical physics of disorder media (Kantor and Kardar (1991), Derr et al (1992), Derrida and Higgs (1994)) and one from probability theory (Kest and Spitzer (1979), Bolthausen (1989))—are very close to one another.         In order to describe the model from statistical physics, let us suppose that $q$ $\{a_k\}_{\infty}^{\infty}$ , is a collection of i.i.d. mean-zero random variables with finite variants   |
| 1       Introduction and the Main Results       2         2       Preliminary Estimates       2         3       Proofs of the Main Results       2         3       Proofs of the Main Results       2         4       Acknowledgement.       2         5       References       2         7       Introduction and the Main Results       2         8       Introduction and the Main Results       2         9       Introduction and the Main Results       2         10       Introduction and the Main Results       2         11       Introduction and the Main Results       2         12       Introduction and the Main Results       2         13       Introduction and the Main Results       2         14       Introduction and Higgs (1994)       2         15       Introduction and Higgs (1989)       2         16       2       2  |
| 1. Introduction and the Main Results<br>The principal goal of this article is to show that two apparently-disparate models<br>one from statistical physics of disorder media (Kantor and Kardar (1991), Derr<br>et al (1992), Derrida and Higgs (1994)) and one from probability theory (Kes<br>and Spitzer (1979), Bolthausen (1989))—are very close to one another.<br>In order to describe the model from statistical physics, let us suppose that $q$<br>$\{a_i\}_{i=1}^{\infty}$ is a collection of i.i.d. mean-zero random variables with finite varian   |
| The principal goal of this article is to show that two apparently-disparate models<br>one from statistical physics of disorder media (Kantor and Kardar (1991), Derr<br>et al (1992), Derrida and Higgs (1994)) and one from probability theory (Kes<br>and Spitzer (1979), Bolthausen (1989))—are very close to one another.<br>In order to describe the model from statistical physics, let us suppose that $q$<br>$\{a_i\}_{i=1}^{\infty}$ is a collection of i.i.d. mean-zero random variables with finite varian   |
| $\sigma^2 > 0$ . For technical reasons, we assume here and throughout that  |
| (1.1) $\mu_6 := \mathrm{E}(q_1^6) < \infty.$  |
| In addition, we let $S := \{S_i\}_{i=0}^{\infty}$ denote a random walk on $\mathbb{Z}^d$ with $S_0 = 0$ the is independent from the collection $q$ . We also rule out the trivial case that $S_1$ is only one possible value.<br>The object of interest to us is the random quantity  |
| (1.2) $H_n := \sum_{1 \le i < j \le n} q_i q_j 1_{\{S_i = S_j\}}.$  |

davar@math.utah.edu  $^{\dagger}\textsc{Research}$  supported in part by NSF grant DMS-0706728. AMS 2000 subject classifications: primary 60K35; secondary 60K37

Keywords and phrases: Polymer measures, random walk in random scenery  In statistical physics,  $H_n$  denotes a random Hamiltonian of spin-glass type that is used to build Gibbsian polymer measures. The  $q_i$ 's are random charges, and each realization of S corresponds to a possible polymer path; see the paper by Kantor and Kardar (1991), its subsequent variations by Derrida et al (1992, 1994) and Wittmer et al (1993), and its predecessos by Garel and Orland (1988) and Obukhov (1986). The resulting Gibbs measure then corresponds to a model for "random walk in random environment." Although we do not consider continuous processes here, the continuum-limit analogue of  $H_n$  has also been studied in the literature (Buffet and Pulé (1997), Martinez and Petritis (1996)).

Kesten and Spitzer (1979) introduced a different model for "random walk in random environment," which they call random walk in random scenery.<sup>1</sup> We can describe that model as follows: Let  $Z := \{Z(x)\}_{x \in \mathbb{Z}^d}$  denote a collection of i.i.d. random variables, with the same common distribution as  $q_1$ , and independent of S. Define

(1.3) 
$$W_n := \sum_{i=1}^n Z(S_i).$$

The process  $W := \{W_n\}_{n=0}^{\infty}$  is called random walk in random scenery, and can be thought of as follows: We fix a realization of the *d*-dimensional random field Z—the "scenery"—and then run an independent walk S on  $\mathbb{Z}^d$ . At time *j*, the walk is at  $S_j$ ; we sample the scenery at that point. This yields  $Z(S_j)$ , which is then used as the increment of the process W at time *j*.

Our goal is to make precise the assertion that if n is large, then

(1.4) 
$$H_n \approx \gamma^{1/2} \cdot W_n$$
 in distribution,

where

(1.5) 
$$\gamma := \begin{cases} 1 & \text{if } S \text{ is recurrent,} \\ \sum_{k=1}^{\infty} P\{S_k = 0\} & \text{if } S \text{ is transient.} \end{cases}$$

Our derivation is based on a classification of recurrence vs. transience for random walks that appears to be new. This classification [Theorem 2.4] might be of independent interest.

We can better understand (1.4) by considering separately the cases that S is transient versus recurrent. The former case is simpler to describe, and appears next.

**Theorem 1.1.** If S is transient, then

(1.6) 
$$\frac{W_n}{n^{1/2}} \xrightarrow{\mathcal{D}} N(0, \sigma^2) \quad and \quad \frac{H_n}{n^{1/2}} \xrightarrow{\mathcal{D}} N(0, \gamma \sigma^2).$$

Kesten and Spitzer (1979) proved the assertion about  $W_n$  under more restrictive conditions on S. Similarly, Chen (2008) proved the statement about  $H_n$  under more hypotheses.

Before we can describe the remaining [and more interesting] recurrent case, we define

(17) 
$$a_{n} := \left(n \sum_{k=0}^{n} P\{S_{k} = 0\}\right)^{1/2}$$

$$\begin{array}{c} _{49} \\ _{50} \end{array} (1.7) \qquad \qquad a_n := \left( n \sum_{k=0}^{\infty} P\{S_k = 0\} \right) \quad . \end{array}$$

<sup>&</sup>lt;sup>1</sup>Kesten and Spitzer ascribe the terminology to Paul Shields.

It is well known (Polya (1921), Chung and Fuchs (1951)) that S is recurrent if and only if  $a_n/n^{1/2} \to \infty$  as  $n \to \infty$ . 

**Theorem 1.2.** If S is recurrent, then for all bounded continuous functions f:  $\mathbf{R}^d \to \mathbf{R}$ .

(1.8) 
$$\mathbf{E}\left[f\left(\frac{W_n}{a_n}\right)\right] = \mathbf{E}\left[f\left(\frac{H_n}{a_n}\right)\right] + o(1),$$

where o(1) converges to zero as  $n \to \infty$ . Moreover, both  $\{W_n/a_n\}_{n\geq 1}$  and  $\{H_n/a_n\}_{n>1}$  are tight.

We demonstrate Theorems 1.1 and 1.2 by using a variant of the replacement method of Liapounov (1900) [pp. 362–364]; this method was rediscovered later by Lindeberg (1922), who used it to prove his famous central limit theorem for triangular arrays of random variables.

It can be proved that when S is in the domain of attraction of a stable law,  $W_n/a_n$  converges in distribution to an explicit law (Kesten and Spitzer (1979), Bolthausen (1989)). Consequently,  $H_n/a_n$  converges in distribution to the same law in that case. This fact was proved earlier by Chen (2008) under further [mild] conditions on S and  $q_1$ .

We conclude the introduction by describing the growth of  $a_n$  under natural conditions on S.

**Remark 1.3.** Suppose S is strongly aperiodic, mean zero, and finite second moments, with a nonsingular covariance matrix. Then, S is transient iff  $d \ge 3$ , and by the local central limit theorem, as  $n \to \infty$ ,

(1.9) 
$$\sum_{k=1}^{n} P\{S_k = 0\} \sim \text{const} \times \begin{cases} n^{1/2} & \text{if } d = 1, \\ \log n & \text{if } d = 2. \end{cases}$$

See, for example (Spitzer (1976) [**P9** on p. 75]). Consequently,

(1.10) 
$$a_n \sim \text{const} \times \begin{cases} n^{3/4} & \text{if } d = 1, \\ (n \log n)^{1/2} & \text{if } d = 2. \end{cases}$$

This agrees with the normalization of Kesten and Spitzer (1979) when d = 1, and Bolthausen (1989) when d = 2.

#### 2. Preliminary Estimates

Consider the local times of S defined by

(2.1) 
$$L_n^x := \sum_{i=1}^n \mathbf{1}_{\{S_i = x\}}.$$

A little thought shows that the random walk in random scenery can be represented compactly as

50 (2.2) 
$$W_n = \sum_{x \in \mathbb{Z}^d} Z(x) L_n^x.$$
 50 51

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

З

There is also a nice way to write the random Hamiltonian  $H_n$  in local-time terms. Consider the "level sets,"

(2.3) 
$$\mathcal{L}_{n}^{x} := \{i \in \{1, \dots, n\} : S_{i} = x\}.$$

It is manifest that if  $j \in \{2, ..., n\}$ , then  $L_j^x > L_{j-1}^x$  if and only if  $j \in \mathcal{L}_n^x$ . Thus, we can write

$$H_n = \frac{1}{2} \left( \sum_{x \in \mathbf{Z}^d} \left| \sum_{i=1}^n q_i \mathbf{1}_{\{S_i = x\}} \right|^2 - \sum_{i=1}^n q_i^2 \right)$$

З

$$=\sum_{x\in\mathbf{Z}^d}h_n^x,$$

where

(2.4)

(2.5) 
$$h_n^x := \frac{1}{2} \left( \left| \sum_{i \in \mathcal{L}_n^x} q_i \right|^2 - \sum_{i \in \mathcal{L}_n^x} q_i^2 \right).$$

We denote by  $\widehat{P}$  the conditional measure, given the entire process S;  $\widehat{E}$  denotes the corresponding expectation operator. The following is borrowed from Chen (2008) [Lemma 2.1].

**Lemma 2.1.** Choose and fix some integer  $n \ge 1$ . Then,  $\{h_n^x\}_{x \in \mathbb{Z}^d}$  is a collection of independent random variables under  $\widehat{P}$ , and

(2.6) 
$$\widehat{\mathrm{E}}h_n^x = 0 \quad and \quad \widehat{\mathrm{E}}\left(\left|h_n^x\right|^2\right) = \frac{\sigma^2}{2}L_n^x\left(L_n^x - 1\right) \qquad \mathrm{P}\text{-}a.s.$$

Moreover, there exists a nonrandom positive and finite constant  $C = C(\sigma)$  such that for all  $n \ge 1$  and  $x \in \mathbb{Z}^d$ ,

(2.7) 
$$\widehat{\mathrm{E}}\left(\left|h_{n}^{x}\right|^{3}\right) \leq C\mu_{6}\left|L_{n}^{x}\left(L_{n}^{x}-1\right)\right|^{3/2} \qquad \mathrm{P}\text{-}a.s.$$

Next we develop some local-time computations.

Lemma 2.2. For all  $n \geq 1$ ,

(2.8) 
$$\sum_{x \in \mathbf{Z}^d} EL_n^x = n \quad and \quad \sum_{x \in \mathbf{Z}^d} E\left(|L_n^x|^2\right) = n + 2\sum_{k=1}^{n-1} (n-k) P\{S_k = 0\}.$$

Moreover, for all integers  $k \geq 1$ ,

$$\sum_{k=1}^{n} \left( |I^{x}|^{k} \right) \leq k \ln \left| \sum_{k=1}^{n} \mathbb{P} \left[ S_{k} = 0 \right] \right|^{k-1}$$

(2.9) 
$$\sum_{x \in \mathbf{Z}^d} \mathbb{E}\left(\left|L_n^x\right|^k\right) \le k! \, n \left|\sum_{j=0} \mathbb{P}\{S_j = 0\}\right| \quad .$$

*Proof.* Since  $EL_n^x = \sum_{j=1}^n P\{S_j = x\}$  and  $\sum_{x \in \mathbb{Z}^d} P\{S_j = x\} = 1$ , we have  $\sum_x EL_n^x = n$ . For the second-moment formula we write

$$E\left(|L_{n}^{x}|^{2}\right) = \sum_{1 \le i \le n} P\{S_{i} = x\} + 2\sum_{1 \le i \le j \le n} P\{S_{i} = S_{j} = x\}$$

(2.10) 
$$1 \le i \le n \qquad 1 \le i < j \le n$$
$$\sum_{i=1}^{n} \mathbb{P}[G_{i} = i] + 2 \sum_{i=1}^{n} \mathbb{P}[G_{i} = i] \mathbb{P}[G_{i} = i]$$

50  
51  
51  

$$= \sum_{1 \le i \le n} P\{S_i = x\} + 2 \sum_{1 \le i < j \le n} P\{S_i = x\} P\{S_{j-i} = 0\}.$$
50  
51  
51  
51

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

We can sum this expression over all  $x \in \mathbf{Z}^d$  to find that  $\sum_{x \in \mathbf{Z}^d} \mathbb{E}\left(|L_n^x|^2\right) = n + 2 \sum_{1 \le i \le j \le n} \mathbb{P}\{S_{j-i} = 0\}.$ (2.11)This readily implies the second-moment formula. Similarly, we write  $\mathrm{E}\left(\left|L_{n}^{x}\right|^{k}\right)$  $\leq k! \sum_{1 \leq i_1 \leq \dots \leq i_k \leq n} \mathbf{P}\{S_{i_1} = \dots = S_{i_k} = x\}$  $=k! \sum_{1 \le i_1 \le \dots \le i_k \le n} \mathbf{P}\{S_{i_1} = x\} \mathbf{P}\{S_{i_2-i_1} = 0\} \cdots \mathbf{P}\{S_{i_k-i_{k-1}} = 0\}$ (2.12) $\leq k! \sum_{i=1}^{n} P\{S_i = x\} \cdot \left| \sum_{i=1}^{n} P\{S_j = 0\} \right|^{k-1}.$ Add over all  $x \in \mathbf{Z}^d$  to finish. Our next lemma provides the first step in a classification of recurrence versus transience for random walks. Lemma 2.3. It is always the case that  $\lim_{n \to \infty} \frac{1}{n} \sum_{x \in \mathbb{Z}^d} \mathbb{E}\left( |L_n^x|^2 \right) = 1 + 2 \sum_{k=1}^{\infty} \mathbb{P}\{S_k = 0\}.$ (2.13)*Proof.* Thanks to Lemma 2.2, for all  $n \ge 1$ ,  $\frac{1}{n} \sum_{x \in I} \mathbb{E}\left(|L_n^x|^2\right) = 1 + 2\sum_{k=1}^{n-1} \left(1 - \frac{k}{n}\right) \mathbb{P}\{S_k = 0\}.$ (2.14)If S is transient, then the monotone convergence theorem ensures that  $\lim_{n \to \infty} \frac{1}{n} \sum_{n \to \infty} E\left( |L_n^x|^2 \right) = 1 + 2 \sum_{k=1}^{\infty} P\{S_k = 0\}.$ (2.15)This proves the lemma in the transient case. When S is recurrent, we note that (2.14) readily implies that for all integers  $m \geq 2$ ,  $\liminf_{n \to \infty} \frac{1}{n} \sum_{x \in \mathbb{Z}^d} \mathbb{E}\left( |L_n^x|^2 \right) \ge 1 + 2 \sum_{k=1}^{m-1} \left( 1 - \frac{k}{m} \right) \mathbb{P}\{S_k = 0\}$ (2.16) $\geq 1 + \sum_{1 \le k \le m/2} \mathbb{P}\{S_k = 0\}.$ Let  $m \uparrow \infty$  to deduce the lemma.

Next we "remove the expectation" from the statement of Lemma 2.3.

З

## 

**Theorem 2.4.** As  $n \to \infty$ ,

(2.17) 
$$\frac{1}{n} \sum_{x \in \mathbf{Z}^d} \left( L_n^x \right)^2 \to 1 + 2 \sum_{k=1}^\infty \mathrm{P}\{S_k = 0\} \quad in \text{ probability.}$$

**Remark 2.5.** The quantity  $I_n := \sum_{x \in \mathbb{Z}^d} (L_n^x)^2$  is the so-called *self-intersection local time* of the walk S. This terminology stems from the following elementary calculation: For all integers n > 1,

(2.18) 
$$I_n = \sum_{1 \le i, j \le n} \mathbf{1}_{\{S_j = S_i\}}.$$

Consequently, Theorem 2.4 implies that a random walk S on  $\mathbf{Z}^d$  is recurrent if and only if its self-intersection local time satisfies  $I_n/n \to \infty$  in probability.

Remark 2.6. Nadine Guillotin–Plantard has kindly pointed out to us that the mode of convergence in Theorem 2.4 can be strengthened to almost-sure convergence. This requires a direct subadditivity argument (Guillotin–Plantard (2004)). It follows also from the estimates that follow, together with a classical blocking argument, which we skip.

*Proof.* First we study the case that  $\{S_i\}_{i=0}^{\infty}$  is transient. Define

(2.19) 
$$Q_n := \sum_{1 \le i < j \le n} \mathbf{1}_{\{S_i = S_j\}}.$$

Then it is not too difficult to see that

(2.20) 
$$\sum_{x \in \mathbf{Z}^d} (L_n^x)^2 = 2Q_n + n \quad \text{for all } n \ge 1.$$

This follows immediately from (2.18), for example. Therefore, it suffices to prove that, under the assumption of transience,

(2.21) 
$$\frac{Q_k}{k} \to \sum_{j=1}^{\infty} \mathbb{P}\{S_j = 0\} \text{ in probability as } k \to \infty.$$

Lemma 2.3 and (2.20) together imply that

(2.22) 
$$\lim_{k \to \infty} \frac{\mathbf{E}Q_k}{k} = \sum_{j=1}^{\infty} \mathbf{P}\{S_j = 0\}.$$

Hence, it suffices to prove that  $\operatorname{Var} Q_n = o(n^2)$  as  $n \to \infty$ . In some cases, this can be done by making an explicit [though hard] estimate for  $\operatorname{Var} Q_n$ ; see, for instance, (Chen (2008) [Lemma 5.1]), and also the technique employed in the proof of Lemma 2.4 of Bolthausen (1989). Here, we opt for a more general approach that is simpler, though it is a little more circuitous. Namely, in rough terms, we write  $Q_n$  as  $Q_n^{(1)} + Q_n^{(2)}$ , where  $EQ_n^{(1)} = o(n)$ , and  $\operatorname{Var} Q_n^{(2)} = o(n^2)$ . Moreover, we will soon see that  $Q_n^{(1)}, Q_n^{(2)} \ge 0$ , and this suffices to complete the proof. 

For all  $m := m_n \in \{1, \ldots, n-1\}$  we write

51 (2.23) 
$$Q_n = Q_n^{1,m} + Q_n^{2,m},$$
 51

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

З



(2.26) 
$$\Upsilon(\Gamma) := \sum_{(i,j)\in\Gamma} \mathbf{1}_{\{S_i=S_j\}}.$$

$$33$$
34
35

Suppose  $\Gamma_1, \Gamma_2, \ldots, \Gamma_{\nu}$  are finite disjoint sets in  $\mathbf{N}^2$ , with common cardinality, and the added property that whenever  $1 \leq a < b \leq \nu$ , we have  $\Gamma_a < \Gamma_b$  in the sense that i < k and j < l for all  $(i, j) \in \Gamma_a$  and  $(k, l) \in \Gamma_b$ . Then, it follows that

(2.27) 
$$\{\Upsilon(\Gamma_{\nu})\}_{\mu=1}^{\nu}$$
 is an i.i.d. sequence.

For all integers  $p \ge 0$  define

(2.28) 
$$B_p^m := \left\{ (i, j) \in \mathbf{N}^2 : (p-1)m < i < j \le pm \right\},$$

n-1

(2.20) 
$$W_p^m := \{(i,j) \in \mathbf{N}^2 : (p-1)m < i \le pm < j \le (p+1)m\}.$$

In Figure 1,  $\{B_p^m\}_{p=1}^{\infty}$  denotes the collection black and  $\{W_p^m\}_{p=1}^{\infty}$  the white triangles that are inside the slanted strip.

We may write

<sup>50</sup> (2.29) 
$$Q_{(n-1)m}^{2,m} = \sum_{p=1} \Upsilon(B_p^m) + \sum_{p=1} \Upsilon(W_p^m).$$
<sup>50</sup> <sup>51</sup>

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

n-1

Consequently,

$$\operatorname{Var} Q_{(n-1)m}^{2,m} \leq 2\operatorname{Var} \sum_{p=1}^{n-1} \Upsilon(B_p^m) + 2\operatorname{Var} \sum_{p=1}^{n-1} \Upsilon(W_p^m).$$

If  $1 \le a < b \le m-1$ , then  $B_a^m < B_b^m$  and  $W_a^m < W_b^m$ . Consequently, (2.27) implies that

(2.31) 
$$\operatorname{Var} Q^{2,m}_{(n-1)m} \le 2(n-1) \left[ \operatorname{Var} \Upsilon(B^m_1) + \operatorname{Var} \Upsilon(W^m_1) \right].$$

Because  $\Upsilon(B_1^m)$  and  $\Upsilon(W_1^m)$  are individually sums of not more than  $\binom{m}{2}$ -many ones.

(2.32) 
$$\operatorname{Var} Q_{(n-1)m}^{2,m} \le 2(n-1)m^2.$$

Let  $Q_n^{(1)} := Q_n^{1,m}$  and  $Q_n^{(2)} := Q_n^{2,m}$ , where  $m = m_n := n^{1/4}$  [say]. Then,  $Q_n = Q_n^{(1)} + Q_n^{(2)}$ , and (2.25) and (2.32) together imply that  $EQ_{(n-1)m}^{(1)} = o((n-1)m)$ . Moreover,  $\operatorname{Var} Q_{(n-1)m}^{(2)} = o((nm)^2)$ . This gives us the desired decomposition of  $Q_{(n-1)m}$ . Now we complete the proof: Thanks to (2.22),

(2.33) 
$$\operatorname{E}Q_{(n-1)m}^{(2)} \sim nm \cdot \sum_{j=1}^{\infty} \operatorname{P}\{S_j = 0\} \text{ as } n \to \infty.$$

Therefore, the variance of  $Q_{(n-1)m}^{\left(2\right)}$  is little-o of the square of its mean. This and the Chebyshev inequality together imply that  $Q_{(n-1)m}^{(2)}/(nm)$  converges in probability to  $\sum_{j=1}^{\infty} P\{S_j = 0\}$ . On the other hand, we know also that  $Q_{(n-1)m}^{(1)}/(nm)$  converges to zero in  $L^1(\mathbf{P})$  and hence in probability. Consequently, we can change variables and note that as  $n \to \infty$ ,

(2.34) 
$$\frac{Q_{nm}}{nm} \to \sum_{j=1}^{\infty} \mathbb{P}\{S_j = 0\} \text{ in probability.}$$

If k is between (n-1)m and nm, then

(2.35) 
$$\frac{Q_{(n-1)m}}{nm} \le \frac{Q_k}{k} \le \frac{Q_{nm}}{(n-1)m}.$$
 38  
39

This proves (2.21), and hence the theorem, in the transient case.

In order to derive the recurrent case, it suffices to prove that  $Q_n/n \to \infty$  in probability as  $n \to \infty$ .

Let us choose and hold an integer  $m \ge 1$ —so that it does not grow with n—and observe that  $Q_n \ge Q_n^{2,m}$  as long as n is sufficiently large. Evidently,

$$EQ_n^{2,m} = \sum \sum P\{S_j = S_i\}$$
<sup>46</sup>
<sub>47</sub>

 (2.36)

$$= (n-1) \sum_{k=1}^{m-1} \mathbf{P}\{S_k = 0\}.$$
<sup>50</sup>
<sub>51</sub>

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

З

We may also observe that (2.32) continues to hold in the present recurrent setting.

Together with the Chebyshev inequality, these computations imply that as  $n \to \infty$ ,

(2.37) 
$$\frac{Q_{n(m-1)}^{2,m}}{n} \to \sum_{k=1}^{m-1} \mathbb{P}\{S_k = 0\} \text{ in probability.}$$

Because  $Q_{n(m-1)} \ge Q_{n(m-1)}^{2,m}$ , the preceding implies that

(2.38) 
$$\lim_{n \to \infty} \Pr\left\{\frac{Q_{n(m-1)}}{n} \ge \frac{1}{2} \sum_{k=1}^{m} \Pr\{S_k = 0\}\right\} = 1.$$

A monotonicity argument shows that  $Q_{n(m-1)}$  can be replaced by  $Q_n$  without altering the end-result; see (2.35). By recurrence, if  $\lambda > 0$  is any predescribed positive number, then we can choose [and fix] our integer m such that  $\sum_{k=1}^{m} \mathbb{P}\{S_k = 0\} \ge 2\lambda$ . This proves that  $\lim_{n \to \infty} \mathbb{P}\{Q_n / n \ge \lambda\} = 1$  for all  $\lambda > 0$ , and hence follows the theorem in the recurrent case.

### 3. Proofs of the Main Results

Now we introduce a sequence  $\{\xi_x\}_{x \in \mathbf{Z}^d}$  of random variables, independent [under P] of  $\{q_i\}_{i=1}^{\infty}$  and the random walk  $\{S_i\}_{i=0}^{\infty}$ , such that

(3.1) 
$$E\xi_0 = 0, \quad E(\xi_0^2) = \sigma^2, \text{ and } \widehat{\mu}_3 := E(|\xi_0|^3) < \infty.$$

Define

(3.2) 
$$\widehat{h}_n^x := \left| \frac{L_n^x \left( L_n^x - 1 \right)}{2} \right|^{1/2} \xi_x \quad \text{for all } n \ge 1 \text{ and } x \in \mathbf{Z}^d.$$

Evidently,  $\{\hat{h}_n^x\}_{x \in \mathbb{Z}^d}$  is a sequence of [conditionally] independent random variables, under  $\widehat{P}$ , and has the same [conditional] mean and variance as  $\{h_n^x\}_{x\in\mathbb{Z}^d}$ .

**Lemma 3.1.** There exists a positive and finite constant  $C_* = C_*(\sigma)$  such that if  $f: \mathbf{R}^d \to \mathbf{R}$  is three time continuously differentiable, then for all  $n \geq 1$ ,

(3.3) 
$$\left| \operatorname{E} f\left(\sum_{x \in \mathbf{Z}^d} \widehat{h}_n^x\right) - \operatorname{E} f(H_n) \right| \le C_* M_f(\widehat{\mu}_3 + \mu_6) n \left| \sum_{j=0}^n \operatorname{P} \{S_j = 0\} \right|^2,$$

with  $M_f := \sup_{x \in \mathbf{R}^d} |f'''(x)|$ .

 $f(H_n)$ 

*Proof.* Temporarily choose and fix some  $y \in \mathbf{Z}^d$ , and notice that

$$(3.4) = f\left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) + f'\left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) h_n^y + \frac{1}{2}f''\left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) |h_n^y|^2 + R_n,$$

$$+R_n,$$

where  $|R_n| \leq \frac{1}{6} ||f'''||_{\infty} |h_n^y|^3$ . It follows from this and Lemma 2.1 that

$$\widehat{\mathrm{E}}f(H_n)$$

$$\begin{array}{ll} {}^{49} & (3.5) \\ {}^{50} & \\ {}^{51} & \\ \end{array} & = \widehat{\mathrm{E}}f\left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) + \frac{\sigma^2}{2} L_n^y \left(L_n^y - 1\right) \widehat{\mathrm{E}}f'' \left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) + R_n^{(1)}, \\ {}^{50} & {}^{50} \\ {}^{51} & {}^{51} & {}^{51} \end{array} \right)$$

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

where

(3.6)

$$\left| R_n^{(1)} \right| \le \frac{CM_f \mu_6}{12} \left| L_n^x \left( L_n^x - 1 \right) \right|^{3/2} \qquad P-$$
  
$$\le \frac{CM_f \mu_6}{12} \left| L_n^y \right|^3.$$

$$\sum_{n=1}^{\infty} (L_n^x - 1)|^{3/2}$$
 P-a.s.

We proceed as in (3.4) and write

$$f\left(\widehat{h}_{n}^{y}+\sum_{x\in\mathbf{Z}^{d}\setminus\{y\}}h_{n}^{x}
ight)$$

$$(3.7) = f\left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) + f'\left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) \widehat{h}_n^y + \frac{1}{2}f''\left(\sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right) \left|\widehat{h}_n^y\right|^2 + \widehat{R}_n,$$

where  $|\hat{R}_n| \leq \frac{1}{6}M_f|\hat{h}_n^y|^3 \leq \frac{1}{12\sqrt{2}}M_f|L_n^y|^3|\xi_y|^3$ . It follows from this and Lemma 2.1 that

(3.8) 
$$\widehat{\mathrm{E}}f\left(\widehat{h}_{n}^{y} + \sum_{x \in \mathbf{Z}^{d} \setminus \{y\}} h_{n}^{x}\right)$$

$$=\widehat{\mathrm{E}}f\left(\sum_{x\in\mathbf{Z}^d\setminus\{y\}}h_n^x\right)+\frac{\sigma^2}{2}L_n^y\left(L_n^y-1\right)\widehat{\mathrm{E}}f''\left(\sum_{x\in\mathbf{Z}^d\setminus\{y\}}h_n^x\right)+R_n^{(2)},$$

where  $|R_n^{(2)}| \leq \frac{1}{12\sqrt{2}} \hat{\mu}_3 M_f |L_n^y|^3$ . Define  $C_* := (C+1)/2$  to deduce from the preceding and (3.5) that P-a.s.,

(3.9) 
$$\left|\widehat{\mathrm{E}}f\left(\widehat{h}_{n}^{y}+\sum_{x\in\mathbf{Z}^{d}\setminus\{y\}}h_{n}^{x}\right)-\widehat{\mathrm{E}}f\left(\sum_{x\in\mathbf{Z}^{d}}h_{n}^{x}\right)\right|\leq\frac{A}{6}|L_{n}^{y}|^{3},$$

where  $A := C_* M_f(\hat{\mu}_3 + \mu_6)$ . The preceding computes the effect of replacing the contribution of  $h_n^x$  to  $H_n$  but the independent quantity  $\hat{h}_n^y$ , for each fixed y, and uses only the fact that the  $\hat{h}$ 's are a conditionally independent sequence with the same means and variances as their corresponding h's. Therefore, if we choose and fix another point  $y \in \mathbb{Z}^d \setminus \{y\}$ , then the very same constant A satisfies the following: Almost surely [P],

$$(3.10) \quad \left|\widehat{\mathrm{E}}f\left(\widehat{h}_n^z + \widehat{h}_n^y + \sum_{x \in \mathbf{Z}^d \setminus \{y,z\}} h_n^x\right) - \widehat{\mathrm{E}}f\left(\widehat{h}_n^y + \sum_{x \in \mathbf{Z}^d \setminus \{y\}} h_n^x\right)\right| \le \frac{A}{6} |L_n^z|^3.$$

And hence, the triangle inequality yields the following: P-a.s.,

$$\left| \widehat{\mathrm{E}}f\left(\widehat{h}_{n}^{z} + \widehat{h}_{n}^{y} + \sum_{n=1, \dots, n} h_{n}^{x} \right) - \widehat{\mathrm{E}}f\left(\sum_{n=1} h_{n}^{x}\right) \right|$$

$$(2.11)$$

$$(2.11)$$

$$\begin{array}{c|cccc}
49 & (3.11) & | & \backslash & x \in \mathbf{Z}^d \setminus \{y, z\} & / & \backslash x \in \mathbf{Z}^d & / & | & 49 \\
50 & & & & \\
51 & & & & \leq \frac{A}{a} \left( |L_n^y|^3 + |L_n^z|^3 \right) . & 51 \\
\end{array}$$

$$\leq \frac{1}{6} \left( |L_n^y|^3 + |L_n^z|^3 \right).$$
 51

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

Because  $\sum_{x \in \mathbb{Z}^d} h_n^x = H_n$ , it is now possible to see how we can iterate the previous inequality to find that P-a.s.,

(3.12) 
$$\left| \widehat{\mathrm{E}}f\left(\sum_{x\in\mathbf{Z}^d}\widehat{h}_n^x\right) - \widehat{\mathrm{E}}f(H_n) \right| \le \frac{A}{6}\sum_{y\in\mathbf{Z}^d} |L_n^y|^3.$$

# We take expectations and appeal to Lemma 2.2 to finish.

# Next, we prove Theorem 1.1.

Proof of Theorem 1.1. We choose, in Lemma 3.1, the collection  $\{\xi_x\}_{x\in \mathbb{Z}^d}$  to be i.i.d. mean-zero normals with variance  $\sigma^2$ . Then, we apply Lemma 3.1 with  $f(x) := g(x/n^{1/2})$  for a smooth bounded function g with bounded derivatives. This yields,

(3.13) 
$$\left| \operatorname{E}g(H_n/n^{1/2}) - \operatorname{E}g\left(\frac{1}{n^{1/2}}\sum_{x\in\mathbf{Z}^d}\widehat{h}_n^x\right) \right| \le \frac{\operatorname{const}}{n^{1/2}}.$$

.

In this way,

$$\sum_{x \in \mathbf{Z}^d} \widehat{h}_n^x \stackrel{\mathcal{D}}{=} \frac{\sigma}{\sqrt{2}} \left| \sum_{x \in \mathbf{Z}^d} L_n^x \left( L_n^x - 1 \right) \right|^{1/2} N(0, 1) \quad \text{under } \widehat{\mathbf{P}}$$

$$= \frac{\sigma}{\sqrt{2}} \left| -n + \sum_{x \in \mathbf{Z}^d} (L_n^x)^2 \right|' \quad N(0, 1),$$
<sup>25</sup>
<sup>26</sup>
<sup>27</sup>
<sup>26</sup>
<sup>27</sup>

where  $\stackrel{\mathcal{D}}{=}$  denotes equality in distribution, and N(0,1) is a standard normal random variable under  $\widehat{P}$  as well as P. Therefore, in accord with Theorem 2.4,

(3.15) 
$$\frac{1}{n^{1/2}} \sum_{x \in \mathbf{Z}^d} \hat{h}_n^x \stackrel{\mathcal{D}}{=} \frac{\sigma}{\sqrt{2}} \left| -1 + \frac{1}{n} \sum_{x \in \mathbf{Z}^d} (L_n^x)^2 \right|^{1/2} N(0, 1)$$

(3.15) 
$$\overline{n^{1/2}} \sum_{x \in \mathbf{Z}^d} h_n^x \stackrel{=}{=} \sqrt{2} \left| -1 + \frac{1}{n} \sum_{x \in \mathbf{Z}^d} (L_n^x)^2 \right| \qquad N(0, 1)$$

$$= o_{\widehat{\mathbf{P}}}(1) + \gamma^{1/2} \cdot N(0, \sigma^2),$$

where  $o_{\widehat{\mathbf{P}}}(1)$  is a term that converges to zero as  $n \to \infty$  in  $\widehat{\mathbf{P}}$ -probability a.s. [P]. Equation (3.13) then completes the proof in the transient case.

#### Theorem 1.2 relies on the following "coupled moderate deviation" result.

**Proposition 3.2.** Suppose that S is recurrent. Consider a sequence  $\{\epsilon_j\}_{j=1}^{\infty}$  of nonnegative numbers that satisfy the following:

(3.16) 
$$\lim_{n \to \infty} \epsilon_n^3 n \left| \sum_{k=1}^n \mathbf{P} \{ S_k = 0 \} \right|^2 = 0.$$

Then for all compactly supported functions  $f: \mathbf{R}^d \to \mathbf{R}$  that are infinitely differen-tiable. 

$$\lim_{51} (3.17) \qquad \lim_{n \to \infty} |\mathbf{E} \left[ f(\epsilon_n W_n) \right] - \mathbf{E} \left[ f(\epsilon_n H_n) \right] | = 0, \qquad 51$$

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

*Proof.* We apply Lemma 3.1 with the  $\xi_x$ 's having the same common distribution as  $q_1$ , and with  $f(x) := g(\epsilon_n x)$  for a smooth and bounded function g with bounded derivatives. This yields, З  $\left| \mathbf{E} \left| g \left( \epsilon_n \sum_{x \in \mathbf{Z}^d} \left| L_n^x \left( L_n^x - 1 \right) \right|^{1/2} Z(x) \right) \right| - \mathbf{E} \left[ g \left( \epsilon_n H_n \right) \right] \right|$ (3.18) $\leq 2C_* M_g \mu_6 n \epsilon_n^3 \left| \sum_{l=0}^n \mathbf{P} \{ S_k = 0 \} \right|^2$ = o(1),owing to Lemma (3.4). According to Taylor's formula,  $g\left(\epsilon_n \sum_{x \in \mathbf{Z}^d} \left|L_n^x \left(L_n^x - 1\right)\right|^{1/2} Z(x)\right)$ (3.19) $=g\left(\epsilon_n\sum_{x\in\mathbf{Z}^d}Z(x)L_n^x\right)+\epsilon_n\sum_{x\in\mathbf{Z}^d}\left(|L_n^x\left(L_n^x-1\right)|^{1/2}-L_n^x\right)Z(x)\cdot R,$ where  $|R| \leq \sup_{x \in \mathbf{R}^d} |g'(x)|$ . Thanks to (2.2), we can write the preceding as follows:  $g\left(\epsilon_n \sum_{x \in \mathbf{Z}^d} \left|L_n^x \left(L_n^x - 1\right)\right|^{1/2} Z(x)\right) - g\left(\epsilon_n W_n\right)$ (3.20) $= \epsilon_n \sum_{n=2^d} \left( |L_n^x (L_n^x - 1)|^{1/2} - L_n^x \right) Z(x) \cdot R.$ Consequently, P-almost surely, ~~  $\left| \widehat{\mathbf{E}} \left| g \left( \epsilon_n \sum_{x \in \mathbf{Z}^d} |L_n^x \left( L_n^x - 1 \right)|^{1/2} Z(x) \right) \right| - \widehat{\mathbf{E}} \left[ g \left( \epsilon_n W_n \right) \right] \right|$ (3.21) $\leq \sup_{x \in \mathbf{R}^d} |g'(x)| \sigma \cdot \epsilon_n \left\{ \widehat{\mathrm{E}} \left( \sum_{x \in \mathbf{R}^d} \left( |L_n^x (L_n^x - 1)|^{1/2} - L_n^x \right)^2 \right) \right\}^{1/2}.$ We apply the elementary inequality  $(a^{1/2}-b^{1/2})^2 \leq |a-b|$  —valid for all  $a,b \geq 0$  —to deduce that P-almost surely,  $\left|\widehat{\mathbf{E}}\left|g\left(\epsilon_{n}\sum_{x\in\mathbf{Z}^{d}}\left|L_{n}^{x}\left(L_{n}^{x}-1\right)\right|^{1/2}Z(x)\right)\right|-\widehat{\mathbf{E}}\left[g\left(\epsilon_{n}W_{n}\right)\right]\right|$ **、 、** 1/2 

$$\begin{array}{l} 47 \\ 48 \\ 49 \\ 50 \end{array} \leq \sup_{x \in \mathbf{R}^d} |g'(x)| \sigma \cdot \epsilon_n \left\{ \widehat{\mathbf{E}} \left( \sum_{x \in \mathbf{Z}^d} L_n^x \right) \right\}^{1/2} \\ 1/2 \\ 1/2 \end{array}$$

$$= \sup_{x \in \mathbf{R}^d} |g'(x)| \sigma \cdot \epsilon_n n^{1/2}.$$

imsart-coll ver. 2008/08/29 file: Khoshnevisan.tex date: March 25, 2009

We take E-expectations and apply Lemma (3.4) to deduce from this and (3.18) that  $|\mathbf{E}[g(\epsilon_n W_n)] - \mathbf{E}[g(\epsilon_n H_n)]| = o(1).$ (3.23)З З This completes the proof. Our proof of Theorem 1.2 hinges on two more basic lemmas. The first is an elementary lemma from integration theory. **Lemma 3.3.** Suppose  $X := \{X_n\}_{n=1}^{\infty}$  and  $Y := \{Y_n\}_{n=1}^{\infty}$  are  $\mathbb{R}^d$ -valued random variables such that: (i) X and Y each form a tight sequence; and (ii) for all bounded infinitely-differentiable functions  $q: \mathbf{R}^d \to \mathbf{R}$ ,  $\lim_{n \to \infty} |\mathrm{E}g(X_n) - \mathrm{E}g(Y_n)| = 0.$ (3.24)Then, the preceding holds for all bounded continuous functions  $q: \mathbf{R}^d \to \mathbf{R}$ . *Proof.* The proof uses standard arguments, but we repeat it for the sake of com-pleteness. Let  $K_m := [-m, m]^d$ , where m takes values in **N**. Given a bounded continuous function  $g: \mathbf{R}^d \to \mathbf{R}$ , we can find a bounded infinitely-differentiable function  $h_m: \mathbf{R}^d \to \mathbf{R}$  such that  $|h_m - g| < 1/m$  on  $K_m$ . It follows that  $|\mathrm{E}g(X_n) - \mathrm{E}g(Y_n)| \le 2/m + |\mathrm{E}h_m(X_n) - \mathrm{E}h_m(Y_n)|$ (3.25) $+ 2 \sup_{x \in \mathbf{R}^d} |g(x)| ( \mathbb{P}\{X_n \notin K_m\} + \mathbb{P}\{Y_n \notin K_m\}).$ Consequently,  $\limsup |\mathrm{E}g(X_n) - \mathrm{E}g(Y_n)|$  $n \rightarrow \infty$ (3.26) $\leq 2/m + 2 \sup_{x \in \mathbf{R}^d} |g(x)| \sup_{j \geq 1} \left( \mathbf{P}\{X_j \notin K_m\} + \mathbf{P}\{Y_j \notin K_m\} \right).$ Let m diverge and appeal to tightness to conclude that the left-had side vanishes. The final ingredient in the proof of Theorem 1.1 is the following harmonic-analytic result. **Lemma 3.4.** If  $\epsilon_n := 1/a_n$ , then (3.16) holds. *Proof.* Let  $\phi$  denote the characteristic function of  $S_1$ . Our immediate goal is to prove that  $|\phi(t)| < 1$  for all but a countable number of  $t \in \mathbf{R}^d$ . We present an argument, due to Firas Rassoul-Agha, that is simpler and more elegant than our original proof. Suppose  $S'_1$  is an independent copy of  $S_1$ , and note that whenever  $t \in \mathbf{R}^d$  is such that  $|\phi(t)| = 1$ ,  $D := \exp\{it \cdot (S_1 - S'_1)\}$  has expectation one. Consequently,  $E(|D-1|^2) = E(|D|^2) - 1 = 0$ , whence D = 1 a.s. Because  $S_1$  is assumed to have at least two possible values,  $S_1 \neq S'_1$  with positive probability, and this proves that  $t \in 2\pi \mathbf{Z}^d$ . It follows readily from this that 

49 (3.27) 
$$\{t \in \mathbf{R}^d : |\phi(t)| = 1\} = 2\pi \mathbf{Z}^d,$$
 50 50

and in particular,  $|\phi(t)| < 1$  for almost all  $t \in \mathbf{R}^d$ .

Chen and Khoshnevisan By the inversion theorem (Spitzer (1976) [**P3**(b), p. 57]), for all  $n \ge 0$ ,  $P\{S_n = 0\} = \frac{1}{(2\pi)^d} \int_{(-\pi,\pi)^d} \{\phi(t)\}^n dt.$ (3.28)This and the dominated convergence theorem together tell us that  $P\{S_n = 0\} =$ o(1) as  $n \to \infty$ , whence it follows that  $\sum_{k=1}^{n} \mathbb{P}\{S_k = 0\} = o(n) \quad \text{as } n \to \infty.$ (3.29)For our particular choice of  $\epsilon_n$  we find that  $\epsilon_n^3 n \left| \sum_{k=1}^n \mathbf{P}\{S_k = 0\} \right|^2 = \left( \frac{1}{n} \sum_{k=1}^n \mathbf{P}\{S_k = 0\} \right)^{1/2},$ (3.30)and this quantity vanishes as  $n \to \infty$  by (3.29). This proves the lemma. *Proof of Theorem 1.2.* Let  $\epsilon_n := 1/a_n$ . In light of Proposition 3.2, and Lemmas 3.3 and 3.4, it suffices to prove that the sequences  $n \mapsto \epsilon_n W_n$  and  $n \mapsto \epsilon_n H_n$  are tight. Lemma 2.2, (2.2), and recurrence together imply that for all n large,  $\mathbf{E}\left(\left|\epsilon_{n}W_{n}\right|^{2}\right) = \sigma^{2}\epsilon_{n}^{2}\sum_{x \in \mathbf{Z}^{d}}\mathbf{E}\left(\left|L_{n}^{x}\right|^{2}\right)$  $\leq \operatorname{const} \cdot \epsilon_n^2 n \sum_{k=1}^n \mathrm{P}\{S_k = 0\}$ (3.31)= const.Thus,  $n \mapsto \epsilon_n W_n$  is bounded in  $L^2(\mathbf{P})$ , and hence is tight. We conclude the proof by verifying that  $n \mapsto \epsilon_n H_n$  is tight. Thanks to (2.4) and recurrence, for all n large,  $\mathbb{E}\left(|\epsilon_n H_n|^2\right) \leq \operatorname{const} \cdot \epsilon_n^2 \mathbb{E}\sum_{x \in \mathbf{Z}^d} (L_n^x)^2$  $\leq \operatorname{const} \cdot \epsilon_n^2 n \sum_{k=1}^n \mathrm{P}\{S_k = 0\}$ (3.32)= const.Confer with Lemma 2.2 for the penultimate line. Thus,  $n \mapsto \epsilon_n H_n$  is bounded in  $L^{2}(\mathbf{P})$  and hence is tight, as was announced. Acknowledgement. We wish to thank Siegfried Hörmann, Richard Nickl, Jon Peterson, and Firas Rassoul-Agha for many enjoyable discussions, particularly on the first portion of

Rassoul-Agha for many enjoyable discussions, particularly on the first portion of
Lemma 3.4. We are grateful to Nadine Guillotin–Plantard for her valuable suggestions and typographical corrections, as well as providing us with references to her
paper (Guillotine–Plantard (2004)). Special thanks are extended to Firas RassoulAgha for providing us with his elegant argument that replaced our clumsier proof
of the first part of Lemma 3.4. Finally, we thank two anonymous referees who made
valuable suggestions and pointed out misprints.

| <ol> <li>BOLTHANSEN, E. (1989). A central limit theorem for two-dimensional random walks in random scener-<br/>iss. Ann. Probab., 17.1, 108–115.</li> <li>BUTET, E. and PULL, J. V. (1977). A model of continuous polymers with random charges. J. Math.<br/>Phys. 38, 5143–512.</li> <li>Churk, N. (2008). Limit laws for the energy of a charged polymer. Annales de l'Institut Henri<br/>Poincaré: Probab. et Statist., 44, 638–672.</li> <li>Churk, K. L. and FOULS, W. H. J. (1951). On the distribution of values of sums of random variables.<br/>Mem. Amer. Math. Soc., 6.</li> <li>DERRIM, B. E. and HEGNS, R. G. (1992). A model of directed walks with random<br/>self interactions. J. Phys. A 27, 558–5493.</li> <li>GARE, T. and ORLAND, H. (1988). Mean-field model for protein folding. Europhys. Lett., 6, 307–3301.</li> <li>GARE, T. and ORLAND, H. (1988). Mean-field model for protein folding. Europhys. Lett., 4, 307–301.</li> <li>GULTON-PLAYNAND, N. (2004). Sur la convergence failed eas systèmes dynamiques chantilleunes.<br/>Annales de l'Institut Fourier, 54.1, 255–278.</li> <li>KANTON, Y. and KARAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 4, 421–426.</li> <li>KANTON, Y. and KARAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 4, 421–426.</li> <li>KANTON, Y. and KARAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 4, 421–426.</li> <li>KANTON, Y. and KARAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 4, 421–426.</li> <li>KANTON, Y. and KARAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 4, 421–426.</li> <li>KANTON, Y. and KARAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 4, 421–426.</li> <li>KARTON, Y. and KARAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 4, 421–426.</li> <li>MARTINEZ, S. and PITTERS, D. (1996). Thermodynamics of a Brownian bridge polymer model in a<br/>random environment. J. Phys. A, 29, 1267–1279.</li> <li>KARTONE, Y. (1993). Proveples of Random Walks, seco</li></ol>   | 1  | References   | 1  |
|--|----|--|----|
| <ol> <li>BOLTAUSSE, E. (1989). A control multi theorem for two-dimensional random while in Fandom scener-<br/>ies. Ann. Probab. 171, 108-118.</li> <li>B. B. M., M. PULS, J. V. (1997). A model of continuous polymers with random charges. J. Math.<br/>50, 108, N. (2008). Limit have for the energy of a charged polymer. Annales de l'Institut Henri<br/>Porncorri. Probab. et Statist., 44, 638-672.</li> <li>C. CHUNG, K. L. and PUCHS, W. H. J. (1951). On the distribution of values of sums of random variables.<br/>70, 400, K. L. and PUCHS, W. H. J. (1951). On the distribution of values of sums of random variables.<br/>71, 400, K. L. and PUCHS, W. H. J. (1951). On the distribution of values of sums of random variables.<br/>73, 400, K. L. and PUCHS, W. H. J. (1994). Low-temporature properties of directed walks with random<br/>74, 51, 51, 51, 51, 51, 51, 51, 51, 51, 51</li></ol>   | 2  |  | 2  |
| <ol> <li>BETTET, E. and PULÉ, L.Y. (1997). A model of continuous polymers with random charges. J. Math.<br/>Phys. 85, 5143-5152.</li> <li>CHEN, N. (2008). Limit have for the energy of a charged polymer. Annales de l'Institut Henri<br/>Poincari. Probab. et Slatist, 44, 638-672.</li> <li>CHUNG, K. L. and FUCHS, W. H. J. (1951). On the distribution of values of sums of random variables.<br/>Mem. Amer. Math. Soc., 6.</li> <li>DERRINA, B., GRIFTTINS, B. B. and HGOS, R. G. (1992). A model of directed walks with random<br/>interactions. Europhys. Lett., 18, 361-366.</li> <li>DERRINA, B. and HHOS, R. G. (1994). Loos-temperature properties of directed walks with random<br/>interactions. Annal. M. (1998). Mean-field model for protein folding. Europhys. Lett., 6, 307-310.</li> <li>G. CAURS, T. and ORLAND, H. (1988). Mean-field model for protein folding. Europhys. Lett., 18, 421-426.</li> <li>KAURON, Y. and Kamata, M. (1990). Polymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>KAURAN, Y. and Kamata, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>KANTOA, Y. and Kamata, M. (1990). Polymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>KANTOA, Y. and Kamata, M. (1990). Sur mas proposition de la théorie des probabilités. Eul. de'l Acad.<br/>Interpreted. de'or. 2023). Einsteinung 154, 363-394.</li> <li>LAROUNKO, A. M. (1990). Sur mas proposition de la théorie des probabilités. Eul. de'l Acad.<br/>Interpreted. de'RETURE, F. (1993). A limit theorem related to a new class of self-similar processes.</li> <li>Maritykez, S. and PETERTS, D. (1996). Thormodynamics of a Brownian bridge polymer model in a<br/>random environment. J. Phys. A. 29, 1207-1279.</li> <li>KWarrich, Verw. Gehrich, 500, 5-23.</li> <li>OUKINOV, S. P. (1996). Thormodynamics of a Brownian bridge polymer model in a<br/>random environment. J. Phys. A. 29, 1207-1279.</li> <li>KWarrich, Verw. Gehrich, 500, 5-25.</li> <li>OUKINOV, S. P. (1996). Thormodynamics of a Brownian</li></ol>                             | 3  | [1] BOLTHAUSEN, E. (1989). A central limit theorem for two-dimensional random walks in random scener-<br>ies. Ann. Probab., <b>17.1</b> , 108–115.   | 3  |
| <ol> <li>Phys. 38, 5143-5152.</li> <li>Cluck, X. (2008). Limit laws for the energy of a charged polymer. Annales de l'Institut Henri Poincard: Probab. et Statist., 44, 638-672.</li> <li>Cluxo, K. L. and Purts, W. H. J. (1951). On the distribution of values of sums of random variables. Mem. Amer. Math. Soc., 6.</li> <li>Dizuno, B. C. aururtus, R. B. and Hörös, R. C. (1992). A model of directed walks with random interactions. Europhys. Lett., 18, 361-366.</li> <li>Dizuno, J. C. Run, A. 27, 5485-5433.</li> <li>C. CLURT, T. and OLARD, H. (1988). Mean-field model for protein folding. Europhys. Lett., 6, 307-310.</li> <li>G. CLURT, T. and OLARD, H. (1988). Mean-field model for protein folding. Europhys. Lett., 6, 307-310.</li> <li>G. KANTOR, V. and KAMDAR, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-420.</li> <li>KANTOR, V. and KAMDAR, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-420.</li> <li>KANTOR, V. and KAMDAR, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-420.</li> <li>Larousawi, A. M. (1990). Sen use proposition de la théorie des probabilités. Bull. de'l Acad. Imperiade des Sci. St. Petersbourg 13.4, 350-366.</li> <li>Larousawi, A. M. (1990). Four met proposition de la théorie des probabilités. Bull. de'l Acad. Imperiade des Sci. St. Petersbourg 13.4, 350-366.</li> <li>Larousawi, A. M. (1990). Sen use proposition de la théorie des probabilités. Bull. de'l Acad. Imperiade des Sci. St. Petersbourg 13.4, 350-366.</li> <li>Larousawi, A. M. (1990). Sen use proposition de la théorie des probabilités. Bull. de'l Acad. Imperiade des Sci. St. Petersbourg 13.4, 360-366.</li> <li>Marimus, S. and Parrize, F. (1975). Thermodynamics of a Brownian bridge polymer model in a 199 random environment. J. Phys. A, 29, 1267-1279.</li> <li>Kabreak, Hann Serrize, F. (1975). A limit theorem related to a new class of self-similar processes. 21, 20, Meanema et al. 199-160.</li> <li>Strizemente, Maid</li></ol>   | 4  | [2] BUFFET, E. and PULÉ, J. V. (1997). A model of continuous polymers with random charges. J. Math.  | 4  |
| <ul> <li>[a] Chink, A. (2009). Limit: have not inte energy of a charged polymer. Anomals in Private 1999.</li> <li>[b] Chink, K. L. and FUCHS, W. H. J. (1951). On the distribution of values of sums of random variables. Man. Marc. Math. Soc., 6</li> <li>[c] DERRIN, B., GRIPTINE, R. E. and HIGES, R. G. (1992). A model of directed walks with random interactions. Envirophys. Lett., 18, 301-366.</li> <li>[c] DERRIN, B. and HIGES, R. G. (1994). Low-temperature properties of directed walks with random self interactions. J. Phys. A. 27, 1485-343.</li> <li>[c] GARE, T. and ORLAND, H. (1988). Mean-field model for protein folding. Europhys. Lett., 6, 307-310.</li> <li>[s] GURDEN-PARTINEN, N. (2004). Sur la convergence fuble des systemes dynamicus channel londs: Annales de l'Institut Fourier, 541, 265-278.</li> <li>[NKANON, Y. and Kamana, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>[10] KINENN, A. M. (1992). Sur une proposition de la théorie des probabilités. Eull. de'l Acad. Imprivale des Sci. 189, 001, 52-53.</li> <li>[11] LANDUNN, A. M. (1992). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeit ascietung. Math. 21, 52-11-29.</li> <li>[12] LUNDERBEG, W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeit ascietung. Math. 21, 52-11-29.</li> <li>[13] MARTINZ, S. and PERTER, F. (1978). A 1101 theorem related to a new class of self-similar processes. Z. Z. Wahrsch. Vers. Gebetz, 60, 15, -23.</li> <li>[14] KINTRE, R. and SPITERE, F. (1978). A 1101 theorem related to a new class of self-similar processes. Z. Z. Wahrsch. Vers. Gebetz, 60, 15, -23.</li> <li>[15] OUNKINN, S. F. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 1365-3664.</li> <li>[16] FÖLLA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betrefiend die Irrfahrt in STITZEN, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg, F. (1979). For</li></ul>   | 5  | Phys. 38, 5143-5152.   | 5  |
| <ol> <li>CHUNG, K. L. and FUCIS, W. H. J. (1951). On the distribution of values of sums of random variables.<br/>Mem. Anew. Anth. Soc., 6</li> <li>DERRINA, B., GREFTINS, R. B. and HUGS, R. G. (1992). A model of directed walks with random<br/>interactions. Everyphys. Lett., 18, 301-366.</li> <li>DERRINA, B., and HUGS, R. G. (1994). Low-temperature properties of directed walks with random<br/>self interactions. J. Phys. A, 27, 1485-5493.</li> <li>GAREL, T. and ORLAND, H. (1988). Mean-field model for protein folding. Everyphys. Lett., 60, 307-310.</li> <li>GULTOTN-PLANTARD, N. (2004). Sur la coveregence faible des systèmes dynamiques échantillomés.<br/>Annales de l'Institut Fourier, 54.1, 255-278.</li> <li>KANTOR, Y. and KALDALM, M. (1991). Folymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>KESTEN, H. and SPUTZER, F. (1979). A limit theorem related to a new class of self-similar processes.<br/>Z. Wahrsch. Veru. Gebiets, 65, 01, 5-25.</li> <li>LANDERBAG, W. (1922). Eine neue Hereliung des Exponentialgesetzes in der Walnscheinlichkeit-<br/>sechnung, Mat. X., 15, 211-225.</li> <li>MARTNEZ, S. and PETERR, P. (1978). A limit theorem related to a new class of self-similar processes.<br/>Z. Wahrsch. Veru. Gebiets, 65, 01, 5-25.</li> <li>MARTNEZ, S. and PETERR, F. (1978). A limit theorem related to a new class of self-similar processes.<br/>Z. Wahrsch. Veru. Gebiets, 65, 01, 5-25.</li> <li>OUKKNOY, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19,<br/>3655-3664.</li> <li>FUTZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-<br/>Heidelberg.</li> <li>WITTNER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> <li>WITTNER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> </ol>   | 6  | [5] CHEN, X. (2008). Limit laws for the energy of a charged polymer. Annues at i Institut Henri<br>Poincaré: Probab. et Statist., 44, 638–672.   | 6  |
| <ol> <li>Menn, Amer. Math. Soc., 6.</li> <li>DERRING, B., GRENTINS, R. B. and HUGS, R. G. (1992). A model of directed walks with random interactions. Europhys. Lett., 18, 301-306.</li> <li>DERRING, B., and HUGS, R. G. (1994). Low-temperature properties of directed walks with random self interactions. J. Phys. A, 27, 5485-5493.</li> <li>GAREL, T. and OLAND, H. (1984). Mean-field model for protein folding. Europhys. Lett., 16, 307-310.</li> <li>GULLOTN-PLANTARD, N. (2004). Sur la convergence faible des systèmes dynamiques échantillonnés. Annales de l'Institut Fourier, 541, 255-278.</li> <li>KANTOR, Y. and KAUDAM, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>KISTEN, H. and SPITZER, F. (1997). A limit theorem related to a new dass of self-similar processes. Z. Wahrsch. Verw. Gebiete, 501, 5-25.</li> <li>LLONDINO, A. M. (1902). Eure use Heristiung des Exponentialgesetzes in der Wahrscheinlichkeitsrechtnung. Math. Z., 15, 211-225.</li> <li>JLONEDEREN, W. 1992). Leine use Heristiung des Exponentialgesetzes in der Wahrscheinlichkeitsrechtnung. Math. Z., 15, 211-225.</li> <li>MARTINZ, S. And PETRITE, D. (1996). Thermodynamics of a Brownian bridge polymer model in a random environment. J. Phys. A, 19, 1267-1279.</li> <li>KISTEN, H. and PETRITE, P. (1997). A limit theorem related to a new class of self-similar processes. Z. Wahrsch. Verw. Gebiete, 501, 5-25.</li> <li>OUKKHOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.</li> <li>Poiva, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Straßennetz. Math. Ann., 84, 149-160.</li> <li>YETZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.</li> <li>Mitther M. 24, 263-268.</li> <li>GUKHOV, S. P. (1996). Thereiples of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.</li> <li>Mitther M. 24, 263-268.</li> <li>GUKHOV, S</li></ol>   | 7  | [4] CHUNG, K. L. and FUCHS, W. H. J. (1951). On the distribution of values of sums of random variables.  | 1  |
| <ul> <li>[9] DERMAR, B., and FRUZE, R. G. (1994). Low-temperature properties of directed walks with random interactions. Europhys. Lett., 18, 361-366.</li> <li>[9] DERMAR, B. and HEGS, R. G. (1994). Low-temperature properties of directed walks with random inself interactions. J. Phys. A. 27, 5485-5483.</li> <li>[9] GAREL, T. and ORLAND, R. (1998). Mean-field model for protein folding. Europhys. Lett., 307-310.</li> <li>[10] CAREL, T. and ORLAND, R. (1998). Near-field model for protein folding systems dynamicus cichantillonics. Interactions Autonom, N. (1990). Sur la convergence faible desystems dynamicus cichantillonics. Annales de Ul'institut Fourier, 54.1, 255-278.</li> <li>[10] KANTOR, Y. and KANDAH, M. (1990). Sur la convergence faible desystems dynamicus cichantillonics. Interactions. Europhys. Lett., 14, 421-426.</li> <li>[11] KANTOR, Y. and KANDAH, M. (1990). Thermodynamics of a new class of self-similar processes. I. Z. Wadrach. Veru. Gebiets, 650, 15-25.</li> <li>[12] LANDERBAG, W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitserschnung, Math. Z., 15, 211-225.</li> <li>[13] MAIRINEZ, S. and PETIURS, D. (1996). Thermodynamics of a Brownian bridge polymer model in a random environment. J. Phys. A. 29, 1267-1279.</li> <li>[14] KESTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes. Z. Wahrsch. Veru. Gebiets, 650, 15, -25.</li> <li>[15] OUKKIOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 365-3664.</li> <li>[16] PÓIXA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Strafementz. Math. Ann., 84, 149-160.</li> <li>[17] SPITZER, F. (1976). Principles of Random Waks, second edition. Springer-Verlag, New York-Heidelberg.</li> <li>[18] WITTMER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.</li> <li>[29]</li> <li>[20]</li> <li>[</li></ul>   | 8  | Mem. Amer. Math. Soc., 6.<br>[5] DEPRIDA B. CRIEFITHE B. B. and HICCE B. C. (1002) A model of directed walks with random   | 8  |
| <ol> <li>[8] DERRIDA, B. and HIGUS, R. G. (1994). Low-temperature properties of directed walks with random<br/>self interactions. J. Phys. A. 27, 5485-5493.</li> <li>[17] GAREE, T. and ORLAND, H. (1988). Mean-field model for protein folding. Europhys. Lett., 6, 307-310.</li> <li>[18] GULLOTN-PLANTARD, N. (2004). Sur La convergence faible des systèmes dynamiques échantillonnés.<br/>Annales de l'Institut Fourier, 54.1, 255-278.</li> <li>[19] KANTOR, Y. and KANDAR, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>[10] KENTEN, H. and SEPTZER, F. (1979). A limit theorem related to a new class of self-similar processes.<br/>I. Wahrsch. Verw. Gebiete, 501, 5-25.</li> <li>[11] LANDUNOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.<br/>Impériade des 561. SP. Petrénourg 134, 330-386.</li> <li>[12] LANDUNOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.<br/>Impériade des 561. SP. Petrénourg 134, 330-386.</li> <li>[13] Matribaz, S. and PETZER, F. (1978). A limit theorem related to a new class of self-similar processes.<br/>Z. Wahrsch. Verw. Gebiete, 501, 5-25.</li> <li>[14] KISTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.<br/>Z. Wahrsch. Verw. Gebiete, 501, 5-26.</li> <li>[15] DUNLENG, S. F. (1986). Comfigurational statistics of a disordered polymer chain. J. Phys. A, 19,<br/>3635-3644.</li> <li>[16] FOINA, G. (1921). Über eine Anfgabe der Wahrscheinlichkeitsrechnung hetreffend die Irrfahrt in<br/>Straßennetz. Math. Ann., 84, 149-160.</li> <li>[17] STITZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-<br/>Heidelberg.</li> <li>[18] WITTER, J. JOINER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> <li>[29]</li> <li>[20] Andren Serlard, Ser</li></ol>                 | 9  | interactions. Europhys. Lett., 18, 361–366.  | 9  |
| <ol> <li>self interactions. J. Phys. A, 27, 5485–5493.</li> <li>[7] GARLE, T. and OURAND, H. (1988). Mean-field model for protein folding. Europhys. Lett., 6, 307–310.</li> <li>[8] GUILLOTN-PLANTARD, N. (2004). Sur Ia convergence fable des systèmes dynamiques échantillonnés.<br/>Arnales de l'INSTRUE, H. and SETTER, F. (1979). A limit theorem related to a new class of self-similar processes.<br/>I. KANTON, Y. and KARDAM, M. (1901). Polymers with self-interactions. Europhys. Lett., 14, 421–426.</li> <li>[10] KANTON, Y. and KARDAM, M. (1991). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.<br/>Impériade des Sci. St. Petershourg 13.4, 359–386.</li> <li>[12] LINDUNON, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.<br/>Impériade des Sci. St. Petershourg 13.4, 359–386.</li> <li>[13] Martínez, S. and PETRITS, D. (1996). Thermodynamics of a Brownian bridge polymer model in a<br/>random environment. J. Phys. A, 19, 1267–1279.</li> <li>[14] KENTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.<br/>I. Wahrsch. Verw. Gebiete, 50.1, 5–25.</li> <li>[15] ODENENOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19,<br/>3655–3664.</li> <li>[16] PÓIMA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrlahrt in<br/>Straßennetz. Math. Ann., 84, 149–160.</li> <li>[17] SPITZER, F. (1976). Principles of Random Waks, second edition. Springer-Verlag, New York-<br/>Heidelberg.</li> <li>[18] WITTMER, J., JOINER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263–268.</li> <li>[29]</li> <li>[20]</li> <li>[21] SURER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263–268.</li> <li>[22]</li> <li>[23] Mathing A. Songer A. And JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263–268.</li> <li>[23]</li> <li>[24] Mathing A</li></ol>                         | 10 | [6] DERRIDA, B. and HIGGS, R. G. (1994). Low-temperature properties of directed walks with random  | 10 |
| <ol> <li>GULLOTN-PLANTARD, N. (2004). Sur las convergence fable des systèmes dynamiques échantillomés.<br/>Annales de l'Institut Fourier, 54.1, 255–278.</li> <li>KARYOR, Y. and KARDAR, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421–426.</li> <li>KERNE, H. and SPITZER, F. (1973). A limit theorem related to a new class of self-similar processes.<br/>E. Wahrsch. Verw. Gebtec. 60.1, 5-25.</li> <li>LAPOUSOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.<br/>Impériade des C.S. St. Petersony 13.4, 350–386.</li> <li>LINDEBEG, W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeit-<br/>stechnung. Math. Z., 15, 211–225.</li> <li>MARTINZ, S. and PERTIN, D. (1996). Thermodynamics of a Brownian bridge polymer model in a<br/>random environment. J. Phys. A, 29, 1267–1279.</li> <li>KENTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.<br/>E. Wahrsch. Verw. Gebiete, 501, 5-25.</li> <li>OBUKHOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19,<br/>3655–3684.</li> <li>PóuxA, G. (1921). Übre eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in<br/>Straßennetz. Math. Ann., 84, 149–160.</li> <li>Straßennetz. Math. Ann., 84, 149–160.</li> <li>Straßennetz. Math. Ann., 84, 149–160.</li> <li>WITTMER, J., JOINNE, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> <li>WITTMER, J., JOINNE, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> </ol>  | 11 | self interactions. J. Phys. A, 27, 5485–5493.<br>[7] CAREL T and ORLAND H (1988) Mean-field model for protein folding. Europhys. Lett. 6, 307–310.   | 11 |
| <ul> <li>Annales de l'Institut Fourier, 54.1, 255-278.</li> <li>[9] KANTOR, Y. and KARDAR, M. (1991). Folymers with self-interactions. Europhys. Lett., 14, 421-426.</li> <li>[10] KESTEN, H. and SPITZER, F. (1979). A limit theorem related to a new class of self-similar processes.</li> <li>[11] LIAPOUNOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.</li> <li>[12] LINDEBERG, W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitsechtung. Math. Z., 15, 211-225.</li> <li>[13] MARTINEZ, S. and PETTRIN, D. (1960). Thermodynamics of a Brownian bridge polymer model in a random environment. J. Phys. A, 29, 1267-1279.</li> <li>[14] KESTEN, H. and SUTZER, F. (1978). A limit theorem related to a new class of self-similar processes. Z. Wahrsch. Verw. Gebiete, 501, 5-25.</li> <li>[15] OUKKINO, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.</li> <li>[16] POIXA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irfahrt in Straßennetz. Math. Ann., 84, 149-160.</li> <li>[17] SPTZR, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelborg.</li> <li>[18] WITTMER, J., JONEER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.</li> <li>[29]</li> <li>[20]</li> <li>[21] KESTEN, H. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.</li> <li>[22]</li> <li>[23]</li> <li>[24] MITTMER, J., JONEER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.</li> <li>[23]</li> <li>[24]</li> <li>[25]</li> <li>[25]</li> <li>[26]</li> <li>[27]</li> <li>[28]</li> <li>[28]</li> <li>[29]</li> <li>[29]</li> <li>[29]</li> <li>[20]</li> <li>[20]</li> <li>[21]</li> <li>[22]</li> <li>[22]</li> <li>[23]</li> <li>[23]</li> <li>[24]</li> <li>[25]</li> <li>[25]</li> <li>[26]</li> <li>[27]</li> <li>[28]</li> <li>[28]</li> <li>[29]</li> <li>[29]</li> <li>[29</li></ul> | 12 | <ul><li>[8] GUILLOTIN-PLANTARD, N. (2004). Sur la convergence faible des systèmes dynamiques échantillonnés.</li></ul>   | 12 |
| <ul> <li>[9] KANTOK, Y. and KARDAK, M. (1991). Polymers with self-interactions. Europhys. Lett., 14, 421-420.</li> <li>[10] KANTOK, Y. and KARDAK, M. (1990). For the over probabilities and Soft Self-similar processes. Z. Wahrsch. Verw. Gebietz, 50, 15-25.</li> <li>[11] LIAPOUNOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad. Impériale des Sci. St. Petérsbourg 13.4, 350-386.</li> <li>[12] LINDERER, V. (1922). Eline neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeit serechtnung. Math. Z., 15, 211-225.</li> <li>[13] MARTINEZ, S. and PETETTE, D. (1996). Thermodynamics of a Brownian bridge polymer model in a random environment. J. Phys. A, 29, 1267-1279.</li> <li>[14] KESTER, H. and SUTZER, F. (1973). A limit theorem related to a new class of self-similar processes. Z. Wahrsch. Verw. Gebietz, 501, 5-25.</li> <li>[15] OBUKHOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.</li> <li>[16] PÓIXA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irfahrt in Strafennetz. Math. Ann., 84, 149-160.</li> <li>[17] SPITZER, F. (1976). Principles of Random Wakks, second edition. Springer-Verlag, New York-Heidelberg.</li> <li>[18] WITTMER, J., JORNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.</li> <li>[29]</li> <li>[30]</li> <li>[31]</li> <li>[31]</li> <li>[32]</li> <li>[33]</li> <li>[34]</li> <li>[34]</li> <li>[35]</li> <li>[36]</li> <li>[36]</li> <li>[36]</li> <li>[36]</li> <li>[36]</li> <li>[36]</li> <li>[36]</li> <li>[36]</li> <li>[37]</li> <li>[38]</li> <li>[39]</li> <li>[30]</li> <li>[31]</li> <li>[31]</li> <li>[32]</li> <li>[32]</li> <li>[33]</li> <li>[33]</li> <li>[34]</li> <li>[35]</li> <li>[36]</li> <li>[36]</li> <li>[37]</li> <li>[38]</li> <li>[39]</li> <li>[39]</li> <li>[30]</li> <li>[31]</li> <li>[32]</li> <li>[33]</li> <li>[34]</li> <li>[35]</li> <li>[35]</li> <li>[36]</li> <li>[36]</li> <li>[37]</li> <li>[38]</li> <li>[39]</li> <li></li></ul> | 13 | Annales de l'Institut Fourier, <b>54.1</b> , 255–278.  | 13 |
| <ol> <li>Z. Wahrsch. Verm. Gebict. 50.1, 5–25.</li> <li>LIAPOUNOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.<br/>Impériale des Sci. St. Petérsbourg 13.4, 359–386.</li> <li>LINDERRG, W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeit-<br/>sechtnung. Math. Z., 15, 211–225.</li> <li>MARTINEZ, S. and PETRITIS, D. (1966). Thermodynamics of a Brownian bridge polymer model in a<br/>random environment. J. Phys. A, 29, 1267–1279.</li> <li>KESTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.<br/>Z. Wahrsch. Verm. Gebicte, 50.1, 5–25.</li> <li>ORUKINOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19,<br/>3655–3664.</li> <li>Pótya, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in<br/>Straßennetz. Math. Ann., 84, 149–160.</li> <li>Straßennetz, Math. Ann., 84, 149–160.</li> <li>Straßennetz, Math. Ann., 84, 149–160.</li> <li>Straßennetz, J., JUINER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263–268.</li> <li>WITTMER, J., JUINER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263–268.</li> <li>Matting Mith. Ann., 84, 149–160.</li> <li>Timer and the straff of the straf</li></ol>                     | 15 | <ul> <li>[9] KANTOR, Y. and KARDAR, M. (1991). Polymers with self-interactions. <i>Europhys. Lett.</i>, 14, 421–426.</li> <li>[10] KESTEN, H. and SPITZER, F. (1979). A limit theorem related to a new class of self-similar processes.</li> </ul> | 15 |
| <ol> <li>LIAPOUNOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.<br/>Inpériela des Sci. St. Petersbourg 13.4, 359-386.</li> <li>LINDERERG, W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeit-<br/>srechtnung. Math. Z., 15, 211-225.</li> <li>MARTÍNEZ, S. and PETRITIS, D. (1996). Thermodynamics of a Brownian bridge polymer model in a<br/>random environment. J. Phys. A, 29, 1267-1279.</li> <li>KESTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.<br/>J. Wahrsch. Verw. Gebiete, 50.1, 5-25.</li> <li>OUKUKOV, S. P. (1986). Onfigurational statistics of a disordered polymer chain. J. Phys. A, 19,<br/>3655-3664.</li> <li>DOUKIKOV, S. C. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in<br/>Straßennetz. Math. Ann., 84, 149-160.</li> <li>Stritzen, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-<br/>Heidelberg.</li> <li>WITTINER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> <li>WITTINER, J. JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> <li>WITTINER, J. JOHNER, M. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-<br/>rophys. Lett., 24, 263-268.</li> </ol>  | 16 | Z. Wahrsch. Verw. Gebiete, <b>50.1</b> , 5–25.   | 16 |
| <ul> <li>Imperate des Sc. 15. Petersourg 15.4, 303-386.</li> <li>[12] LINDERERO, V. (1922). Etien neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeitsrechtung, Math. Z., 15, 211-225.</li> <li>[13] MARTINEZ, S. and PETRITTS, D. (1996). Thermodynamics of a Brownian bridge polymer model in a random environment. J. Phys. A, 29, 1267-1279.</li> <li>[14] KESTEN, H. and SPITZER, F. (1975). A limit theorem related to a new class of self-similar processes. Z. Wahrsch. Verw. Gebiete, 50.1, 5-25.</li> <li>[15] OBUKHOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.</li> <li>[16] PÓIXA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Straßennetz. Math. Ann., 84, 149-160.</li> <li>[17] SPITZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.</li> <li>[18] WITTIKER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.</li> <li>[29]</li> <li>[30]</li> <li>[41] MUTTIKER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.</li> <li>[42]</li> <li>[43]</li> <li>[44]</li> <li>[44]</li> <li>[44]</li> <li>[45]</li> <li>[46]</li> &lt;</ul>   | 17 | [11] LIAPOUNOV, A. M. (1900). Sur une proposition de la théorie des probabilités. Bull. de'l Acad.   | 10 |
| 19       I33       Strethtung, Math. Z., 15, 211–225.       19         19       IAMATNEX, S. and PERTRIS, D. (1996). Thermodynamics of a Brownian bridge polymer model in a random environment. J. Phys. A, 29, 1267–1279.       20         21       I4       KESTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.       21         21       I4       KESTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.       21         23       OBUKINOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.       23         24       I16       PóIXA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Straßennetz. Math. Ann. 84, 149-160.       26         217       SPITZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.       28         29       WITTMER, J., JUBNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263–268.       29         30       30       31       31         31       32       32       33         32       33       33       33         34       34       34       34         35       36       36       36         36       37       37       38   | 18 | <ul><li>[12] LINDEBERG, W. (1922). Eine neue Herleitung des Exponentialgesetzes in der Wahrscheinlichkeit-</li></ul>   | 18 |
| [15] MARINEZ, S. and PERGIN, D. (1990). Internotynamic of a Fromman bridge polymer indder in a random environment. J. Phys. A, 29, 1267-1279.       20         [14] KESTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.       21         [25] ORUKIOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.       23         [15] ORUKIOV, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.       23         [16] PÓIXA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Straßennetz. Math. Ann., 84, 149-160.       26         [17] SPITZER, F. (1976). Principles of Randorn Walks, second edition. Springer-Verlag, New York-Heidelberg.       26         [18] WITTMER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.       29         20       30       30         31       31       31         32       32       32         33       34       34         34       34       34         35       36       36         36       36       36         37       38       38         38       39       39         39       39       40         41       41         42       41   | 19 | srechtnung. Math. Z., <b>15</b> , 211–225.   | 19 |
| 14       KENTEX, H. and SUTTZER, F. (1978). A limit theorem related to a new class of self-similar processes.       21         21       J. Wahrsch. Verw. Gebiete, 50.1, 5–25.       23         23       J. Sof-Sa64.       23         24       III.       Póirx, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Straßennetz. Math. Ann. 84, 149–160.       25         26       J. Sof-Sa64.       26         27       III.       Stritzer, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.       26         28       J. J. JUHER, J., JUHER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263–268.       28         29       30       30       31         31       31       31       31         32       J. JUHER, J., JUHER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263–268.       28         33       33       33       33         34       34       34       34         35       36       36       36         36       37       37       38         36       36       36       36         37       38       38       38         38       39       41   | 20 | random environment. J. Phys. A, 29, 1267–1279.   | 20 |
| 22 <i>I.S.</i> Dauxkiov, S. P. (1986). Configurational statistics of a disordered polymer chain. J. Phys. A, 19, 3655-3664.       23         23       16] PótxA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Straßennetz. Math. Ann., 84, 149-160.       24         26       17] SPITZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.       26         27       18] WITTMER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.       29         28       29       30         31       31       31         32       33       30         33       34       34         34       34       34         35       36       36         36       37       37         38       39       39         44       44       44         45       45       45         46       44       44         47       44       44         48       44       44         49       44       44         41       44       44         42       44       44         43       44       44         44   | 21 | [14] KESTEN, H. and SPITZER, F. (1978). A limit theorem related to a new class of self-similar processes.  | 21 |
| [16]       Göshör, S. F. (1980). Comgutational seasates of a usordered polymer tham 5. Fuge. A, 13, 23         24       [16]       Pótxa, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in Straßennetz. Math. Ann., 84, 149–160.       25         27       [17]       Sprizer, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.       26         28       (18)       Wirtriker, J. JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263–268.       29         29       30       30         31       31       31         32       33       33         34       44       34         35       36       36         36       36       36         37       38       38         38       39       39         39       39       39         44       44       44         45       43         46       44         47       44       44         48       44       44         49       49       49         41       41       41         42       42       43         43       43       43<  | 22 | Z. Wahrsch. Verw. Gebiete, <b>50.1</b> , 5–25.   | 22 |
| 24       [16] PÓIVA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in       24         25       [17] SPITZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.       25         27       [18] WITTMER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263-268.       29         30       30         31       31         32       33         33       34         34       34         35       36         36       36         37       38         38       34         39       31         34       34         35       36         36       36         37       38         38       38         39       39         40       41         41       42         42       42         43       43         44       44         45       45         46       47         47       48         48       49         49       49         41       49   | 23 | 3655–3664.   | 23 |
| 25       Straßennetz. Math. Ann., 84, 149–160.       25         [17] SPITZER, F. (1976). Principles of Random Walks, second edition. Springer-Verlag, New York-Heidelberg.       26         27       [18] WITTMER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263–268.       29         29       29         30       30         31       31         32       33         33       31         34       34         35       36         36       36         37       37         38       36         39       39         40       41         41       41         42       42         43       44         44       44         45       45         46       46         47       48         48       49         49       49         50       50         51       50  | 24 | [16] PÓLYA, G. (1921). Über eine Aufgabe der Wahrscheinlichkeitsrechnung betreffend die Irrfahrt in  | 24 |
| 26       III Original in (10) Friendball of Flamboli Frank, etcolin Carlon, opinger Fried, Refer Flamboli, 100, 100, 100, 100, 100, 100, 100, 10   | 25 | Straßennetz. Math. Ann., 84, 149–160.<br>[17] SPITZER F (1976) Principles of Random Walks second edition Springer-Verlag New York-   | 25 |
| 27       [18] WITTMER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Europhys. Lett., 24, 263–268.       28         29       29         30       30         31       31         32       32         33       31         34       34         35       35         36       36         37       38         38       38         39       39         40       40         41       41         42       42         43       44         44       45         45       46         46       47         47       48         48       49         49       49         50       50  | 26 | Heidelberg.  | 26 |
| 28     rophys. Lett., 24, 203-208.     28       29     30     30       31     31     31       32     33     33       33     33     33       34     35     36       35     36     36       36     37     37       38     39     39       40     40       41     41       42     42       43     44       44     45       45     46       46     47       47     48       48     48       49     40       50     50  | 27 | [18] WITTMER, J., JOHNER, A. and JOANNY, J. F. (1993). Random and alternating polyampholytes. Eu-  | 27 |
| 29     29       30     30       31     31       32     32       33     33       34     33       35     36       36     37       37     38       39     38       40     40       41     40       42     41       43     41       44     45       45     41       46     41       47     48       48     48       49     41       41     41       42     41       43     41       44     45       45     41       46     41       47     41       48     41       49     41       40     41       41     41       42     41       43     41       44     45       45     41       46     41       47     41       48     41       49     41       41     41       42     41       43     41       44     41       45   | 28 | rophys. Lett., 24, 263–268.  | 28 |
| 30     30       31     31       32     32       33     33       34     33       35     36       36     36       37     38       39     39       40     40       41     41       42     42       43     44       44     45       45     46       46     47       47     48       49     49       50     50       51     50  | 29 |  | 29 |
| 31     31       32     32       33     33       34     34       35     36       37     37       38     39       40     40       41     41       42     43       43     44       44     45       45     46       46     47       47     48       48     49       50     50       51     51  | 30 |  | 30 |
| 32       32         33       33         34       34         35       35         36       36         37       38         39       39         40       40         41       40         42       41         43       41         44       42         45       41         46       41         47       41         48       41         49       41         49       41         41       41         42       41         43       41         44       45         45       46         46       47         47       48         48       49         50       51  | 31 |  | 31 |
| 33     33       34     34       35     35       36     36       37     37       38     39       40     40       41     41       42     42       43     43       44     45       45     46       47     48       48     49       49     49       41     49       42     41       43     45       44     45       45     46       46     47       47     48       48     49       50     50  | 32 |  | 32 |
| 34     34       35     35       36     36       37     37       38     39       40     40       41     41       42     42       43     43       44     43       45     45       46     47       47     48       49     49       50     50  | 33 |  | 33 |
| 35     36       36     36       37     37       38     38       39     39       40     40       41     41       42     42       43     43       44     45       45     46       46     47       47     48       48     49       50     51  | 34 |  | 34 |
| 36     36       37     37       38     39       40     40       41     41       42     43       43     43       44     45       45     46       47     47       48     49       50     50  | 35 |  | 35 |
| 37     37       38     38       39     39       40     40       41     41       42     42       43     43       44     43       45     46       46     47       47     48       49     49       50     51  | 36 |  | 36 |
| 38       38       38         39       39         40       40         41       41         42       42         43       43         44       45         45       46         46       47         47       48         48       49         50       50   | 37 |  | 37 |
| 39       39         40       40         41       41         42       42         43       43         44       43         45       46         46       46         47       47         48       49         50       50  | 38 |  | 38 |
| 40       40         41       41         42       42         43       43         44       43         45       44         46       45         47       47         48       49         50       50  | 39 |  | 39 |
| 41     41       42     42       43     43       44     43       45     45       46     46       47     47       48     49       50     50  | 40 |  | 40 |
| 42     42       43     43       44     45       45     45       46     46       47     47       48     49       50     50  | 41 |  | 41 |
| 43     44       44     44       45     45       46     46       47     47       48     49       50     50  | 42 |  | 42 |
| 11     11       45     45       46     46       47     47       48     49       50     50  | 40 |  | 40 |
| 1-     10       46     46       47     47       48     48       49     49       50     50  | 45 |  | 45 |
| 1-1     10       47     47       48     48       49     49       50     50   | 46 |  | 46 |
| 48 49 49 50 50 50  | 47 |  | 47 |
| 49 49 50 50 50 50 51 51 51 51 51 51 51 51 51 51 51 51 51   | 48 |  | 48 |
| 50 50 50 T   | 49 |  | 49 |
|  | 50 |  | 50 |
| 51 51  | 51 |  | 51 |

IMS Collections Vol. 0 (2009) 252-265 © Institute of Mathematical Statistics, 2009 arXiv: math.PR/0000011

1

|                  | Robert M. Mnatsakanov <sup>1,*</sup> and Artak S. Hakobyan <sup>2,*</sup>   |
|------------------|---|
|                  | West Virginia University  |
|                  | <b>Abstract:</b> The problem of recovering a cumulative distribution function (cdf) and corresponding density function from its moments is studied. This problem is a special case of the classical moment problem. The results obtained within the moment problem can be applied in many indirect models, e.g., those based on convolutions, mixtures, multiplicative censoring, and right-censoring, where the moments of unobserved distribution of actual interest can be easily estimated from the transformed moments of the observed distributions. Nonparametric estimation of a quantile function via moments of a target distribution represents another very interesting area where the moment problem arises. In all such models one can apply the present results to recover a function via its moments. In this article some properties of the approximation of cdf, its density function, quantile and quantile density function are obtained as well. |
| Co               | ntents  |
|                  |   |
| 1<br>2<br>3<br>4 | Introduction       251         Notation and Assumptions       253         Asymptotic Properties of $F_{\alpha,\nu}$ and $f_{\alpha,\nu}$ 254         Some Applications and Examples       257   |
| Acl              | nowledgments  |
| Ref              | erences   |
|                  |   |
| _                | <b>T</b> . <b>1</b> . <b>1</b>  |
| 1.               | Introduction  |
|                  |   |
| The              | e probabilistic Stielties moment problem can be described as follows: let a se-   |
| que              | nce $\nu = \{\mu_j, j = 0, 1,\}$ of real numbers be given. Find a probability distribu-   |
| U10I             | 1 on the non-negative real line $\mathbb{R}_+ = [0, \infty)$ , such that $\mu_j = \int t^j dF(t)$ for $j \in [0, 1, \infty)$ .  |
| C+:-             | $-10, 1, \dots$ The classical sciently moment problem was introduced first by sltigg [21]. When the support of the distribution $E$ is correct set  |
| 5016             | $F_{1,2} = [0, T]$ when the support of the distribution $F$ is compact, say, $m(F) = [0, T]$ with $T < \infty$ then the corresponding problem is known as a   |
| əuμ<br>Hər       | $p_{1}r_{1} - [0, r_{1}]$ with $r < \infty$ , then the corresponding problem is known as a usdorff moment problem   |
| 114)<br>1        | Consider two important questions related to the Stielties (or Hausdorff) moment   |
| nro              | hlem.   |
| 019<br>/         | (i) If the distribution F exists is it uniquely determined by the moments $\{u, v\}$  |
| (                | (i) How is this uniquely defined distribution F reconstructed?  |
| (                |   |
| 1                | Department of Statistics, West Virginia University, Morgantown, WV 26506, email: rmnatsak@  |
| sta              | t.wvu.edu   |
| 212-             | work supported by the National Institute for Occupational Safety and Health, contract no. 2005-M-12857.   |
| - 1 2-2          | Popartment of Industrial and Management Systems Engineering, West Virginia University.  |
| Mor              | gantown, WV 26506, email: artakhak@yahoo.com  |
|                  | AMS 2000 subject classifications: Primary 62G05: secondary 62G20  |
|                  | Kannada and anagan Drababilistin many problem. Mataminata distribution Manual   |

If there is a positive answer to question (i) we say that a distribution F is moment-determinate (*M*-determinate), otherwise it is *M*-indeterminate.

In this paper we mainly address the question of recovering the M-determinate distribution (density and quantile functions) via its moments in the Hausdorff moment problem, i.e., we study question (*ii*). Another question we focus on here is the estimation of an unknown distribution and its quantile function, given the estimated moments of the target distribution.

It is known from the probabilistic moment problem that under suitable condi-tions an *M*-determinate distribution is uniquely defined by its moments. There are many articles that investigated the conditions (for example, the Carlemann's and the Krein's conditions), under which the distributions are either *M*-determinate or M-indeterminate. See, e.g., Akhiezer [2], Feller [6], Lin [10-11], and Stoyanov [22-24] among others. However, there are very few works dealing with the reconstruction of distributions via their moments. Several inversion formulas were obtained by invert-ing the moment generating function and Laplace transform (Shohat and Tamarkin [20], Widder [27], Feller [6], Chauveau et al. [4], and Tagliani and Velasquez [25]). These methods are too restrictive, since there are many distributions for which the moment generating function does not exist even though all the moments are finite.

The reconstruction of an M-determinate cdf by means of mixtures having the same assigned moments as the target distribution have been proposed in Lindsay et al. [12]. Note that this procedure requires calculations of high-order Hankel de-terminants, and due to ill-conditioning of the Hankel matrices this method is not useful when the number of assigned moments is large. The reconstruction of an unknown density function using the Maximum Entropy principle with the specified ordinary and fractional moments has been studied in Kevasan and Kapur [9] and Novi Inverardi et al. [18], among others. 

In Mnatsakanov and Ruymgaart [17] the constructions (2.2) and (3.13) (see Sections 2 and 3 below) have been introduced, and only their convergence has been established.

Different types of convergence of maximum entropy approximation have been studied by Borwein and Lewis [3], Frontini and Tagliani [7], and Novi Inverardi et al. [18], but the rates of approximations have not been established yet. Our construction enables us to derive the uniform rate of convergence for moment-recovered cdfs  $F_{\alpha,\nu}$ , corresponding quantile function  $Q_{\alpha}$ , and the uniform convergence of the moment-recovered density approximation  $f_{\alpha,\nu}$ , as the parameter  $\alpha \to \infty$ . Other constructions of moment-recovered cdfs and pdfs (see, (3.13) and (3.14) in Remark (3.2) were proposed in Mnatsakanov [13-14], where the uniform and  $L_1$ -rates of the approximations were established. 

The paper is organized as follows: in Section 2 we introduce the notation and assumptions, while in Section 3 we study the properties of  $F_{\alpha,\nu}$  and  $f_{\alpha,\nu}$ . Note that our construction also gives a possibility to recover different distributions through the simple transformations of moment sequences of given distributions (see Theorem 3.1 in Section 3 and similar properties derived in Mnatsakanov [13]: Theorem 1 and Corollary 1). In Theorem 3.2 we state the uniform rate of convergence for moment-recovered cdfs. In Theorem 3.3 as well as in Corollaries 3.1 and 3.2 we apply the constructions (2.2) and (3.11) to recover the pdf f, the quantile function Q, and the corresponding quantile density function q of F given the moments of F. In Section 4 some other applications of the constructions (2.2) and (3.11) are discussed: the uniform convergence of the empirical counterpart of (2.2), the rate of approximation of moment-recovered quantile function (see (4.4) in Section 4) along with the demixing and deconvolution problems in several particular models. 

Note that our approach is particularly applicable in situations where other estimators cannot be used, e.g., in situations where only moments (empirical) are available. The results obtained in this paper will not be compared with similar results derived by other methods. We only carry out the calculations of moment-recovered cdfs, pdfs, and quantile functions, and compare them with the target distributions via graphs in several simple examples. We also compare the performances of  $F_{\alpha,\nu}$ and  $f_{\alpha,\nu}$  with the similar constructions studied in Mnatsakanov [13-14] (see, Figures 1 (b) and 3 (b)). The moment-estimated quantile function  $\hat{Q}_{\alpha}$  and well known Harrell-Davis quantile function estimator  $\hat{Q}_{HD}$  (Sheather and Marron [19]) defined in (4.6) and (4.7), respectively, are compared as well (see, Figure 2 (b)).

## 2. Notation and Assumptions

Suppose that the *M*-determinate cdf *F* is absolute continuous with respect to the Lebesgue measure and has support [0, T],  $T < \infty$ . Denote the corresponding density function by *f*. Our method of recovering the cdf F(x),  $0 \le x \le T$ , is based on an inverse transformation that yields a solution of the Hausdorff moment problem.

Let us denote the moments of F by

(2.1) 
$$\mu_{j,F} = \int t^j dF(t) = (\mathcal{K}F)(j), j \in \mathbb{N},$$

and assume that the moment sequence  $\nu = (\mu_{0,F}, \mu_{1,F}, \dots)$  determines F uniquely. An approximate inverse of the operator  $\mathcal{K}$  from (2.1) constructed according to

(2.2) 
$$\left(\mathcal{K}_{\alpha}^{-1}\nu\right)(x) = \sum_{k=0}^{\left[\alpha x\right]} \sum_{j=k}^{\infty} \frac{(-\alpha)^{j-k}}{(j-k)!} \frac{\alpha^k}{k!} \mu_{j,F}, \ 0 \leqslant x \leqslant T, \ \alpha \in \mathbb{R}_+,$$

is such that  $\mathcal{K}_{\alpha}^{-1}\mathcal{K}F \to_w F$ , as  $\alpha \to \infty$  (see, Mnatsakanov and Ruymgaart [17]). Here  $\to_w$  denotes the weak convergence of cdfs, i.e. convergence at each continuity point of the limiting cdf. The success of the inversion formula (2.2) hinges on the convergence

$$P_{\alpha}(t,x) = \sum_{k=0}^{[\alpha x]} \frac{(\alpha t)^k}{k!} e^{-\alpha t} \to \begin{cases} 1, & t < x \\ 0, & t > x \end{cases},$$

(2.3)

as  $\alpha \to \infty$ . This result is immediate from a suitable interpretation of the left hand side as a sum of Poisson probabilities.

For any moment sequence  $\nu = \{\nu_j, j \in \mathbb{N}\}$ , let us denote by  $F_{\nu}$  the cdf recovered via  $F_{\alpha,\nu} = \mathcal{K}_{\alpha}^{-1}\nu$  according to (2.2), when  $\alpha \to \infty$ , i.e.

(2.4) 
$$F_{\alpha,\nu} \to_w F_{\nu}, \text{ as } \alpha \to \infty.$$

Note that if  $\nu = {\mu_{j,F}, j \in \mathbb{N}}$  is the moment sequence of F, the statement (2.4) with  $F_{\nu} = F$  is proved in Mnatsakanov and Ruymgaart [17].

To recover a pdf f via its moment sequence  $\{\mu_{j,F}, j \in \mathbb{N}\}$ , consider the ratio:

where  $\Delta F_{\alpha,\nu}(x) = F_{\alpha,\nu}(x+\Delta) - F_{\alpha,\nu}(x)$  and  $\alpha \to \infty$ .

In the sequel the uniform convergence on any bounded interval in  $\mathbb{R}_+$  will be denoted by  $\longrightarrow_u$ , while the sup-norm between two functions  $f_1$  and  $f_2$  by  $|| f_1 - f_2 ||$ . Note also that the statements from Sections 3 and 4 are valid for distributions defined on any compact  $[0,T], T < \infty$ . Without loss of generality we assume that F has support [0, 1].

## 3. Asymptotic Properties of $F_{\alpha,\nu}$ and $f_{\alpha,\nu}$

In this Section we present asymptotic properties of the moment-recovered cdf  $F_{\alpha,\nu}$ and pdf  $f_{\alpha,\nu}$  functions based on the transformation  $\mathcal{K}_{\alpha}^{-1}\nu$  (2.2). The uniform approximation rate of  $F_{\alpha,\nu}$  and the uniform convergence of  $f_{\alpha,\nu}$  are derived as well.

Denote the family of all cdfs defined on [0,1] by  $\mathbb{F}$ . The construction (2.2) gives us the possibility to recover also two non-linear operators  $\mathcal{A}_k : \mathbb{F} \times \mathbb{F} \to \mathbb{F}, k =$ 1, 2, defined as follows: denote the convolution with respect to the multiplication operation on  $\mathbb{R}_+$  by

(3.1) 
$$F_1 \otimes F_2(x) = \int F_1(x/\tau) \, dF_2(\tau) := \mathcal{A}_1(F_1, F_2)(x), \ 0 \le x \le 1,$$

while the convolution with respect to the addition operation is denoted by

$$F_1 \star F_2(x) = \int F_1(x-\tau) \, dF_2(\tau) := \mathcal{A}_2(F_1, F_2)(x), \ 0 \le x \le 2.$$

For any two moment sequences  $\nu_1 = \{\mu_{j,F_1}, j \in \mathbb{N}\}\$  and  $\nu_2 = \{\mu_{j,F_2}, j \in \mathbb{N}\},\$ define  $\nu_1 \odot \nu_2 = \{\mu_{j,F_1} \times \mu_{j,F_2}, j \in \mathbb{N}\}$  and  $\nu_1 \oplus \nu_2 = \{\bar{\nu}_j, j \in \mathbb{N}\}$ , where

(3.2) 
$$\bar{\nu}_j = \sum_{m=0}^j {j \choose m} \mu_{m,F_1} \times \mu_{j-m,F_2} \,,$$

while  $\mu_F^{\odot k} = \{\mu_{i,F}^k, j \in \mathbb{N}\}$  and  $F^{\otimes k} = F \otimes \cdots \otimes F$  for the corresponding k-fold convolution (cf. (3.1)). Also denote by  $F \circ \phi^{-1}$  the composition  $F(\phi^{-1}(x)), x \in [0, 1]$ , with  $\phi$  - continuous and increasing function  $\phi: [0,1] \to [0,1]$ .

Since cdfs  $\mathcal{A}_1(F_1, F_2) = F_1 \otimes F_2$ ,  $\mathcal{A}_2(F_1, F_2) = F_1 \star F_2$ , and  $F \circ \phi^{-1}$  have compact support, they all are *M*-determinate and have the moment sequences  $\nu_1 \odot \nu_2$ ,  $\nu_1 \oplus \nu_2$ , and  $\nu = \{\bar{\mu}_j, j \in \mathbb{N}\},$  with

(3.3) 
$$\bar{\mu}_j = \int [\phi(t)]^j \, dF(t),$$

respectively. Hence, applying Theorem 3.1 from Mnatsakanov and Ruymgaart [17] a statement similar to the one in Mnatsakanov [13] (see Theorem 1 and Corollary 1, where T = T' = 1) is obtained. Besides, the following statement is true:

THEOREM 3.1. If 
$$\nu = \sum_{k=1}^{m} \beta_k \mu_F^{\odot k}$$
, where  $\sum_{k=1}^{m} \beta_k = 1, \beta_k > 0$ , then (2.4) holds with

(3.4) 
$$F_{\nu} = \sum_{k=1}^{m} \beta_k F^{\otimes k} \,.$$

*Proof.* The equation (3.4) follows from Theorem 1 (i) (Mnatsakanov [13]) and the linearity of  $\mathcal{K}_{\alpha}^{-1}\nu$ . 

The construction (2.2) is also useful when recovering the quantile function Q(t) = $\inf\{x: F(x) \ge t\}$  via moments (see (4.5) in Section 4). Define  $Q_{\alpha} = F_{\alpha,\nu_{Q}}$ , where З  $\nu_Q = \left\{ \int_{-1}^{1} [F(u)]^j \, du, j \in \mathbb{N} \right\}.$ (3.5)The following statement is true: COROLLARY 3.1. If F is continuous, then  $Q_{\alpha} \to_w Q$ , as  $\alpha \to \infty$ . *Proof.* Replacing the functions  $\phi$  and F in (3.3) by F and the uniform cdf on [0,1], respectively, we obtain from Theorem 1 (iv) (Mnatsakanov [13]) that  $Q_{\alpha} =$  $F_{\alpha,\nu_Q} \to_w F_{\nu} = F^{-1}$  as  $\alpha \to \infty$ . Under additional conditions on the smoothness of F one can obtain the uniform rate of convergence in (2.4) and, hence, in Theorem 3.1 too. Consider the following condition F'' = f' is bounded on [0, 1]. (3.6)THEOREM 3.2. If  $\nu = \{\mu_{j,F}, j \in \mathbb{N}\}$ , and (3.6) holds, we have  $\sup_{0 \le x \le 1} \left| F_{\alpha,\nu}(x) - F(x) \right| = O\left(\frac{1}{\alpha}\right), \quad \text{as } \alpha \to \infty.$ (3.7)*Proof.* Let us use the following representation  $P_{\alpha}(t,x) = \mathbf{P}\{N_{\alpha t} \leqslant \alpha x\} = \mathbf{P}\{S_{[\alpha x]} \ge \alpha t\}.$ Here  $\{N_{\alpha t}, t \in [0, 1]\}$  is a Poisson process with intensity  $\alpha t, S_m = \sum_{k=0}^m \xi_k, S_0 = 0$ , with  $\xi_k$  being *iid* Exp(1) random variables. Integration by parts gives (3.8)  $F_{\alpha,\nu}(x) = (\mathcal{K}_{\alpha}^{-1}\nu)(x) = \int_{0}^{1} \sum_{k=0}^{\lfloor \alpha x \rfloor} \frac{(\alpha t)^{k}}{k!} \sum_{k=0}^{\infty} \frac{(-\alpha t)^{j-k}}{(j-k)!} dF(t)$  $= \int_{\alpha}^{1} P_{\alpha}(t, x) dF(t) = \int_{\alpha}^{1} \mathbf{P}\{S_{[\alpha x]} \ge \alpha t\} dF(t)$  $= F(t) \mathbf{P}\{S_{[\alpha x]} \ge \alpha t\}\Big|_{0}^{1} - \int_{0}^{1} F(t) d\mathbf{P}\{S_{[\alpha x]} \ge \alpha t\}\Big|_{0}^{1}$  $= \mathbf{P}\{S_{[\alpha x]} \ge \alpha\} + \int_{\alpha}^{1} F(t) \, d\, \mathbf{P}\{S_{[\alpha x]} \le \alpha t\} = \int_{\alpha}^{\infty} F(t) \, d\, \mathbf{P}\{S_{[\alpha x]} \le \alpha t\}.$ Thus, (3.6) and the argument used in Adell and de la Cal [1] yield (3.7). REMARK 3.1. When  $supp\{F\} = \mathbb{R}_+, F_{\alpha,\nu}(x) = \int_0^\infty P_\alpha(t,x)dF(t)$  (cf. with (3.8)). According to Mnatsakanov and Klaassen [16] (see the proof of Theorem 3.1), one can derive the exact rate of approximation of  $F_{\alpha,\nu}$  in the space  $L_2(\mathbb{R}_+, dF)$ . Namely, if the pdf f is bounded, say by C > 0, then 

 $\int_0^\infty \left(F_{\alpha,\nu}(x) - F(x)\right)^2 dF(x) \le \frac{2C}{\alpha}.$ <sup>50</sup>
<sub>51</sub>

imsart-coll ver. 2008/08/29 file: Mnatsak\_Revised.tex date: April 10, 2009

Now let us consider the moment-recovered density function  $f_{\alpha,\nu}$  defined in (2.5)

and denote by  $\Delta(f, \delta) = \sup_{|t-s| \leq \delta} |f(t) - f(s)|$  the modulus of continuity of f,

where  $0 < \delta < 1$ . THEOREM 3.3. If the pdf f is continuous on [0, 1], then  $f_{\alpha,\nu} \longrightarrow_u f$  and (3.9)  $|| f_{\alpha,\nu} - f || \leq \Delta(f, \delta) + \frac{2 || f ||}{\alpha \delta^2} + o\left(\frac{1}{\alpha}\right), \text{ as } \alpha \to \infty.$ Proof. Since  $[\alpha(x + 1/\alpha)] = [\alpha x] + 1$ , for any  $x \in [0, 1]$ , we have (3.10)  $f_{\alpha,\nu}(x) = \alpha \left[\sum_{k=0}^{[\alpha x]+1} \sum_{j=k}^{\infty} \frac{(-\alpha)^{j-k}}{(j-k)!} \frac{\alpha^k}{k!} \mu_{j,F} - \sum_{k=0}^{[\alpha x]} \sum_{j=k}^{\infty} \frac{(-\alpha)^{j-k}}{(j-k)!} \frac{\alpha^k}{k!} \mu_{j,F} \right],$ and, after some algebra (3.10) yields (3.11)  $f_{\alpha,\nu}(x) = \frac{\alpha^{[\alpha x]+2}}{\Gamma([\alpha x]+2)} \cdot \sum_{m=0}^{\infty} \frac{(-\alpha)^m}{m!} \mu_{m+[\alpha x]+1,F}.$ 

Let g(t, a, b) denote a gamma pdf with shape and scale parameters a and b, respectively. Substitution of (2.1) into the right hand side of (3.11) gives

(3.12) 
$$f_{\alpha,\nu}(x) = \frac{\alpha^{[\alpha x]+2}}{\Gamma([\alpha x]+2)} \int_0^1 \sum_{m=0}^\infty \frac{(-\alpha t)^m}{m!} t^{[\alpha x]+1} dF(t)$$

To show (3.9), note that the pdf g in (3.12) has mean  $([\alpha x] + 2)/\alpha$  and variance  $([\alpha x] + 2)/\alpha^2$ , respectively. The rest of the proof is similar to the lines of Theorem 1 (i) (Mnatsakanov [14]).

REMARK 3.2. In Mnatsakanov [13-14] the uniform and  $L_1$ -rates of momentrecovered approximations of F and f defined by

(3.13) 
$$F_{\alpha,\nu}^{*}(x) = \sum_{k=0}^{[\alpha x]} \sum_{j=k}^{\alpha} {\alpha \choose j} {j \choose k} (-1)^{j-k} \mu_{j,F}$$

and

(3.14) 
$$f^*_{\alpha,\nu}(x) = \frac{\Gamma(\alpha+2)}{\Gamma([\alpha x]+1)} \sum_{m=0}^{\alpha-[\alpha x]} \frac{(-1)^m \mu_{m+[\alpha x],F}}{m! (\alpha-[\alpha x]-m)!}, \ x \in [0,1], \ \alpha \in \mathbb{N},$$

are established. In Section 4, see Example 4.2, the cdf  $F(x) = x^3 - 3x^3 \ln x$  and its density function  $f(x) = -9x^2 \ln x, 0 \leq t \leq 1$ , are recovered using  $F_{\alpha,\nu}$  and  $F^*_{\alpha,\nu}$ , and  $f_{\alpha,\nu}$  and  $f^*_{\alpha,\nu}$  constructions, (see Figures 1 (b) and 3 (b), respectively).

The formulas (3.11) and (3.14) with  $\nu = \nu_Q$  defined according to (3.5) can be used to recover a quantile density function

$$q(x) = Q'(x) = \frac{1}{f(F^{-1}(x))}, \quad x \in [0, 1].$$

imsart-coll ver. 2008/08/29 file: Mnatsak\_Revised.tex date: April 10, 2009

З

 $F^{-1}$  instead of F yields  $q_{\alpha}(x) = \int_{0}^{1} g(F(u), [\alpha x] + 2, \alpha) \, du$ 

and corresponding moment-recovered quantile density function

$$q_{\alpha,\beta}(x) = \int_0^1 g(F_{\beta,\nu}(u), [\alpha x] + 2, \alpha) \, du \,, \ \alpha, \beta \in \mathbb{N}.$$

Here  $F_{\beta,\nu}$  is a moment-recovered cdf of F. As a consequence of Theorem 3.3 we have the following

COROLLARY 3.2. If  $q_{\alpha} = f_{\alpha,\nu_Q}$ , with  $\nu_Q$  defined in (3.5), and f is continuous on [0, 1] with  $\inf_{0 \leq x \leq 1} f(x) > \gamma > 0$ , then  $q_{\alpha} \longrightarrow_{u} q$  and

$$|| q_{\alpha} - q || \leq \frac{\Delta(f, \delta)}{\gamma^2} + \frac{2 || f ||}{\alpha \, \delta^2 \, \gamma^2} + o\left(\frac{1}{\alpha}\right), \quad as \quad \alpha \to \infty.$$

Finally, note that taking  $\nu = \nu_Q$  in (3.14), we derive another approximation  $q_{\alpha}^* =$  $f^*_{\alpha,\nu_Q}$  of q based on Beta densities  $\beta(\cdot, a, b)$  with the shape parameters  $a = [\alpha x] + 1$ and  $b = \alpha - [\alpha x] + 1$ :

(3.16) 
$$q_{\alpha}^{*}(x) = \int_{0}^{1} \beta(F(u), [\alpha x] + 1, \alpha - [\alpha x] + 1) \, du \, .$$

# 4. Some Applications and Examples

In this Section the construction of the moment-recovered cdf  $F_{\alpha,\nu}$  is applied to the problem of nonparametric estimation of a cdf, its density and a quantile functions as well as to the problem of demixing in exponential, binomial and negative binomial mixtures, and deconvolution in error-in-variable model. In Theorems 4.1 we derive the uniform rate of convergence for the empirical counterpart of  $F_{\alpha,\nu}$  denoted by  $F_{\alpha}$ , i.e. for  $F_{\alpha} = F_{\alpha,\hat{\nu}}$ , where  $\hat{\nu}$  is the sequence of all empirical moments of the sample from F. In Theorem 4.2 the uniform rate of approximation for moment-recovered quantile function of F is obtained. Finally, the graphs of moment-recovered cdfs, pdfs, and quantile functions are presented in Figures 1-3.

Direct model. Let  $X_1, \ldots, X_n$  be a random sample from F defined on [0, 1]. Denote by  $F_n$  the empirical cdf (ecdf) of the sample  $X_1, \ldots, X_n$ :

$$\hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n I_{[0,t]}(X_i), \quad 0 \le t \le 1.$$

Substitution of the empirical moments

$$\hat{\nu}_j = \frac{1}{n} \sum_{i=1}^n X_i^j, \, j \in \mathbb{N} \,,$$

instead of  $\mu_{j,F}$  into (2.2) yields

$$\widetilde{F}_{\alpha}(x) = F_{\alpha,\hat{\nu}}(x) = \int_0^1 P_{\alpha}(t,x) d\, \hat{F}_n(t) = \int_0^1 \mathbf{P}\{S_{[\alpha x]} \ge \alpha t\} d\, \hat{F}_n(t) \,. \tag{50}$$

imsart-coll ver. 2008/08/29 file: Mnatsak\_Revised.tex date: April 10, 2009

(3.15)

З

Furthermore, the empirical analogue of (3.8) admits a similar representation

$$\widetilde{F}_{\alpha}(x) = \int_{0}^{\infty} \widehat{F}_{n}(t) \, d \, \mathbf{P}\{S_{[\alpha x]} \leqslant \alpha t\}.$$

The application of the Theorem 3.2 and the asymptotic properties of  $\hat{F}_n$  yield

THEOREM 4.1. If  $\nu = {\mu_{j,F}, j \in \mathbb{N}}$ , then under condition (3.6) we have

(4.1) 
$$\sup_{0 \leqslant x \leqslant 1} \left| \widetilde{F}_{\alpha}(x) - F(x) \right| = O\left(\frac{1}{\sqrt{n}}\right) + O\left(\frac{1}{\alpha}\right) \quad \text{a.s., as} \quad \alpha, n \to \infty.$$

REMARK 4.1. In Mnatsakanov and Ruymgaart [17] the weak convergence of the moment-empirical processes  $\{\sqrt{n}\{\tilde{F}_n(t) - F(t)\}, t \in [0, 1]\}$  to the Brownian bridge is obtained.

Of course, when the sample is directly drawn from the cdf F of actual interest, one might use the ecdf  $\hat{F}_n$  and empirical process  $U_n = \sqrt{n} (\hat{F}_n - F)$ . The result mentioned in Remark 4.1 yields, that even if the only information available is the empirical moments, we still can construct different test statistics based on the moment-empirical processes  $\tilde{U}_n = \sqrt{n} (\tilde{F}_n - F)$ .

On the other hand, using the construction (3.11), one can estimate the density function f given only the estimated or empirical moments in:

 $\infty$ 

(4.2) 
$$f_{\alpha,\hat{\nu}}(x) = \frac{\alpha^{\lfloor \alpha x \rfloor + 2}}{\Gamma(\lfloor \alpha x \rfloor + 2)} \sum_{m=0} \frac{(-\alpha)^m}{m!} \hat{\nu}_{m+\lfloor \alpha x \rfloor + 1}, x \in [0,1].$$

 $[\alpha r] \perp 2$ 

REMARK 4.2. In practice, the parameter  $\alpha$  as well as the number of summands in (4.2) (and the number of summands in the inner summation of  $F_{\alpha,\hat{\nu}}$ ) can be chosen as the functions of  $n: \alpha = \alpha(n) \to \infty$  and  $M = M(n) \to \infty$  as  $n \to \infty$ , that optimize the accuracy of corresponding estimates. Further analysis is required to derive the asymptotic forms of  $\alpha(n)$  and M(n) as  $n \to \infty$ . This question is currently under investigation and is beyond the scope of the present article.

Note that the construction (4.2) yields the estimator  $\hat{f}_{\alpha}(x) = f_{\alpha,\hat{\nu}}$  with  $\hat{\nu} = \{\hat{\nu}_j, j \in \mathbb{N}\}$ :

$$\hat{f}_{\alpha}(x) = \frac{\alpha}{n} \sum_{i=1}^{n} \frac{(\alpha X_i)^{[\alpha x]+1}}{([\alpha x]+1)!} e^{-\alpha X_i} = \frac{1}{n} \sum_{i=1}^{n} g(X_i, [\alpha x]+2, \alpha), \ x \in [0, 1].$$

Here  $g(\cdot, [\alpha x] + 2, \alpha)$  is defined in (3.12). The estimator  $\hat{f}_{\alpha}$  does not represent a traditional kernel density estimator of f. It is defined by a  $\delta$ -sequence, which consists of the gamma density functions of varying shapes (the shape and the rate parameters are equal to  $[\alpha x] + 2$  and  $\alpha$ , respectively). It is natural to use this estimator when  $supp\{F\} = [0, \infty)$ , since, in this case, the supports of f and gamma kernel densities coincide and one avoids the boundary effect of  $\hat{f}_{\alpha}$  (cf. Chen [5]).

Some asymptotic properties such as the convergence in probability of  $\hat{f}_{\alpha}$  uniformly on any bounded interval and the Integrated Mean Squared Error (*IMSE*) of  $\hat{f}_{\alpha}$  have been studied in Mnatsakanov and Ruymgaart [17] and Chen [5], respectively.

<sup>50</sup> Applying the results from Mnatsakanov and Khmaladze [15], where the neces-<sup>51</sup> sary and sufficient conditions for  $L_1$ -consistency of general kernel density estimates З

are established, one can prove in a similar way (cf. Mnatsakanov [14], Theorem 3) that if f is continuous on [0,1], then  $E \mid |\hat{f}_{\alpha} - f \mid |_{L_1} \to 0$ , as  $\sqrt{\alpha}/n \to 0$  and  $\alpha, n \to \infty$ .

З

Exponential mixture model. Let  $Y_1, \ldots, Y_n$  be a random sample from the mixture of exponentials

$$G(x) = \int_0^T (1 - e^{-x/\tau}) \, dF(\tau) \,, \ x \ge 0 \,.$$

The unknown cdf F can be recovered according to the construction  $F_{\alpha,\nu} = \mathcal{K}_{\alpha}^{-1} \nu$ with  $\nu = \{\mu_{j,G}/j!, j \in \mathbb{N}\}$ . Similarly, given the sample  $Y_1, \ldots, Y_n$  from G and taking  $F_{\alpha,\hat{\nu}} = \mathcal{K}_{\alpha}^{-1}\hat{\nu}$ , where  $\hat{\nu} = \{\hat{\mu}_{j,G}/j!, j \in \mathbb{N}\}$ , we obtain the estimate of F. Here  $\{\hat{\mu}_{j,G}, j \in \mathbb{N}\}$  are the empirical moments of the sample  $Y_1, \ldots, Y_n$ . The regularized inversion of the noisy Laplace transform and the  $L_2$ -rate of convergence were obtained in Chauveau *et al.* [4].

Binomial and negative binomial mixture models. When  $Y_1, \ldots, Y_n$  is a random sample from the binomial or negative binomial mixture distributions, respectively:

$$p(x) := P(Y = x) = \int_0^1 {\binom{m}{x}} \tau^x (1 - \tau)^{m-x} dF(\tau), \ x = 0, \dots m,$$

$$p(x) := P(Y = x) = \int_0^1 \frac{\Gamma(r + x)}{\Gamma(r) \, x!} \left(\frac{1}{1 + \tau}\right)^r \left(\frac{\tau}{1 + \tau}\right)^x dG(\tau) \,, \ x = 0, 1, \dots \,,$$

where m and r are given positive integers. Assume that the unknown mixing cdfs Fand G are such that F has at most  $\frac{m+1}{2}$  support points in (0, 1), while G is a right continuous cdf on (0, 1). In both models the mixing distributions are identifiable (see, for example, Teicher [26] for binomial mixture model). Note also that the *j*th moments of F and G are related to the *j*th factorial moments of corresponding  $Y_i$ 's in the following ways:

$$\mu_{j,F} = \frac{1}{m^{[j]}} E(Y_1^{[j]}) \text{ and } \mu_{j,G} = \frac{1}{r_{(j)}} E(Y_1^{[j]}).$$

Here  $y^{[j]} = y(y-1)\cdots(y-j+1)$  and  $r_{(j)} = r(r+1)\cdots(r+j-1)$ . To estimate F and G one can use the moment-recovered formulas (2.2) or (3.13) with  $\mu_{j,F}$  and  $\mu_{j,G}$  defined in previous two equations where the theoretical factorial moments are replaced by corresponding empirical counterparts. The asymptotic properties of the derived estimators of F and G will be studied in a separate work.

Deconvolution problem: error-in-variable model. Consider the random variable Y = X + U, with cdf G, where U (the error) has some known symmetric distribution  $F_2$ , X has cdf  $F_1$  with a support [0, T], and U and X are independent. This model, known as an error-in-variable model, corresponds to the convolution  $G = F_1 \star F_2$ . Assuming that all moments of X and U exist, the moments  $\{\bar{\nu}_j, j \in \mathbb{N}\}$  of Y are described by (3.2). Hence, given the moments of U (with E(U) = 0), we can recalculate the moments of  $F_1$  as follows:  $\mu_{1,F_1} = \bar{\nu}_1$ ,  $\mu_{2,F_1} = \bar{\nu}_2 - \mu_{2,F_2}$ , and so on. So that, assuming that we already calculated  $\mu_{k,F_1}$ , or estimated them by  $\mu_{k,F_1}^*$  for  $1 \leq k \leq j-2$ , we will have, for any  $j \geq 1$ :

$$\mu_{j,F_1} = \bar{\nu}_j - \sum_{m=2}^{j} \binom{j}{m} \mu_{m,F_2} \times \mu_{j-m,F_1}$$
<sup>50</sup>
<sub>51</sub>

imsart-coll ver. 2008/08/29 file: Mnatsak\_Revised.tex date: April 10, 2009



FIG 1. (a) Approximation of  $G(x) = x - x \ln x$  by  $F_{\alpha,\nu}$  and (b) Approximation of  $G(x^3)$  by  $F_{\alpha,\nu}$ and by  $F^*_{\alpha,\nu}$ 

or, respectively,

$$\mu_{j,F_1}^* = \hat{\mu}_{j,G} - \sum_{m=2}^{j} {j \choose m} \mu_{m,F_2} \times \mu_{j-m,F_1}^*$$

given the sample  $Y_1, \ldots, Y_n$  from cdf G. Now the moment-recovered estimate of  $F_1$  will have the form  $F_{\alpha,\hat{\nu}} = \mathcal{K}_{\alpha}^{-1}\hat{\nu}$ , where  $\hat{\nu} = \{\mu_{j,F_1}^*, j \in \mathbb{N}\}$ . The alternative construction of the kernel type estimate of  $F_1$  based on the Fourier transforms is studied in Hall and Lahiri [8], where the  $\sqrt{n}$ -consistency and other properties of the estimated moments  $\mu_{j,F_1}^*, j \in \mathbb{N}$ , are derived as well.

Example 4.1. Consider the moment sequence  $\mu = \{1/(j+1), j \in \mathbb{N}\}$ . The corresponding moment-recovered distribution  $F_{\alpha,\mu} = \mathcal{K}_{\alpha}^{-1} \mu$  is a good approximation of F(x) = x already with  $\alpha = 50$  and M = 100.

Assume now that we want to recover the distribution G with corresponding moments  $\nu_{j,G} = 1/(j+1)^2, j \in \mathbb{N}$ . Since we can represent  $\nu_G = \mu \odot \mu$ , we conclude from Theorem 1 (i) in Mnatsakanov [13], that  $G = F \otimes F$ , with F(x) = x, and hence  $G(x) = x - x \ln x, 0 \leq x \leq 1$ . We plotted the curves of  $F_{\alpha,\nu_G}$  (the solid line) and G (the dashed line) on Figure 1 (a). We took  $\alpha = 50$  and M = 200, the number of terms in the inner summation of the formula (2.2). From Figure 1 (a) we can see that the approximation of G by  $F_{\alpha,\nu_G}$  at x = 0 is not as good as inside of the interval [0, 1]. This happened because the condition (3.6) from Theorem 3.2 is not valid for g'(x) = G''(x) = -1/x.

Example 4.2. To recover the distribution F via moments  $\nu_j = 9/(j+3)^2, j \in \mathbb{N}$ , note that  $\nu_j = \nu_{aj,G}$ , with a = 1/3. Hence,  $F(x) = G(x^3) = x^3 - x^3 \ln(x^3), 0 \le x \le 1$ (Theorem 1 (iii), Mnatsakanov [13]). We conducted computations of moment-recovered cdf  $F_{\alpha,\nu}$  when  $\alpha = 50$  and the number of terms in the inner summation of the formula (2.2) is equal to 200. Also, we calculated  $F^*_{\alpha,\nu}$  defined in (3.13) with  $\alpha = 32$ . See Figure 1 (b), where we plotted  $F_{\alpha,\nu}$  (the solid blue line),  $F^*_{\alpha,\nu}$  (the 

Mnatsakanov and Hakobyan

solid red line), and F (the dashed line), respectively. These two approximations of cdf F justify a good fit already with  $\alpha = 50$  and M = 200 for the first one and with  $\alpha = 32$  for the second one. From Figure 1 (b) we can see that the performance of  $F_{\alpha,\nu}^*$  is slightly better compared to  $F_{\alpha,\nu}$ :  $F_{\alpha,\nu}^*$  does not have the "boundary" effect around x = 1.

Estimation of a quantile function Q and quantile density function q. Assume that a random variable X has a continuous cdf F defined on [0, 1]. To approximate (estimate) the quantile function Q given only the moments (estimated moments) of F, one can use Corollary 3.1. Indeed, after some algebra, we have

(4.3) 
$$Q_{\alpha}(x) = F_{\alpha,\nu_Q}(x) = \int_0^1 P_{\alpha}(F(u), x) \, du \,, \ 0 \le x \le 1 \,,$$

where  $\nu_Q$  and  $P_{\alpha}(\cdot, \cdot)$  are defined in (3.5) and in (2.3), respectively. Comparing (4.3) and (3.8) we can prove in a similar way (see, the proof of Theorem 3.2) the following

THEOREM 4.2. If f' is bounded and  $\inf_{0 \le x \le 1} f(x) > \gamma > 0$ , then

(4.4) 
$$\sup_{0 \le x \le 1} \left| Q_{\alpha}(x) - Q(x) \right| = O\left(\frac{1}{\alpha}\right), \quad \text{as } \alpha \to \infty.$$

Now, given only the moment sequence  $\nu$  of F, one can construct the approximation  $Q_{\alpha,\beta}$  of Q by substituting the moment-recovered cdf  $F_{\beta,\nu}$  (instead of F) in the right hand side of (4.3). Let us denote the corresponding approximation of Q by

(4.5) 
$$Q_{\alpha,\beta}(x) = \int_0^1 P_\alpha(F_{\beta,\nu}(u), x) \, du \,, \ \alpha \,, \beta \in \mathbb{N}.$$

Figure 2 (a) shows the cdf  $F(x) = x^3 - x^3 \ln(x^3)$  (the dashed line), introduced in Example 4.2, and its quantile approximation  $Q_{\alpha,\beta}$  (the solid line), when  $\nu = \{9/(j+3)^2, j \in \mathbb{N}\}, \alpha = \beta = 100$ , and M = 200.

Replacing F by the empirical  $\hat{F}_n$  in (4.3), (3.15), and in (3.16) yields the following estimators, respectively, based on the spacings  $\Delta X_{(i)} = X_{(i)} - X_{(i-1)}, i = 1, \ldots, n+1$ :

(4.6) 
$$\hat{Q}_{\alpha}(x) = F_{\alpha,\hat{\nu}_Q}(x) = \int_0^1 P_{\alpha}(\hat{F}_n(u), x) \, du = \sum_{i=1}^{n+1} \Delta X_{(i)} \, P_{\alpha}(\frac{i-1}{n}, x) \,,$$

$$\hat{q}_{\alpha}(x) = \int_{0}^{1} g(\hat{F}_{n}(u), [\alpha x] + 2, \alpha) \, du = \sum_{i=1}^{n+1} \Delta X_{(i)} \, g\left(\frac{i-1}{n}, [\alpha x] + 2, \alpha\right),$$

and

$$\hat{q}_{\alpha}^{*}(x) = \sum_{i=1}^{n+1} \Delta X_{(i)} \beta\left(\frac{i-1}{n}, [\alpha x] + 1, \alpha - [\alpha x] + 1\right).$$

Here  $\hat{\nu}_Q = \{\int_0^1 [\hat{F}_n(u)]^j du, j \in \mathbb{N}\}$ , while  $X_{(i)}, i = 1, ..., n, X_{(0)} = 0, X_{(n+1)} = 1$ , are the order statistics of the sample  $X_1, ..., X_n$ .

Now, let us compare the curves of  $\hat{Q}_{\alpha}$  and the well known Harrell-Davis estimator

49  
50 (4.7) 
$$\hat{Q}_{HD}(x) = \sum_{i=1}^{n} X_{(i)} \Delta Beta\left(\frac{i}{n}, (n+1)x, (n+1)(1-x)\right),$$
51  
51  
51  
51  
51  
51

imsart-coll ver. 2008/08/29 file: Mnatsak\_Revised.tex date: April 10, 2009



FIG 2. (a) Approximation of Q by  $Q_{\alpha,\beta}$  and (b) Estimation of  $Q(x) = x^{1/3}$  by  $\hat{Q}_{\alpha}$  and by  $\hat{Q}_{HD}$ 

where  $Beta(\cdot, a, b)$  denotes the cdf of a Beta distribution with the shape parameters a > 0 and b > 0. For asymptotic expressions of MSE and the bias term of  $\hat{Q}_{HD}$  we refer the reader to Sheather and Marron [19]. Let us generate n = 100 independent random variables  $X_1, \ldots, X_n$  from  $F(x) = x^3, 0 \le x \le 1$ . Taking  $\alpha = 100$ , we estimate (see, Figure 2 (b)) the corresponding quantile function  $Q(x) = x^{1/3}, 0 \le x \le 1$ , (the dashed line) by means of  $\hat{Q}_{\alpha}$  (the solid line) and by  $\hat{Q}_{HD}$  (the dashed-dotted line), defined in (4.6) and (4.7), accordingly. Through simulations we conclude that the asymptotic behavior of the moment-recovered estimator  $\hat{Q}_{\alpha}$  and the Harrell-Davis estimator  $\hat{Q}_{HD}$  are similar. The MSE and other properties of  $\hat{Q}_{\alpha}, \hat{q}_{\alpha}$ , and  $\hat{q}^*_{\alpha}$  will be presented in a separate article.

Example 4.1 (continued). Assume now that we want to recover pdf of the distribution G studied in the Example 4.1 via the moments  $\nu_{j,G} = 1/(j+1)^2, j \in \mathbb{N}$ . On the Figure 3 (a) we plotted the curves of the moment-recovered density  $f_{\alpha,\nu}$ (the solid line) defined by (3.11) and  $g(x) = G'(x) = -\ln x, \ 0 \leq x \leq 1$  (the dashed line), respectively. Here we took  $\alpha = 50$  and M = 200.

Example 4.2 (continued). Now let us recover the pdf  $f(x) = -9x^2 \ln x, 0 \le x \le 1$ , of distribution F defined in Example 4.2 where  $\nu_{j,F} = 9/(j+3)^2, j \in \mathbb{N}$ . We applied the approximations  $f_{\alpha,\nu}$  and  $f^*_{\alpha,\nu}$  defined in (3.11) and (3.14), respectively, by calculating the values of  $f_{\alpha,\nu}$  and  $f^*_{\alpha,\nu}$  at the points  $x = k/\alpha, k = 1, 2, \ldots, \alpha$ . Figure 3 (b) shows the curves of  $f_{\alpha,\nu}$  (the blue dashed-dotted line), and  $f^*_{\alpha,\nu}$  (the red solid line), and f (the black dashed line). Here, we took  $\alpha = 50$  and M = 200when calculating  $f_{\alpha,\nu}$  and  $\alpha = 32$  in  $f^*_{\alpha,\nu}$ . One can see that the performance of  $f^*_{\alpha,\nu}$ with  $\alpha = 32$  is better than the performance of  $f_{\alpha,\nu}$  with  $\alpha = 50$  and M = 200.

After conducting many calculations of moment-recovered approximants for several models we conclude that the accuracy of the formulas (2.2) and (3.11) are not as good as the ones defined in (3.13) and (3.14) in the Hausdorff case. On the other hand, the constructions (2.2) and (3.11) could be useful in the Sieltjes moment problem as well. 

FIG 3. (a) Approximation of  $g(x) = -\ln x$  by  $f_{\alpha,\nu}$  and (b) Approximation of  $f(x) = -9x^2 \ln x$  by  $f_{\alpha,\nu}$  and  $f^*_{\alpha,\nu}$ 

#### Acknowledgments

The authors are thankful to Estate Khmaladze, E. James Harner, and Cecil Burchfiel for helpful discussions. The work of the first author was supported by the National Institute for Occupational Safety and Health, contract no. 212-2005-M-12857. The findings and conclusions in this paper are those of the authors and do not necessarily represent the views of the National Institute for Occupational Safety and Health.

### References

[1] ADELL, J. A. and DE LA CAL, J. (1993). On the uniform convergence of normalized Poisson mixtures to their mixing distribution. J. Statist. Prob. Letters, 18, 227-232. AKHIEZER, N. I. (1965). The Classical Moment Problem and Some Related Questions in Analysis. [2]Oliver & Boyd, Edinburgh. BORWEIN, J. M. and LEWIS, A. S. (1993). A survey of convergence results for maximum entropy [3] methods. In Djafari, M. and Demoments, G., eds. Maximum Entropy and Bayesian Methods. Kluwer Academic Publishers, 39-48. CHAUVEAU, D. E., VAN ROOLJ, A. C. M. and RUYMGAART, F. H. (1994). Regularized inversion of noisy [4]Laplace transforms. Advances in Appl. Math., 15, 186–201. CHEN, S. X. (2000). Probability density function estimation using gamma kernels. Ann. Inst. [5]Statist. Math., 52, 471-480. [6] FELLER, W. (1971). An Introduction to Probability Theory and its Applications. vol. II. Wiley, New York. FRONTINI, M. and TAGLIANI, A. (1997). Entropy-convergence in Stieltejs and Hamburger moment [7]problem. Applied Math. and Computation, 88, 39-51. HALL, P. and LAHIRI, S. N. (2008). Estimation of distributions, moments and quantiles in deconvolution problems. Ann. Statist., 36, 2110-2134. KEVASAN, H. K. and KAPUR, J. N. (1992). Entropy Optimization Principles with Applications. Academic Press. New York. [10] LIN, G. D. (1992). Characterizations of distributions via moments. Sankhya, 54, Series A, 128–132. [11]LIN, G. D. (1997). On the moment problem. J. Statist. Prob. Letters, 35, 85-90. [12] LINDSAY, B. G., PILLA, R. S. and BASAK, P. (2000). Moment-based approximations of distributions using mixtures: theory and applications. Ann. Inst. Statist. Math., 52, 215-230. [13]MNATSAKANOV, R. M. (2008). Hausdorff moment problem: reconstruction of distributions. J. Statist. Prob. Letters, 78, 1612-1618.

| 1        | [14] | MNATSAKANOV, R. M. (2008). Hausdorff moment problem: Reconstruction of probability density   | 1         |
|----------|------|--|-----------|
| 2        | [15] | functions. J. Statist. Prob. Letters, <b>78</b> , 1869–1877.<br>MNATSAKANOV R M and KHMALADZE E V (1981) On L <sub>1</sub> -convergence of statistical kernel esti-  | 2         |
| 3        | [10] | mators of distribution densities. Soviet Math. Dokl., 23, 633–636.   | 3         |
| 4        | [16] | MNATSAKANOV, R. M. and KLAASSEN, C. A. J. (2003). Estimation of the mixing distribution in   | 4         |
| 5        |      | ference on Statistics and Related Fields, Honolulu, Hawaii, June 4-8, 2003, 1-18, CD: ISSN#  | 5         |
| 6        |      | 1539–7211.   | 6         |
| 7        | [17] | MNATSAKANOV, R. M. and RUYMGAART, F. H. (2003). Some properties of moment-empirical cdf's  | 7         |
| 8        | [18] | with application to some inverse estimation problems. <i>Math. Meth. Statist.</i> , <b>12</b> , 478–495.<br>NOVI INVERARDI, P. L., PETRI, A., PONTUALE, G. and TAGLIANI, A. (2003). Hausdorff moment prob- | 8         |
| 9        | [=~] | lem via fractional moments. Applied Mathematics and Computation, 144, 61–74.   | 9         |
| 10       | [19] | SHEATHER, S. J. and MARRON, J. S. (1990). Kernel quantile estimators. JASA, 85, No 410, 410–416.   | 10        |
| 11       | [20] | SHOHAT, J. A. and TAMARKIN, J. D. (1943). The Problem of Moments. Amer. Math. Soc., New York.  | 11        |
| 12       | [21] | STIELTJES, T. (1894). Recherches sur les fractions continues. Anns Sci. Fac. Univ. Toulouse (1894-   | 12        |
| 13       | [99] | 1895), 8 J1-J122; A5-A47.  | 13        |
| 14       | [22] | STOYANOV, J. (1997). Counterexamples in Probabilistic moment problems. Bernoulli, 6, 939–949.  | 14        |
| 15       | [24] | STOYANOV, J. (2004). Stieltjes classes for moment indeterminate probability distributions. J. Appl.  | 15        |
| 16       | [95] | Probab., <b>41A</b> , 281–294.   | 16        |
| 17       | [25] | moments. Applied Mathematics and Computation, 143, 99–107.   | 17        |
| 18       | [26] | TEICHER, H. (1963). Identifiability of finite mixtures. Ann. Math. Statist., 34, 1265–1269.  | 18        |
| 19       | [27] | WIDDER, D. V. (1934). The inversion of the Laplace integral and the related moment problem.  | 19        |
| 20       |      | Transactions of the Amer. Math. Soc., 30, 107–200.   | 20        |
| 21       |      |  | 21        |
| 22       |      |  | 22        |
| 23       |      |  | 23        |
| 24       |      |  | 24        |
| 25       |      |  | 25        |
| 26       |      |  | 26        |
| 27       |      |  | 27        |
| 28       |      |  | 28        |
| 29       |      |  | 29        |
| 30       |      |  | 30        |
| 31       |      |  | 31        |
| 32       |      |  | 32        |
| 33       |      |  | 33        |
| 34<br>25 |      |  | 34        |
| 20       |      |  | 20        |
| 37       |      |  | 30        |
| 38       |      |  | 30        |
| 20       |      |  | 20        |
| 40       |      |  | 40        |
| 41       |      |  | 41        |
| 40       |      |  | 41        |
| 43       |      |  | 42        |
| 44       |      |  | 44        |
| 45       |      |  | 45        |
| 46       |      |  | 46        |
| 47       |      |  | -10<br>47 |
| 48       |      |  | 48        |
| 49       |      |  | 49        |
| 50       |      |  | 50        |
| 51       |      |  | 51        |
|          |      |  |           |
**IMS** Collections Vol. 0 (2009) 266-275 © Institute of Mathematical Statistics, 2009 arXiv: math.PR/0000013

# Asymptotic Efficiency of Simple Decisions for the Compound Decision Problem

З

## Eitan Greenshtein<sup>1,\*</sup> and Ya'acov Ritov<sup>2,\*</sup>

Duke University and Jerusalem, Israel

Abstract: We consider the compound decision problem of estimating a vector of n parameters, known up to a permutation, corresponding to n independent observations, and discuss the difference between two symmetric classes of estimators. The first and larger class is restricted to the set of all permutation invariant estimators. The second class is restricted further to simple symmetric procedures. That is, estimators such that each parameter is estimated by a function of the corresponding observation alone. We show that under mild conditions, the minimal total squared error risks over these two classes are asymptotically equivalent up to essentially O(1) difference. -266Dense  $\mu$ 's  $\ldots \ldots 270$ 

# Contents

 $\mathbf{2}$ 

 $\mathbf{5}$ 

| 1. | Introduction |
|----|--------------|
|----|--------------|

Let  $\mathcal{F} = \{F_{\mu} : \mu \in \mathcal{M}\}$  be a parameterized family of distributions. Let  $Y_1, Y_2...$  be a sequence of independent random variables, where  $Y_i$  takes value in some space  $\mathcal{Y}$ , and  $Y_i \sim F_{\mu_i}$ ,  $i = 1, 2, \ldots$  For each n, we suppose that the sequence  $\mu_{1:n}$  is known up to a permutation, where for any sequence  $x = (x_1, x_2, ...)$  we denote the subsequence  $x_s, \ldots, x_t$  by  $x_{s:t}$ . We denote by  $\boldsymbol{\mu} = \boldsymbol{\mu}_n$  the set  $\{\mu_1, \ldots, \mu_n\}$ , i.e.,  $\boldsymbol{\mu}$  is  $\mu_{1:n}$ without any order information. We consider in this note the problem of estimating  $\mu_{1:n}$  by  $\hat{\mu}_{1:n}$  under the loss  $\sum_{i=1}^{n} (\hat{\mu}_i - \mu_i)^2$ , where  $\hat{\mu}_{1:n} = \Delta(Y_{1:n})$ . We assume that the family  $\mathcal{F}$  is dominated by a measure  $\nu$ , and denote the corresponding densities simply by  $f_i = f_{\mu_i}$ , i = 1, ..., n. The important example is, as usual,  $F_{\mu_i} = N(\mu_i, 1).$ 

Let  $\mathcal{D}^S = \mathcal{D}_n^S$  be the set of all simple symmetric decision functions  $\Delta$ , that is, all  $\Delta$  such that  $\Delta(Y_{1:n}) = (\delta(Y_1), \ldots, \delta(Y_n))$ , for some function  $\delta : \mathcal{Y} \to \mathcal{M}$ . In particular, the best simple symmetric function is denoted by  $\Delta^{S}_{\mu} = (\delta^{S}_{\mu}(Y_{1}), \dots, \delta^{S}_{\mu}(Y_{n}))$ :

$$\Delta^{S}_{\boldsymbol{\mu}} = \operatorname*{arg\,min}_{\boldsymbol{\Delta} \in \mathcal{D}^{S}} \mathrm{E} \, ||\boldsymbol{\Delta} - \boldsymbol{\mu}_{1:n}||^{2},$$

<sup>1</sup>Department of Statistical Sciences, Duke University, Durham, NC 27708-0251, USA, email: eitan.greenshtein@gmail.com 

<sup>&</sup>lt;sup>2</sup>Jerusalem, Israel, email: yaacov.ritov@gmail.com 

<sup>\*</sup>Partially supported by NSF grant DMS-0605236, and an ISF Grant.

AMS 2000 subject classifications: Primary 62C25; secondary 62C12, 62C07

Keywords and phrases: Compound decision, Simple decision rules, Permutation invariant rules

and denote

$$r_n^S = \mathbf{E} ||\Delta^S_{\mu}(Y_{1:n}) - \mu_{1:n}||^2,$$
 2

where, as usual,  $||a_{1:n}||^2 = \sum_{i=1}^n a_i^2$ .

The class of simple rules may be considered too restrictive. Since the  $\mu$ s are known up to a permutation, the problem seems to be of matching the Ys to the  $\mu$ s. Thus, if  $Y_i \sim N(\mu_i, 1)$ , and n = 2, a reasonable decision would make  $\hat{\mu}_1$  closer to  $\mu_1 \wedge \mu_2$  as  $Y_2$  gets larger. The simple rule clearly remains inefficient if the  $\mu$ s are well separated, and generally speaking, a bigger class of decision rules may be needed to obtain efficiency. However, given the natural invariance of the problem, it makes sense to be restricted to the class  $\mathcal{D}^{PI} = \mathcal{D}_n^{PI}$  of all permutation invariant decision functions, i.e, functions  $\Delta$  that satisfy for any permutation  $\pi$  and any  $(Y_1, \ldots, Y_n)$ :

$$\Delta(Y_1,\ldots,Y_n) = (\hat{\mu}_1,\ldots,\hat{\mu}_n) \quad \Longleftrightarrow \quad \Delta(Y_{\pi(1)},\ldots,Y_{\pi(n)}) = (\hat{\mu}_{\pi(1)},\ldots,\hat{\mu}_{\pi(n)}).$$

Let

$$\Delta_{\boldsymbol{\mu}}^{PI} = \operatorname*{arg\,min}_{\Delta \in \mathcal{D}^{PI}} \mathrm{E}\, ||\Delta(Y^n) - \mu_{1:n}||^2$$

be the optimal permutation invariant rule under  $\mu$ , and denote its risk by

$$r_n^{PI} = E ||\Delta_{\mu}^{PI}(Y_{1:n}) - \mu_{1:n}||^2.$$

Obviously  $\mathcal{D}^S \subset \mathcal{D}^{PI}$ , and whence  $r_n^S \geq r_n^{PI}$ . Still, 'folklore', theorems in the spirit of De Finetti, and results like Hannan and Robbins (1955), imply that asymptotically (as  $n \to \infty$ )  $\Delta_{\mu^n}^{PI}$  and  $\Delta_{\mu^n}^S$  will have 'similar' mean risks:  $r_n^S - r_n^{PI} = o(n)$ . Our main result establishes conditions that imply the stronger claim,  $r_n^S - r_n^{PI} = O(1)$ .

To repeat,  $\mu$  is assumed known in this note. In the general decision theory framework the unknown parameter is the order of its member to correspond with  $Y_{1:n}$ , and the parameter space, therefore, corresponds to the set of all the permutations of  $1, \ldots, n$ .

An asymptotic equivalence as above implies, that when we confine ourselves to the class of permutation invariant procedures, we may further restrict ourselves to the class of simple symmetric procedures, as is usually done in the standard analysis of compound decision problems. The later class is smaller and simpler.

The motivation for this paper stems from the way the notion of oracle is used in some sparse estimation problems. Consider two oracles, both know the value of  $\mu$ . Oracle I is restricted to use only a procedure from the class  $\mathcal{D}^{PI}$ , while Oracle II is further restricted to use procedures from  $\mathcal{D}^S$ . Obviously Oracle I has an advantage, our results quantify this advantage and show that it is asymptotically negligible. Furthermore, starting with Robbins (1951) various oracle-inequalities were obtained showing that one can achieve nearly the risk of Oracle II, by a 'legitimate' statistical procedure. See, e.g., the survey Zhang (2003), for oracle-inequalities regarding the difference in risks. See also Brown and Greenshtein (2007), and Wenuha and Zhang (2007) for oracle inequalities regarding the ratio of the risks. However, Oracle II is limited, and hence, these claims may seem to be too weak. Our equivalence results, extend many of those oracle inequalities to be valid also with respect to Oracle I. We needed a stronger result than the usual objective that the mean risks are equal up to o(1) difference. Many of the above mentioned recent applications of the compound decision notion are about sparse situations when most of the  $\mu$ s are in fact 0, the mean risk is o(1), and the only interest is in total risk. 

З

model under which  $(\pi, Y_{1:n}), \pi$  a random permutation, have a distribution given by

Let  $\mu_1, \ldots, \mu_n$  be some arbitrary ordering of  $\mu$ . Consider now the Bayesian

З

 $\pi$  is uniformly distributed over  $\mathcal{P}(1:n)$ ; (1.1)Given  $\pi$ ,  $Y_{1:n}$  are independent,  $Y_i \sim F_{\mu_{\pi(i)}}, i = 1, \ldots, n$ , where for every s < t,  $\mathcal{P}(s : t)$  is the set of all permutations of  $s, \ldots, t$ . The above description induces a joint distribution of  $(M_1, \ldots, M_n, Y_1, \ldots, Y_n)$ , where  $M_i \equiv \boldsymbol{\mu}_{\pi(i)}$ , for a random permutation  $\pi$ . The first part of the following proposition is a simple special case of general theorems representing the best invariant procedure under certain groups as the Bayesian decision with respect to the appropriate Haar measure; for background see, e.g., Berger (1985), Chapter 6. The second part of the proposition was derived in various papers starting with Robbins (1951). In the following proposition and proof,  $E_{\mu_{1:n}}$  is the expectation under the model in which the observations are independent,  $Y_i \sim F_{\mu_i}$ , and  $E_{\mu}$  is the expectation under the above joint distribution of  $Y_{1:n}$  and  $M_{1:n}$ . Note that under the latter model, for any  $i = 1, \ldots, n$ , marginally  $M_i \sim \mathbb{G}_n$ , the empirical measure defined by the vector  $\boldsymbol{\mu}$ , and conditional on  $M_i = m, Y_i \sim F_m$ . **Proposition 1.1.** The best simple and permutation invariant rules are given by (i)  $\Delta_{\mu}^{PI}(Y_{1:n}) = E_{\mu}(M_{1:n}|Y_{1:n}).$ 

(*ii*) 
$$\Delta_{\mu}^{S}(Y_{1:n}) = (E_{\mu}(M_{1}|Y_{1}), \dots, E_{\mu}(M_{n}|Y_{n})).$$

(*iii*) 
$$r_n^S = r_n^{PI} + \mathbf{E}_{\boldsymbol{\mu}} \| \Delta_{\boldsymbol{\mu}}^S - \Delta_{\boldsymbol{\mu}}^{PI} \|^2$$

*Proof.* We need only to give the standard proof of the third part. First, note that by invariance  $\Delta_{\mu}^{PI}$  is an equalizer (over all the permutations of  $\mu$ ), and hence  $E_{\mu_{1:n}}(\Delta_{\mu}^{PI}-\mu_{1:n})^2 = E_{\mu}(\Delta_{\mu}^{PI}-M_{1:n})^2$ . Also  $E_{\mu_{1:n}}(\Delta_{\mu}^{S}-\mu_{1:n})^2 = E_{\mu}(\Delta_{\mu}^{S}-M_{1:n})^2$ . Then, given the above joint distribution,

$$r_n^S = \mathbf{E}_{\boldsymbol{\mu}} \| \Delta_{\boldsymbol{\mu}}^S - M_{1:n} \|^2$$

$$= \mathbf{E}_{\mu} \mathbf{E}_{\mu} \{ \| \Delta_{\mu}^{S} - M_{1:n} \|^{2} | Y_{1:n} \}$$

$$= \mathbf{E}_{\mu} \mathbf{E}_{\mu} \{ \| \Delta_{\mu}^{S} - \Delta_{\mu}^{PI} \|^{2} + \| \Delta_{\mu}^{PI} - M_{1:n} \|^{2} |Y_{1:n} \}$$

$$= r_n^{PI} + \mathbf{E}_{\boldsymbol{\mu}} \| \Delta_{\boldsymbol{\mu}}^S - \Delta_{\boldsymbol{\mu}}^{PI} \|^2.$$

We now briefly review some related literature and problems. On simple symmetric functions, compound decision and its relation to empirical Bayes, see Samuel (1965), Copas (1969), Robbins (1983), Zhang (2003), among many other papers.

Hannan and Robbins (1955) formulated essentially the same equivalence problem in testing problems, see their Section 6. They show for a special case an equivalence up to o(n) difference in the 'total risk' (i.e., non-averaged risk). Our results for estimation under squared loss are stated in terms of the total risk and we obtain O(1) difference.

Our results have a strong conceptual connection to De Finetti's Theorem. The exchangeability induced on  $M_1, \ldots, M_n$ , by the Haar measure, implies 'asymptotic independence' as in De Finetti's theorem, and consequently asymptotic indepen-dence of  $Y_1, \ldots, Y_n$ . Thus we expect  $E(M_1|Y_1)$  to be asymptotically similar to  $E(M_1|Y_1,\ldots,Y_n)$ . Quantifying this similarity as n grows, has to do with the rate of convergence in De Finetti's theorem. Such rates were established by Diaconis and Freedman (1980), but are not directly applicable to obtain our results. 

imsart-coll ver. 2008/08/29 file: Ritov.tex date: April 10, 2009

After quoting a simple result in the following section, we consider in Section 3the special important, but simple, case of two-valued parameter. In Section 4 we obtain a strong result under strong conditions. Finally, the main result is given in Section 5, it covers the two preceding cases, but with some price to pay for the generality.

## 2. Basic Lemma and Notation

The following lemma is standard in comparison of experiments theory; for background on comparison of experiments in testing see Lehmann (1986), p-86. The proof follows a simple application of Jensen's inequality.

**Lemma 2.1.** Consider two pairs of distributions,  $\{G_0, G_1\}$  and  $\{\tilde{G}_0, \tilde{G}_1\}$ , such that the first pair represents a weaker experiment in the sense that there is a Markov kernel K, and  $G_i(\cdot) = \int \mathbb{K}(y, \cdot) dG_i(y), i = 1, 2$ . Then

$$\operatorname{E}_{G_0}\psi\big(\frac{dG_1}{dG_0}\big) \le \operatorname{E}_{\tilde{G}_0}\psi\big(\frac{d\tilde{G}_1}{d\tilde{G}_0}\big)$$

for any convex function  $\psi$ 

For simplicity denote  $f_i(\cdot) = f_{\mu_i}(\cdot)$ , and for any random variable X, we may write  $X \sim g$  if g is its density with respect to a certain dominating measure. Finally, for simplicity we use the notation  $y_{-i}$  to denote the sequence  $y_1, \ldots, y_n$  without its i member, and similarly  $\boldsymbol{\mu}_{-i} = \{\mu_1, \dots, \mu_n\} \setminus \{\mu_i\}$ . Finally  $f_{-i}(Y_{-j})$  is the marginal density of  $Y_{-j}$  under the model (1.1) conditional on  $M_j = \mu_i$ .

## 3. Two Valued Parameter

We suppose in this section that  $\mu$  can get one of two values which we denote by  $\{0,1\}$ . To simplify notation we denote the two densities by  $f_0$  and  $f_1$ .

**Theorem 3.1.** Suppose that either of the following two conditions holds:

- (i)  $f_{1-\mu}(Y_1)/f_{\mu}(Y_1)$  has a finite variance under both  $\mu \in \{0, 1\}$ . (ii)  $\sum_{i=1}^{n} \mu_i/n \to \gamma \in (0, 1)$ , and  $f_{1-\mu}(Y_1)/f_{\mu}(Y_1)$  has a finite variance under one of  $\mu \in \{0, 1\}$ .

Then 
$$E_{\mu} \| \hat{\mu}^{S} - \hat{\mu}^{PI} \|^{2} = O(1)$$

*Proof.* Suppose condition (i) holds. Let  $K = \sum_{i=1}^{n} \mu_i$ , and suppose, WLOG, that  $K \leq n/2$ . Consider the Bayes model of (1.1). By Bayes Theorem

$$P(M_1 = 1|Y_1) = \frac{Kf_1(Y_1)}{Kf_1(Y_1) + (n - K)f_0(Y_1)}.$$

On the other hand

$$P(M_1 = 1|Y_{1:n})$$

$$=\frac{Kf_1(Y_1)f_{K-1}(Y_{2:n})}{Kf_2(Y_2)f_{X-1}(Y_{2:n}) + (n-K)f_2(Y_2)f_{X}(Y_{2:n})}$$

$$Kf_{1}(Y_{1})f_{K-1}(Y_{2:n}) + (n-K)f_{0}(Y_{1})f_{K}(Y_{2:n})$$

$$Kf_{1}(Y_{1}) \qquad (n-K)f_{0}(Y_{1}) \qquad (f_{K-1}(Y_{1}) - 1)^{-1}$$

$$= \frac{1}{Kf_1(Y_1) + (n-K)f_0(Y_1)} \left(1 + \frac{(N-M)f_0(Y_1)}{Kf_1(Y_1) + (n-K)f_0(Y_1)} \left(\frac{f_K}{f_{K-1}}(Y_{2:n}) - 1\right)\right)$$
48
49

$$= P(M_1 = 1|Y_1) \left( 1 + \gamma \left( \frac{f_K}{f_{K-1}}(Y_{2:n}) - 1 \right) \right)^{-1},$$

where, with some abuse of notation  $f_k(Y_{2:n})$  is the joint density of  $Y_{2:n}$  conditional on  $\sum_{j=2}^{n} \mu_j = k$ , and the random variable  $\gamma$  is in [0,1]. We prove now that  $f_K/f_{K-1}(Y_{2:n})$  converges to 1 in the mean square. We use Lemma 2.1 (with  $\psi$  the square) to compare the testing of  $f_K(Y_{2:k})$  vs.  $f_{K-1}(Y_{2:k})$  to an easier problem, from which the original problem can be obtained by adding a random permutation. Suppose for simplicity and WLOG that in fact  $Y_{2:K}$  are i.i.d. under  $f_1$ , while  $Y_{K+1:n}$  are i.i.d. under  $f_0$ . Then we compare  $g_{K-1}(Y_{2:n}) = \prod_{j=2}^{K} f_1(Y_j) \prod_{j=K+1}^{n} f_0(Y_j),$ the true distribution, to the mixture  $g_K(Y_{2:n}) = g_{K-1}(Y_{2:n}) \frac{1}{n-K} \sum_{i=K+1}^n \frac{f_1}{f_0}(Y_j).$ However, the likelihood ratio between  $g_K$  and  $g_{K-1}$  is a sum of n - K terms, each with mean 1 (under  $g_{K-1}$ ) and finite variance. The ratio between the gs is, therefore,  $1 + O_p(n^{-1/2})$  in the mean square. By Lemma 2.1, this applies also to the fs' ratio. Consider now the second condition. By assumption, K is of the same order as n, and we can assume, WLOG, that the  $f_1/f_0$  has a finite variance under  $f_0$ . With this understanding, the above proof holds for the second condition. The condition of the theorem is clearly satisfied in the normal shift model:  $F_i =$  $N(\mu_i, 1), i = 1, 2$ . It is satisfied for the normal scale model,  $F_i = N(0, \sigma_i^2), i = 1, 2,$ if K is of the same order as n, or if  $\sigma_0^2/2 < \sigma_1^2 < 2\sigma_0^2$ .

## 4. Dense $\mu$ 's

We consider now another simple case in which  $\mu$  can be ordered  $\mu_{(1)}, \ldots, \mu_{(n)}$  such that the difference  $\mu_{(i+1)} - \mu_{(i)}$  is uniformly small. This will happen if, for example,  $\mu$  is in fact a random sample from a distribution with density with respect to Lebesgue measure, which is bounded away from 0 on its support, or more generally, if it is sampled from a distribution with short tails. Denote by  $Y_{(1)}, \ldots, Y_{(n)}$  and  $f_{(1)}, \ldots, f_{(n)}$  the Ys and fs ordered according to the  $\mu$ s. We assume in this section

(B1) For some constants  $A_n$  and  $V_n$  which are bounded by a slowly converging to infinite sequence:

$$\max|\mu_i - \mu_j| = A_n,$$

 $\operatorname{Var}\left(\frac{f_{(j+1)}}{f_{(j)}}(Y_{(j)})\right) \leq \frac{V_n}{n^2}.$ 

Note that condition (B1) holds for both the normal shift model and the normal scale model, if  $\mu$  behaves like a sample from a distribution with a density as above.

**Theorem 4.1.** If Assumption (B1) holds then

$$=O_{r}(A^{2}V^{2}/n).$$

$$\sum_{i=1}^{n} |\hat{\mu}_{i}^{PI} - \hat{\mu}_{i}^{S}|^{2} = O_{p}(A_{n}^{2}V_{n}^{2}/n).$$

З

*Proof.* By definition

$$\hat{\mu}_1^S = \frac{\sum_{i=1}^n \mu_i f_i(Y_i)}{\sum_{i=1}^n f_i(Y_i)},$$

$$\hat{\mu}_{1}^{PI} = \frac{\sum_{i=1}^{n} \mu_{i} f_{i}(Y_{1}) f_{-i}(Y_{2:n})}{\sum_{i=1}^{n} f_{i}(Y_{1}) f_{-i}(Y_{2:n})},$$

where 
$$f_{-i}$$
 is the density of  $Y_{2:n}$  under  $\mu_{-i}$ :

$$f_{-i}(y_{2:m}) = \frac{1}{(n-1)!} \sum_{\pi \in \mathcal{P}(2:n)} \prod_{j=2}^{n} f_{\pi(j)}(y_j) \frac{f_1}{f_i}(y_i).$$

The result will follow if we argue that

(4.1) 
$$|\mu_1^{PI} - \mu_1^S| \le \max_{i,j} |\mu_i - \mu_j| \Big( \max_{i,j} \frac{f_{-i}}{f_{-j}}(Y_{2:n}) - 1 \Big) = O_p(A_n V_n/n).$$

That is,  $\max_i |f_{-i}(Y_{2:n})/f_{-1}(Y_{2:n}) - 1| = O_p(V_n/n)$ . In fact we will establish a slightly stronger claim that  $||f_{-i} - f_{-1}||_{TV} = O_p(V_n/n)$ , where  $|| \cdot ||_{TV}$  denotes the total variation norm.

We will bound this distance by the distance between two other densities. Let  $g_{-1}(y_{2:n}) = \prod_{j=2}^{n} f_j(y_j)$ , the true distribution of  $Y_{2:n}$ . We define now a similar analog of  $f_{-i}$ . Let  $r_j$  and  $y_{(r_j)}$  be defined by  $f_j = f_{(r_j)}$  and  $y_{(r_j)} = y_j$ ,  $j = 1, \ldots, n$ . Suppose, for simplicity, that  $r_i < r_1$ . Let

$$g_{-i}(y_{2:n}) = g_{-1}(y_{2:n}) \prod_{j=r_i}^{r_1-1} \frac{f_{(j+1)}}{f_{(j)}}(y_{(j)}).$$
<sup>25</sup>
<sup>26</sup>
<sup>27</sup>

The case  $r_1 < r_i$  is defined similarly. Note that  $g_{-i}$  depends only on  $\mu_{-i}$ . Moreover, if  $\tilde{Y}_{2:n} \sim g_{-j}$ , then one can obtain  $Y_{2:n} \sim f_{-j}$  by the Markov kernel that takes  $\tilde{Y}_{2:n}$ to a random permutation of itself. It follows from Lemma 2.1

$$\|f_{-i} - f_{-1}\|_{TV} \le \|g_{-i} - g_{-1}\|_{TV}$$

$$= \mathbf{E}_{\mu_{2:n}} \left| \frac{g_{-i}}{g_{-1}} (Y_{2:n}) - 1 \right|$$

$$r_{1} - 1 \quad \epsilon \qquad 35$$

$$= \mathbf{E}_{\mu_{2:n}} \Big| \prod_{j=k}^{r_1-1} \frac{f_{(j+1)}}{f_{(j)}} (Y_{(j)}) - 1 \Big|$$

But, by assumption

$$R_k = \prod_{j=k}^{r_1-1} \frac{f_{(j+1)}}{f_{(j)}}(Y_{(j)})$$

is a reversed  $L_2$  martingale, and it follows from Assumption (B1) that

$$\max_{k < r_1} |R_k - 1| = O_p(A_n V_n / n).$$
45
46

Similar argument applies to 
$$i, r_i > r_1$$
, yielding

$$\max_{i} \|f_{-i} - f_{-1}\|_{TV} = O_p(A_n V_n / n)$$
49

51 We established (4.1). The theorem follows.

| Wo posumo                  |   |
|----------------------------|---|
| we assume:                 |   |
| (G1) For some $C$          | $C < \infty: \max_{i \in \{1, \dots, n\}}  \mu_i  < C,$   |
| and $\max_{i,j}$           | $\in 1,, n \to L_{\mu_i}(f_{\mu_j}(Y_1)/f_{\mu_i}(Y_1))^2 < C.$ Also, there is $\gamma > 0$ such that       |
| $\min_{i,j\in 1,\ldots,n}$ | $P_{\mu_i}(f_{\mu_j}(Y_1)/f_{\mu_i}(Y_1) > \gamma) \ge 1/2.$  |
| (G2) The random            | n variables   |
|                            | $f_i(Y_i)$  |
|                            | $p_j(Y_i) = \frac{1}{\sum_{k=1}^n f_k(Y_i)},  i, j = 1, \dots, n.$  |
| are bounde                 | d in expectation by   |
|                            | $\prod_{n=1}^{n} \sum_{j=1}^{n} (-(\mathbf{x}_{j}))^{2} \rightarrow C$                                      |
|                            | $\mathrm{E}\sum\sum\left(p_{j}(Y_{i}) ight) \ < C$  |
|                            | i=1 $j=1$   |
|                            | $\sum_{n=1}^{n} \mathbf{F} = \frac{1}{2}$   |
|                            | $\sum_{i=1}^{L} \frac{1}{n \min_{j} p_{j}(Y_{i})} < Cn$   |
|                            | $n \sum^{n} ((\mathbf{x}))^2$   |
|                            | $\mathbf{E}\sum_{i=1}^{n} \frac{\sum_{j=1}^{n} (p_j(Y_i))}{C_i} < C_i$                                      |
|                            | $\sum_{i=1}^{n} n \min_{j} p_{j}(Y_{i}) $   |
|                            |   |
| Both assumpt               | ions describe a situation where the $\mu$ s do not "separate". They   |
| $(\mathbf{C1})$            | ar one from another, geometrically or statistically (Assumption   |
| (GI), and they             | are dense in the sense that each Y can be explained by many of $(Go)$                                       |
| the $\mu s$ (Assumpt       | ion (G2)). The conditions hold for the normal shift model if $\mu_n$  |
| are uniformly bo           | unded: Suppose the common variance is 1 and $ \mu_j  < A_n$ . Then  |
| $\frac{n}{n}$              | $f_{i}(V_{i}) \rightarrow 2 \qquad \sum_{i=1}^{n} f_{i}^{2}(Y_{i})$   |
| $E\sum_{n}$                | $\frac{JJ(1_{j})}{n} = E \frac{2J(1_{j})}{(\nabla^{n} \mathbf{f}(\mathbf{V}))^{2}}$                         |
| $\sum_{j=1}$               | $\sum_{k=1}^{k} J_k(I_1)^2 \qquad (\sum_{k=1}^{k} J_k(I_1))^2$  |
|                            | $ne^{-Y_1^2+2A_n Y_1 -A_n^2}$   |
|                            | $\leq \mathrm{E} \frac{1}{(ne^{-(Y_1^2 - 2A_n Y_1  + A_n^2)/2)2}}$  |
|                            | ( <i>IIC</i> ) 1 ( <i>IIC</i> )   |
|                            | $= \frac{1}{n} \mathbf{E}  e^{4A_n  Y_1 }$  |
|                            | n   |
|                            | $= \frac{1}{n} \left( e^{8A_n^2 + 4A_n\mu_1} + e^{8A_n^2 - 4A_n\mu_1} \right) \le \frac{1}{n} e^{12A_n^2}.$ |
|                            |   |
| and the first part         | G of (G2) hold. The other parts follow a similar calculations.  |
| Theorem 5.1.               | Assume that (G1) and (G2) hold. Then  |
| (i)                        | $\mathbf{E} \ \Delta_{\boldsymbol{\mu}}^{S} - \Delta_{\boldsymbol{\mu}}^{PI}\ ^{2} = O(1)$                  |
| (ii)                       | $r_n^S - r_n^{PI} = O(1).$  |
| Corollary 5 2              | Summary $\mathcal{F} = \int N(\mu, 1) \cdot  \mu  < c \int for some c < \infty$ then the                    |
| conclusions of th          | suppose $J = \{1, (\mu, 1),  \mu  \leq c\}$ for some $c \leq \infty$ , then the e theorem follow            |
| conclusions of th          |   |
| Proof. It was me           | ntioned already in the introduction that when we are restricted to  |
| permutation inva           | riant procedure we can consider the Bayesian model under which  |
| $(\pi V_1)$ $\pi$ a rate   | ndom permutation have a distribution given by $(11)$ Fix now  |

 $(\pi, Y_{1:n}), \pi$  a random permutation, have a distribution given by (1.1). Fix now  $i \in \{1, \ldots, n\}$ . Under this model we want to compare  $\mu^S_{s}$  $E(\dots | V)$ 

$$\mu_i^{\rm S} = E(\mu_{\pi(i)}|Y_i), \quad i = 1, \dots, n$$

49

50

imsart-coll ver. 2008/08/29 file: Ritov.tex date: April 10, 2009

272

49

50

 $\mathrm{to}$ 

$$\mu_i^{PI} = E(\mu_{\pi(i)}|Y_{1:n}), \quad i = 1, \dots, n.$$

More explicitly:

$$\mu_i^S = \frac{\sum_{j=1}^n \mu_j f_j(Y_i)}{\sum_{j=1}^n f_j(Y_i)}$$

$$=\sum_{j=1}^{n} \mu_{i} p_{i}(Y_{i}), \quad i=1,\ldots,n$$

(5.1) 
$$\sum_{j=1}^{n} \mu_{j} f_{j}(Y_{j}) = 0 \quad 11$$

$$\mu_i^{PI} = \frac{\sum_{j=1}^{n} \mu_j f_j(Y_i) f_{-j}(Y_{-i})}{\sum_{j=1}^{n} f_j(Y_i) f_{-j}(Y_{-i})}$$
<sup>13</sup>
<sup>14</sup>

$$=\sum_{j=1}^{n}\mu_{j}p_{j}(Y_{i})W_{j}(Y_{-i},Y_{i}), \quad i=1,\ldots,n,$$
15
16
17

where for all  $i, j = 1, ..., n, f_j(Y_i)$  was defined in Section 2, and

$$f(Y_{-i}, Y_i) = \frac{f_{-j}(Y_{-i})}{\sum^n x_i(Y_i) f_{-i}(Y_{-i})}$$
23
24

$$W_j(Y_{-i}, Y_i) = \frac{J_{-j}(Y_{-i})}{\sum_{k=1}^n p_k(Y_i) f_{-k}(Y_{-i})}$$

Note that  $\sum_{k=1}^{n} p_k(Y_i) = 1$ , and  $W_j(Y_{-i}, Y_i)$  is the likelihood ratio between two (conditional on  $Y_i$ ) densities of  $Y_{-i}$ , say  $g_{j0}$  and  $g_1$ . Consider two other densities (again, conditional on  $Y_i$ ):

$$\tilde{g}_{j0}(Y_{-i}|Y_i) = f_i(Y_j) \prod_{m \neq i,j} f_m(Y_m),$$

$$\tilde{g}_{j1}(Y_{-i}|Y_i) = \tilde{g}_{j0}(Y_{-i}|Y_i) \Big(\sum_{k \neq i,j} p_k(Y_i) \frac{f_j}{f_k}(Y_k) + p_i(Y_i) \frac{f_j}{f_i}(Y_j) + p_j(Y_i)\Big)$$

Note that  $g_{j0} = \tilde{g}_{j0} \circ \mathbb{K}$  and  $g_1 = \tilde{g}_{j1} \circ \mathbb{K}$ , where  $\mathbb{K}$  is the Markov kernel that takes  $Y_{-i}$  to a random permutation of itself. It follows from Lemma 2.1 that

$$\mathbf{E}(|W_j(Y_{-i}, Y_i) - 1|^2 | Y_i) \le \mathbf{E}_{\tilde{g}_{j1}} \left(\frac{\tilde{g}_{j0}}{\tilde{g}_{j1}} - 1\right)^2$$

(5.2)

$$= \mathbf{E}_{\tilde{g}_{j0}} \Big( \frac{\tilde{g}_{j0}}{\tilde{g}_{j1}} - 2 + \frac{\tilde{g}_{j1}}{\tilde{g}_{j0}} \Big).$$

This expectation does not depend on i except for the value of  $Y_i$ . Hence, to simplify notation, we take WLOG i = j. Denote

$$L = \frac{\tilde{g}_{j1}}{\tilde{g}_{i0}} = p_j(Y_j) + \sum_{k \neq i} p_k(Y_j) \frac{f_j}{f_k}(Y_k)$$
<sup>47</sup>
<sup>48</sup>

$$g_{j0}$$
  $f_{k\neq i}$   $f_{k}$  (49)

$$V = \frac{n}{4} \gamma \min_{k} p_k(Y_j),$$
50
51

imsart-coll ver. 2008/08/29 file: Ritov.tex date: April 10, 2009



З

where  $\gamma$  is as in (G1). Then by (5.2)

$$\mathbb{E}(|W_{j}(Y_{-j}, Y_{j}) - 1|^{2} | Y_{j}) \le \mathbb{E}_{\tilde{g}_{j0}}\left(\frac{1}{L} - 2 + L\right)$$

$$=\mathrm{E}_{ ilde{g}_{j0}}\,rac{(L-1)^2}{L}$$

$$\leq \frac{1}{V} \mathcal{E}_{\tilde{g}_{j0}}(L-1)^2 \mathbf{I}(L>V) + \mathcal{E}_{\tilde{g}_{j0}} \frac{\mathbf{I}(L\leq V)}{L}$$

$$\leq \operatorname{E}_{\tilde{g}_{j0}} \frac{\mathbf{I}(L \leq V)}{L} + \frac{1}{V} \sum_{k=1}^{n} p_k^2(Y_j),$$

by G1. Bound

$$L \ge \gamma \min_{k} p_k(Y_j) \sum_{k=1}^n \mathbf{1}\left(\frac{f_j}{f_k}(Y_k) > \gamma\right) \ge \gamma \min_{k} p_k(Y_j)(1+U),$$

where  $U \sim B(n-1, 1/2)$  (the 1 is for the *i*th summand). Hence

$$\mathcal{E}_{\tilde{g}_{j0}} \frac{\mathbf{I}(L \le V)}{L} \le \frac{1}{\gamma \min_k p_k(Y_j)} \sum_{k=0}^{\lceil n/4 \rceil} \frac{1}{k+1} \binom{n-1}{k} 2^{-n+1}$$

(5.4) 
$$= \frac{1}{2^{m\min(n)} (V)} \sum_{k=1}^{\lceil n/4 \rceil} {n \choose k+1} 2^{-n+1}$$

$$\gamma n \min_{k} p_{k}(Y_{j}) \xrightarrow{k=0} (k+1)$$
$$= O(e^{-n}) \frac{1}{\gamma n \min_{k} p_{k}(Y_{j})}$$

by large deviation.

From (G1), (G2), (5.1), (5.3), and (5.4):

$$EE((\mu_{i}^{S} - \mu_{i}^{PI})^{2} | Y_{i}) = EE\left(\left(\sum_{j=1}^{n} \mu_{j} p_{j}(Y_{i}) (W_{j}(Y_{-i}, Y_{i}) - 1)\right)^{2} | Y_{i}\right)$$

$$\leq \max_{j} |\mu_{j}|^{2} \operatorname{E} \sum_{j=1}^{n} p_{j}(Y_{i}) \operatorname{E} \left( W_{j}(Y_{-i}, Y_{i}) - 1 \right)^{2} \right) |Y_{i} \right)$$
  
 
$$\leq \kappa C^{3}/n,$$

$$C^3/n,$$

for some  $\kappa$  large enough. Claim (i) of the theorem follows. Claim (ii) follows (i) by Proposition 1.1.

## References

| [1] | BERGER, J. O. (1985). Statistical Decision Theory and Bayesian Analysis, 2 <sup>nd</sup> edition. Springer- |
|-----|---|
|     | Verlag, New York.   |
| [2] | BROWN, L. D. and GREENSHTEIN, E. (2007). Non parametric empirical Bayes and compound decision               |
|     | approaches to estimation of a high dimensional vector of normal means. Manuscript.                          |
| [3] | COPAS, J. B. (1969). Compound decisions and empirical Bayes (with discussion). JRSSB <b>31</b> 397–         |

[6] Immand, or 1 and Robbins, in (1999) insymptotic controls of one compound decision products for two completely specified distributions. Ann. Math. Stat., 26.1, 37–51.
[6] LEHMANN, E. L. (1986). Testing Statistical Hypothesis, 2<sup>nd</sup> edition. Wiley & Sons, New York.

З

| 1        | [7]   | ROBBINS, H. (1951). Asymptotically subminimax solutions of compound decision problems. Proc.   | 1         |
|----------|-------|--|-----------|
| 2        | [9]   | Third Berkeley Symp., 157–164.   | 2         |
| 3        | [8]   | SAMUEL, E. (1965). On simple rules for the compound decision problem. JRSSB, 27, 238–244.  | 3         |
| 4        | [10]  | ZHANG, C. H. (2003). Compound decision theory and empirical Bayes methods. (Invited paper.   | 4         |
| 5        | [1 1] | Ann. Stat., <b>31</b> , 379–390.   | 5         |
| 6        | [11]  | wenoha, J. and Zhang, C. H. (2007) General maximum likelihood empirical Bayes estimation of normal means. Manuscript.  | 6         |
| 7        |       | Final States of Final States o | 7         |
| 8        |       |  | 8         |
| 9        |       |  | 9         |
| 10       |       |  | 10        |
| 11       |       |  | 11        |
| 12       |       |  | 12        |
| 13       |       |  | 13        |
| 14       |       |  | 14        |
| 15       |       |  | 15        |
| 16       |       |  | 16        |
| 17       |       |  | 17        |
| 18       |       |  | 18        |
| 19       |       |  | 19        |
| 20       |       |  | 20        |
| 21       |       |  | 21        |
| 22       |       |  | 22        |
| 23       |       |  | 23        |
| 24       |       |  | 24        |
| 25       |       |  | 25        |
| 26       |       |  | 26        |
| 27       |       |  | 27        |
| 28       |       |  | 28        |
| 29       |       |  | 29        |
| 30       |       |  | 30        |
| 31       |       |  | 31        |
| 32       |       |  | 32        |
| 33       |       |  | 33        |
| 34       |       |  | 34        |
| 35       |       |  | 35        |
| 36       |       |  | 36        |
| 37       |       |  | 37        |
| 38       |       |  | 38        |
| 39       |       |  | 39        |
| 40       |       |  | 40        |
| 41       |       |  | 41        |
| 42       |       |  | 42        |
| 43       |       |  | 43        |
| 44       |       |  | 44        |
| 45       |       |  | 45        |
| 46       |       |  | 46        |
| 47       |       |  | 40<br>47  |
| 48       |       |  | 19        |
| 49       |       |  | 10        |
| 50       |       |  | -10<br>50 |
| 51       |       |  | 50        |
| <u> </u> |       |  | 01        |

# Large Sample Statistical Inference for Skew-Symmetric Families on the Real Line

# Rolando Cavazos–Cadena<sup>1,\*</sup> and Graciela González–Farías<sup>2,\*</sup>

Universidad Autónoma Agraria Antonio Narro and Centro de Investigación en Matemáticas A. C.

Abstract: For a general family of one-dimensional skew-symmetric probability densities, the application of the maximum likelihood method to the estimation of the asymmetry parameter  $\lambda$  is studied. Under mild conditions, the existence and consistency of a sequence  $\{\hat{\lambda}_n\}$  of maximum likelihood estimators is established, and the limit distributions of  $\{\hat{\lambda}_n\}$  and the sequence of likelihood ratios are determined under the null hypothesis  $\mathcal{H}_0: \lambda = 0$ . These latter conclusions, which hold under differential singularity of the likelihood function at  $\lambda = 0$ , extend to the present framework results recently obtained for general statistical models with null Fisher information.

Dedicated to Professor O. Hernández-Lerma, on the occasion of his sixtieth birthday

| 1        | Introduction                               |
|----------|--|
| 2        | Identifiability                            |
| 3        | Existence of Maximum Likelihood Estimators |
| 4        | Consistency                                |
| <b>5</b> | Asymptotic Distribution                    |
| 6        | Technical Preliminaries                    |
| 7        | Proof of Theorem 5.1                       |
| Re       | ferences                                   |

## 1. Introduction

This work concerns likelihood inference for the general one-dimensional skewsymmetric family, which is constructed as follows: Given two symmetric densities

\*This research was supported by the PSF Organization under Grant 2007-3, and by CONACYT under Grants 25357 and 45974-F.

AMS 2000 subject classifications: Primary 62F10; secondary 62F12

Keywords and phrases: Boundedness of maximum likelihood estimators, Kullback's inequality,
 Lateral Taylor series, Strong law of large numbers, Central limit theorem, Asymptotic normality
 51

<sup>&</sup>lt;sup>1</sup>Departamento de Estadística y Cálculo, Universidad Autónoma Agraria Antonio Narro, Buenavista, Saltillo COAH, 25315, **MÉXICO**, email: rcavazos@uaaan.mx

<sup>&</sup>lt;sup>2</sup>Centro de Investigación en Matemáticas A. C., Apartado Postal 402, Guanajuato, GTO, 36240, **MÉXICO**, email: farias@cimat.mx

f and g on the real line—that is, f(x) = f(-x) and g(x) = g(-x) for all  $x \in \mathbb{R}$ — let  $G(x) := \int_{-\infty}^{x} g(z) \, dz$ (1.1)

be the cumulative distribution function of density g. Notice that

$$\begin{split} &\int_{\mathbb{R}} [1 - G(\lambda w)] f(w) \, dw = \int_{\mathbb{R}} [1 - G(-\lambda w)] f(-w) \, dw = \int_{\mathbb{R}} G(\lambda w) f(w) \, dw \text{ for every} \\ &\lambda \in \mathbb{R}, \text{ so that } \int_{\mathbb{R}} 2f(w) G(\lambda w) \, dw = 1. \text{ This argument, due to Azzalini (1985),} \end{split}$$
shows that for each  $\lambda \in \mathbb{R}$ 

(1.2) 
$$\rho(x;\lambda) := 2f(x)G(\lambda x), \quad x \in \mathbb{R}$$

is a genuine density, and the collection

is the skew-symmetric family determined by f and g. When the parametrization  $\lambda \mapsto \rho(\cdot; \lambda)$  is one-to-one,  $f(\cdot) = \rho(\cdot; 0)$  is the unique symmetric density in S(f, q). and  $\lambda$  can be considered as a measure of the asymmetry (or skewness) of density  $\rho(\cdot; \lambda)$ . The first systematic treatment of a skew-symmetric family was presented in Azzalini (1985, 1986) for the case in which f and g coincide with the standard nor-mal density  $\varphi$ , and location-scale and regression models based on  $S(\varphi, \varphi)$  as well as multivariate extensions have been intensively studied during the last twenty years; see Azzalini and Dalla Valle (1996), Azzalini and Capitanio (1999), Pewsey (2000), Genton (2004) and the references therein. The analysis of the maximum likelihood method for the location-scale model based on  $S(\varphi, \varphi)$  has proved most challeng-ing since, in that context, the Fisher information matrix has incomplete rank at  $\lambda = 0$ , a problem that also arises for the skew exponential family, which includes the skew normal location-scale model as a particular case (Azzalini, 1986, DiCiccio and Monti, 2004). The singularity of the information matrix has been analyzed via the centered parametrization introduced in Azzalini (1985) and asymptotic results are based on the recent work by Rotzintzky et al. (2000). Using the conclusions in this latter paper, rates of convergence for maximum likelihood estimators are derived in Chiogna (2005), and Sartori (2006) studied the finiteness of the estimator of the asymmetry parameter.

As suggested by the previous comments, family  $S(\varphi, \varphi)$  has been intensively studied, and a great effort has been done on generalizations of that model (Genton, 2004), so that looking for inference results applicable to a broad class of skew families is, certainly, an interesting problem. This work is a first step in this direction since, although no scale or location parameters will be introduced, the maximum likelihood method applied to the estimation of the asymmetry parameter  $\lambda$  will be studied under rather minimal conditions on densities f and q, making the likelihood inference problem a very interesting one. The first objective of this note is

(i) To establish the existence of a sequence  $\{\hat{\lambda}_n\}$  of maximum likelihood estimators of  $\lambda$  and to prove its consistency. 

To see the interest behind this problem, denote the (kernel log-) likelihood corre-sponding to a single observation x by 

$$\ell_{48}^{47}$$
 (1.4)  $\ell(\lambda; x) := \log(G(\lambda x)),$ 

and observe that if g is the standard normal density, then  $\ell(\cdot; x)$  is strictly concave for  $x \neq 0$ , a property that yields the existence and consistency of maximum likeli-hood estimators (Newey and McFadden, 1993). However, strict concavity of  $\ell(\cdot; x)$ 

З

is far from being a general property, and may fail in common cases, for instance, if q is the Laplace density. In this work, problem (i) above will be studied under the minimal assumption that the parametrization  $\lambda \mapsto \rho(\cdot; \lambda)$  is identifiable and, since the parameter space is not compact for the S(f,g) family, to the best of the authors' knowledge, in this context the existence and consistency of maximum likelihood estimators can not be directly obtained form general available results. The second problem studied in this work concerns the asymptotic distribution of  $\{\lambda_n\}$ and, as usual, the analysis below requires differentiability assumptions on  $\ell(\cdot; x)$ , and then, on density q. This problem will be studied under conditions allowing qto be non smooth at x = 0, which can be roughly described as follows:

**A1:** q is continuous on  $\mathbb{R}$ , is 'smooth' outside 0, and has lateral derivatives at zero. 

Under this requirement, it is not difficult to see that if the true parameter value, say  $\nu$ , is non-null, then  $\ell(\cdot; x)$  is smooth on a neighborhood of  $\nu$  if  $\ell(\nu; x) < \infty$ . Thus, when the information number at  $\nu$ , given by 

 $\mathcal{I}(\nu) = \int_{\mathbb{R}} (\partial_{\lambda} \ell(\nu; x))^2 \rho(\nu; x) \, dx (> 0)$ , is finite, under standard regularity condi-tions Wald's classical results yield that  $\sqrt{n\mathcal{I}(\nu)}(\hat{\lambda}_n - \nu)$  has a standard normal distribution at the limit (Lehmann and Casella, 1998, Section 6.3, Shao, 1999, Section 4.4). However, under the null hypothesis  $\mathcal{H}_0: \lambda = 0$  such a direct conclu-sion is not possible, since  $q(\cdot)$  is not necessarily differentiable in a neighborhood of zero under condition A1 above. Also, observe that  $\partial_{\lambda}\ell(0;x) = 2g(0)x$ , so that  $\mathcal{I}(0) = 4g(0)^2 \int_B x^2 f(x) \, dx$ ; thus,  $\mathcal{I}(0)$  is null if g(0) = 0 and, again, in this case the asymptotic distribution of  $\{\hat{\lambda}_n\}$  can not be obtained from the results in the afore-mentioned references. The case of null information at  $\lambda = 0$  was recently studied in a general context by Rotnitzky et al. (2000) under several assumptions, includ-ing (a) compactness of the parameter space, and (b) the existence of higher order derivatives  $\partial_{\lambda}^{k}\ell(\lambda;x)$  in a neighborhood of zero. Since in the present context nei-ther the parameter space is compact, nor higher order derivatives  $\partial_{\lambda}^{k}\ell(\lambda;x)$  exist around zero, if q(0) = 0 then the limiting distribution can not be solved by direct application of available results under condition A1. Therefore, the second problem considered in this note is 

(ii) to determine both the limit distribution of the (appropriately normalized) se-quence  $\{\hat{\lambda}_n\}$  under the hypothesis  $\mathcal{H}_0: \lambda = 0$ , and the asymptotic null distribution of the likelihood ratio statistic. 

The approach used below to study problems (i) and (ii) can be described as follows: The monotonicity of the mapping  $\lambda \mapsto \ell(\lambda; x)$  is used to establish the existence of maximizers  $\hat{\lambda}_n$  of the observed likelihood when the sample size n is large enough, whereas the proof for the consistency of  $\{\hat{\lambda}_n\}$  follows the ideas in Cavazos-Cadena and Gonzalez-Farías (2007) where, under mild conditions, it is shown that sequence  $\{\hat{\lambda}_n\}$  is consistent if and only if it is bounded with probability 1; here, after estab-lishing the boundedness property, a direct proof of consistency is given via a simple consequence of Kullback's inequality. Concerning problem (ii), as in the results pre-sented in Lehmann and Casella (1998, Section 6.3), Shao (1999, Section 4.4), or in Rotnitzky et al. (2000), the analysis is based on Taylor series expansions for the observed likelihood and its derivative around zero, which in the present context are lateral expansions; they are used to show that the maximum likelihood estimator is no null with probability increasing to 1, and the asymptotic distributions are obtained via the central limit theorem and the strong law of large numbers. 

The organization of the paper is as follows: Problem (i) is analyzed in the following

three sections. Thus, in Section 2 the identification property of the parametrization  $\lambda \mapsto \rho(\lambda; \cdot)$  is discussed, the existence of a sequence  $\{\hat{\lambda}_n\}$  of maximum likelihood estimators is established in Section 3, and the consistency of  $\{\hat{\lambda}_n\}$  is proved in Section 4. After this point, the reminder of the paper concerns problem (ii). In Section 5 the basic dominance and smoothness assumptions are formally introduced, and the main asymptotic result is stated as Theorem 5.1. Next, in Section 6 the necessary technical tools concerning *lateral Taylor series* for the likelihood function and its first derivative are established and, finally, the exposition concludes in Section 7 with a proof of Theorem 5.1.

NOTATION. Throughout the remainder  $X_1, X_2, X_3, \ldots$  stands for a sequence of independent and identically distributed random variables with common density belonging to family S(f,g) in (1.3): For each  $n = 1, 2, \ldots$ , set

(1.5) 
$$X_1^n := (X_1, \dots, X_n).$$

The distribution of the sequence  $(X_1, X_2, ...)$  when  $\nu$  is the true parameter value is denoted by  $P_{\nu}[\cdot]$ , whereas  $E_{\nu}[\cdot]$  stands for the corresponding expectation operator. On the other hand, if  $H(\cdot)$  is a function defined around zero,  $H(0+) := \lim_{x \searrow 0} H(x)$ , and  $H(0-) := \lim_{x \nearrow 0} H(x)$ , whereas given  $A \subset \mathbb{R}$ ,  $I_A$  is the indicator function of set A, that is,  $I_A(x) = 1$  if  $x \in A$ , whereas  $I_A(x) = 0$  when  $x \notin A$ . Finally, the following convention is enforced:  $\sum_{i=a}^{b} C_i = 0$  when b < a.

## 2. Identifiability

In this section the identifiability of the parametrization  $\lambda \mapsto \rho(\cdot; \lambda)$  is briefly discussed. This condition establishes that different parameters correspond to different densities, and plays a fundamental role in parametric estimation (Newey and Mc-Fadden, 1993).

**Assumption 2.1.** The mapping  $\lambda \mapsto \rho(\cdot; \lambda)$  is one-to-one, that is,

$$\int_{\mathbb{R}} |\rho(x;\lambda) - \rho(x;\nu)| \, dx = \int_{\mathbb{R}} 2f(x) |G(\lambda x) - G(\nu x)| \, dx \neq 0, \quad \text{if } \lambda \neq \nu;$$

see (1.2).

Some primitive conditions ensuring this requirement are now given.

**Lemma 2.1.** Assumption 2.1 holds under either of the following conditions (i)-(iii).

(i) For every nonempty open interval  $J \subset \mathbb{R}$ ,  $\int_J f(x) dx > 0$ ;

(ii)  $\int_J g(x) dx > 0$  for each nonempty open interval  $J \subset \mathbb{R}$ ;

(iii) There exists  $\delta > 0$  such that

$$\int_J g(x) \, dx > 0 \quad and \quad \int_J f(x) \, dx > 0,$$

for each nonempty open interval  $J \subset (0, \delta)$ .

<sup>48</sup> <sup>49</sup> <sup>50</sup> Proof. Let  $\lambda$  and  $\nu$  be fixed and different real numbers, and suppose that

(2.1) 
$$\int_{\mathbb{R}} f(x) |G(\lambda x) - G(\nu x)| \, dx = 0.$$

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

A contradiction will be obtained under each of the three conditions in the lemma. Assume that condition (i) holds. Given an open interval  $J \subset \mathbb{R}$  with positive length, (2.1) and  $\int_{I} f(x) dx > 0$  together yield that  $|G(\lambda x) - G(\nu x)| = 0$  for some  $x \in J$ , so that  $\{x \mid G(\lambda x) - G(\nu x) = 0\}$  is dense in **R**. Since  $G(\cdot)$  is continuous, this set is also closed. so that 

(2.2) 
$$G(\lambda x) = G(\nu x) \text{ for all } x \in \mathbb{R}.$$

This fact implies that  $\lambda \neq 0$  and  $\nu \neq 0$ . Indeed, if  $\nu = 0$ , it follows that  $\lambda \neq 0$  and the right-hand side equals 1/2 for all x, whereas the values of left-hand side cover the whole interval (0,1) as x moves on IR. Thus,  $\nu \neq 0$  and, similarly,  $\lambda \neq 0$ . Moreover, (2.2) also yields that  $\lambda$  and  $\nu$  have the same sign since, otherwise, as  $x \to \infty$  one side of the equality converges to 1 and the other converges to 0. Therefore, recalling that  $\lambda \neq \nu$ , it follows that  $|\lambda| \neq |\nu|$ , and without loss of generality it can be assumed that  $\beta = |\lambda|/|\nu| = \lambda/\nu \in (0,1)$ . Replacing x by  $y/\nu$ , (2.2) yields that, for all  $y \in \mathbb{R}$ ,  $G(y) = G(\beta y)$ , and then

$$G(y) = G(\beta^n y), \quad y \in \mathbb{R}, \quad n = 1, 2, 3, \dots$$

Letting n go to  $\infty$ , the continuity of  $G(\cdot)$  and the inclusion  $\beta \in (0,1)$  yield that G(y) = G(0) = 1/2 for all  $y \in \mathbb{R}$ , in contradiction with the basic properties of a distribution function.

Under condition (ii),  $G(b) - G(a) = \int_a^b g(x) \, dx > 0$  for a < b, so that  $G(\cdot)$  is strictly increasing. Thus, since  $\lambda \neq \nu$ ,  $|G(\lambda x) - G(\nu x)| > 0$  for all  $x \neq 0$ , and it follows that (2.1) is equivalent to  $\int_{\mathbb{R}} f(x) dx = 0$ , which is not possible, since  $f(\cdot)$  is a density. Suppose that condition (iii) occurs. In this context,  $G(\cdot)$  is strictly increasing on the interval  $(0, \delta)$  and then on  $(-\delta, \delta)$ , by symmetry. Next, define  $\delta_1 := \delta/(|\lambda| + \delta_1)$  $|\nu|+1$ ) and, recalling that  $\lambda \neq \nu$ , notice that if  $x \in [\delta_1/2, \delta_1]$  then  $\lambda x$  and  $\nu x$  are different points in  $(-\delta, \delta)$ , so that  $|G(\lambda x) - G(\nu x)| > 0$ . Thus, by continuity of  $G(\cdot)$ ,  $\min_{x \in [\delta_1/2, \delta_1]} |G(\lambda x) - G(\nu x)| =: \varepsilon > 0.$  Consequently,

$$\int_{[\delta_1/2,\delta_1]} f(x) |G(\lambda x) - G(\nu x)| \, dx \ge \varepsilon \int_{[\delta_1/2,\delta_1]} f(x) \, dx > 0,$$

where the inclusion  $[\delta_1/2, \delta_1] \subset (0, \delta)$  was used to set the last inequality. Therefore, (2.1) can not occur under condition (iii).

According to the previous result, Assumption 2.1 is valid under mild requirements on densities f and q, and it is interesting to observe that if conditions (i)–(iii) in Lemma 2.1 do not hold, then identifiability may fail.

**Example 2.1.** Let the symmetric densities f and g be such that f(x) = 0 for  $x \in [-1,1]$ , whereas q(x) = 0 for |x| > 1. In this case, it is not difficult to see that the general density  $\rho(\cdot; \lambda) \in S(f, g)$  satisfies  $\rho(\cdot; \lambda) = \rho(\cdot; 1)$  for  $\lambda \geq 1$ , and  $\rho(\cdot; \lambda) = \rho(\cdot; -1)$  when  $\lambda \leq -1$ , so that Assumption 2.1 does not hold.

The basic consequence of Assumption 2.1, which plays a central role in the subsequent development, is established in the following lemma. Firstly, recall that  $G(\cdot)$  is increasing and notice that if  $\lambda_2 > \lambda_1$  then

50  
51  

$$\int_{0}^{\infty} f(x) |G(\lambda_{2}x) - G(\lambda_{1}x)| \, dx = \int_{0}^{\infty} f(x) (G(\lambda_{2}x) - G(\lambda_{1}x)) \, dx,$$

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

and

$$\int_{-\infty}^0 f(x) |G(\lambda_2 x) - G(\lambda_1 x)| \, dx \quad = \quad \int_{-\infty}^0 f(x) (G(\lambda_1 x) - G(\lambda_2 x)) \, dx$$

$$= \int_0^\infty f(x)(G(\lambda_2 x) - G(\lambda_1 x)) \, dx,$$

where the second equality comes from  $\int_{\mathbb{R}} f(x) G(\lambda_i x) dx = 1/2$  for i = 1, 2, so that

(2.3) 
$$\int_{R} f(x) |G(\lambda_2 x) - G(\lambda_1 x)| \, dx$$

$$2\int_0^\infty f(x)(G(\lambda_2 x) - G(\lambda_1 x))\,dx$$

$$=2\int_{-\infty}^{0}f(x)(G(\lambda_{1}x)-G(\lambda_{2}x))\,dx,\quad\lambda_{1}<\lambda_{2}.$$

**Lemma 2.2.** Under the identifiability Assumption 2.1, the following assertions (i) and (ii) hold:

(i) For each  $\lambda \in \mathbb{R}$ ,  $\int_0^\infty f(x)G(\lambda x) \, dx > 0$  and  $\int_{-\infty}^0 f(x)G(\lambda x) \, dx > 0$ . (ii) There exists a function  $c : \mathbb{R} \to (0, \infty)$  such that

=

$$P_{\lambda}[X_1 < -c(\lambda)] > 0 \quad and \quad P_{\lambda}[X_1 > c(\lambda)] > 0, \quad \lambda \in \mathbb{R}.$$

Proof. (i) Notice that

$$0 \le \int_0^\infty f(x) G(\lambda_1 x) \, dx \le \int_0^\infty f(x) G(\lambda_2 x) \, dx, \quad \lambda_1 < \lambda_2,$$

since  $G(\cdot)$  is increasing. Now, suppose that  $\int_0^\infty f(x)G(\lambda_2 x) dx = 0$  for some  $\lambda_2 \in \mathbb{R}$ . In this case the above display yields  $\int_0^\infty f(x)G(\lambda_1 x) dx = 0$  for every  $\lambda_1 \leq \lambda_2$ , and then,

$$\int_{R} f(x) |G(\lambda_{2}x) - G(\lambda_{1}x)| \, dx = 0, \quad \lambda_{1} \le \lambda_{2},$$

by the first equality in (2.3), contradicting Assumption 2.1; consequently,  $\int_0^\infty f(x)G(\lambda x) \, dx > 0$  for all  $\lambda \in \mathbb{R}$ , whereas the other part of the conclusion can be obtained along similar lines.

(ii) By the monotone convergence theorem, as  $\varepsilon \searrow 0$ ,  $\int_{\varepsilon}^{\infty} f(x)G(\lambda x) dx \nearrow \int_{0}^{\infty} f(x)G(\lambda x) dx$  and  $\int_{-\infty}^{-\varepsilon} f(x)G(\lambda x) dx \nearrow \int_{-\infty}^{0} f(x)G(\lambda x) dx$  for each  $\lambda$ . There-fore, by part (i), there exists  $c(\lambda) > 0$  such that  $\int_{c(\lambda)}^{\infty} f(x)G(\lambda x) dx > 0$  and  $\int_{-\infty}^{-c(\lambda)} f(x)G(\lambda x) dx > 0$ , and the conclusion follows.

## 

## 3. Existence of Maximum Likelihood Estimators

The objective of this section is to establish the existence of a sequence of maximum likelihood estimators, an idea that is formally stated below. To begin with, given a fixed sample size n > 0, for each possible sample  $\mathbf{x} = (x_1, \ldots, x_n) \in \mathbb{R}^n$  define the average (kernel log-)likelihood function  $L_n(\cdot; \mathbf{x})$  by

$$\lambda \in \mathbb{R}$$

(3.1) 
$$L_n(\lambda; \mathbf{x}) := \frac{1}{n} \sum_{k=1}^n \ell(\lambda; x_k) = \frac{1}{n} \sum_{k=1}^n \log(G(\lambda x_k)), \quad \lambda \in \mathbb{R},$$

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

З

(see (1.1) and (1.4)), where the usual convention  $\log(0) := -\infty$  is enforced. Since  $G(\cdot)$  is continuous and takes values in [0, 1],  $L_n(\cdot; \cdot)$  is continuous function from  $\mathbb{R} \times \mathbb{R}^n$  into  $[-\infty, 0]$ ; moreover, for every  $Q \subset \{1, 2, \ldots, n\}$ 

(3.2) 
$$L_n(\lambda; \mathbf{x}) \le \frac{1}{n} \sum_{i:i \in Q} \log(G(\lambda x_i)) \le 0, \quad \lambda \in \mathbb{R}, \quad \emptyset \ne Q.$$

**Definition 3.1.** Let  $\{\lambda_n : \mathbb{R}^n \to \mathbb{R}\}$  be a sequence of (Borel) measurable functions and set

$$\hat{\lambda}_n := \lambda_n(X_1^n).$$

In this case,  $\{\hat{\lambda}_n\}$  is a sequence of maximum likelihood of estimators of  $\lambda$  if

(3.3) 
$$P_{\lambda_0}\left[\bigcup_{n=1}^{\infty}\bigcap_{k=n}^{\infty} [L_k(\hat{\lambda}_k; X_1^k) \ge L_k(\lambda; X_1^k) \text{ for all } \lambda]\right] = 1, \quad \lambda_0 \in \mathbb{R}.$$

**Remark 3.1.** (i) By continuity,  $L_n(\hat{\lambda}_n; X_1^n) \ge L_n(\lambda; X_1^n)$  occurs for every  $\lambda \in \mathbb{R}$  if and only if it holds for each rational number, so that  $[L_n(\hat{\lambda}_n; X_1^n) \ge L_n(\lambda; X_1^n), \lambda \in \mathbb{R}]$  is an event.

(ii) In words,  $\{\hat{\lambda}_n\}$  is a sequence of maximum likelihood estimators of  $\lambda$  if, with probability 1 and regardless of the true parameter value,  $\hat{\lambda}_n$  maximizes the observed average likelihood function  $L_n(\cdot; X_1^n)$  whenever n is large enough. The event within brackets in (3.3) is the inferior limit of the events  $[L_k(\hat{\lambda}_k; X_1^k) \geq L_k(\lambda; X_1^k), \lambda \in$  $\mathbb{R}]$ , and when (3.3) holds then, as  $k \to \infty$ ,  $P_{\lambda_0}[L_k(\hat{\lambda}_k; X_1^k) \geq L_k(\lambda; X_1^k)$  for all  $\lambda] \to$ 1; see, for instance, Billingsley (1995, Section 4).

The main objective of this section is to prove the following result.

**Theorem 3.1.** Under Assumption 2.1, there exists a sequence of maximum likelihood estimators of  $\lambda$ .

The proof of this theorem has been divided into three simple lemmas involving the following notation: For each  $\mathbf{x} \in \mathbb{R}^n$ , set

(3.4) 
$$m_n(\mathbf{x}) := \sup_{\lambda \in \mathbb{R}} L_n(\lambda; \mathbf{x}),$$

so that (3.1) and (3.2) lead to

(3.5) 
$$0 \ge m_n(\mathbf{x}) \ge L_n(0; \mathbf{x}) = -\log(2), \quad \mathbf{x} \in \mathbb{R}^n,$$

since G(0) = 1/2. Next, define

(3.6) 
$$\mathcal{M}_n(\mathbf{x}) := \{ \nu \in \mathbb{R} \, | \, L_n(\nu; \mathbf{x}) = m_n(\mathbf{x}) \}, \quad \mathbf{x} \in \mathbb{R}^n,$$

which is a closed subset of  $\mathbb{R}$ , by the continuity of  $L_n(\cdot; \mathbf{x})$ . As can be seen from the monotonicity of  $\ell(\cdot; x)$ , the set  $\mathcal{M}_n(\mathbf{x})$  may be empty if the observed sample  $\mathbf{x} \in \mathbb{R}^n$  does not contain observations of different sign. The first step to the proof of Theorem 3.1 is the following lemma, showing that  $\mathcal{M}_n(\mathbf{x})$  is nonempty and compact if  $\mathbf{x} \in \mathbb{R}^n$  contains components with opposite signs. For each integer  $n \geq 2$ , set

(3.7) 
$$\mathcal{S}_n := \{ \mathbf{x} \in \mathbb{R}^n \mid x_i x_j < 0 \text{ for some } i \text{ and } j \text{ with } 1 \le i \ne j \le n \}$$

and observe that  $S_n$  is an open subset of  $\mathbb{R}^n$ .

З

З

**Lemma 3.1.** For each integer  $n \geq 2$  and  $\mathbf{x} \in S_n$ , the set of maximizers  $\mathcal{M}_n(\mathbf{x})$  is nonempty and compact. *Proof.* Given  $\mathbf{x} \in S_n$ , select indexes  $i^*$  and  $j^*$  such that  $x_{i^*} < 0$  and  $x_{j^*} > 0$ , so that  $\lim_{\lambda \to \infty} \log(G(\lambda x_{i^*})) = -\infty = \lim_{\lambda \to -\infty} \log(G(\lambda x_{i^*}))$ . After setting  $Q = \{i^*\}$ and  $Q = \{j^*\}$  in (3.2), these convergences yield that  $\lim_{|\lambda| \to \infty} L_n(\lambda; \mathbf{x}) = -\infty,$ and it follows that the set  $\mathcal{M}_n(\mathbf{x})$ —consisting of the maximizers of the continuous function  $L_n(\cdot; \mathbf{x})$ —is nonempty and compact. Next, for each  $n \geq 2$ , define  $\lambda_n^+ : S_n \to \mathbb{R}$  by so that  $\lambda_n^+(\mathbf{x})$  is the largest element in  $\mathcal{M}_n(\mathbf{x})$ . Observe that  $\lambda_n^+(\mathbf{x}) \in \mathcal{M}_n(\mathbf{x})$  is a well-defined finite number for each  $\mathbf{x} \in S_n$ , by Lemma 3.1. As it is shown below, this function  $\lambda_n^+(\cdot)$  is upper semi-continuous. **Lemma 3.2.** Let the integer  $n \geq 2$  be arbitrary but fixed, and suppose that  $\{\mathbf{x}_k\} \subset$  $S_n$  is such that  $\lim_{k\to\infty} \mathbf{x}_k = \mathbf{y} \in S_n$ . In this context, (i) If  $\nu_k \in \mathcal{M}_n(\mathbf{x}_k)$  for each k, then sequence  $\{\nu_k\}$  is bounded, and (ii) Every limit point of  $\{\nu_k\}$  belongs to  $\mathcal{M}_n(\mathbf{y})$ . Consequently, (iii)  $\lambda_n^+(\cdot)$  is upper semi-continuous. *Proof.* To begin with, write  $\mathbf{x}_k = (x_{k1}, \ldots, x_{kn})$  and notice that, since  $\nu_k \in \mathcal{M}_n(\mathbf{x}_k)$ , for ever  $\lambda \in \mathbb{R}$ ; in particular,  $-\log(2) = L_n(0; \mathbf{x}_k) \le \frac{1}{n} \sum_{i=1}^n \log(G(\nu_k x_{ki})).$ (3.10)(i) Assume now that  $\limsup_{k\to\infty} \nu_k = \infty$  and, taking a subsequence if necessary, without loss of generality suppose that  $\nu_k \to \infty$ . In this context, select an index i<sup>\*</sup> such that  $y_{i^*} < 0$ , which is possible since  $\mathbf{y} \in \mathcal{S}_n$  (see (3.7)) and observe that the convergence  $x_{ki^*} \to y_{i^*} < 0$  leads to  $\nu_k x_{ki^*} \to -\infty$  as  $k \to \infty$ , so that  $\log(G(\nu_k x_{ki^*}) \to -\infty)$ , and then, via (3.2) with  $Q = \{i^*\}$ , this yields  $\lim_{k \to \infty} \frac{1}{n} \sum_{i=1}^{n} \log(G(\nu_k x_{ki})) = -\infty,$ which contradicts (3.10); it follows that  $\limsup_k \nu_k < \infty$ , whereas the inequality  $\liminf_k \nu_k > -\infty$  can be established along similar lines. (ii) Let  $\nu^*$  be an arbitrary limit point of  $\{\nu_k\}$ , and notice that, by part (i),  $\nu^*$ is finite; selecting a subsequence, if necessary, assume that  $\nu_k \to \nu^*$ . In this case, taking the limit as k goes to  $\infty$  in (3.9), it follows that

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

(3.8) 
$$\lambda_n^+(\mathbf{x}) := \max \mathcal{M}_n(\mathbf{x}), \quad \mathbf{x} \in \mathcal{S}_n,.$$

(3.9) 
$$\frac{1}{n}\sum_{i=1}^{n}\log(G(\lambda x_{ki})) = L_n(\lambda; \mathbf{x}_k) \le L_n(\nu_k; \mathbf{x}_k) = \frac{1}{n}\sum_{i=1}^{n}\log(G(\nu_k x_{ki})),$$

*i.e.*,  $\nu^* \in \mathcal{M}_n(\mathbf{y})$ . (iii) Let  $\mathbf{y} \in \mathcal{S}_n$  be arbitrary. If  $\{\mathbf{x}_k\} \subset \mathcal{S}_n$  is such that  $\lim_k \mathbf{x}_k = \mathbf{y}$ , recalling that  $\lambda_n^+(\mathbf{x}_k) \in \mathcal{M}_n(\mathbf{x}_k)$ , part (ii) with  $\nu_k = \lambda_n^+(\mathbf{x}_k)$  yields that  $\limsup_k \lambda_n^+(\mathbf{x}_k) \in$  $\mathcal{M}_n(\mathbf{y})$ , and then  $\limsup_k \lambda_n^+(\mathbf{x}_k) \leq \lambda_n^+(\mathbf{y})$ , by (3.8), so that  $\lambda_n^+(\cdot)$  is upper semi-continuous.

The last step before the proof of Theorem 3.1 is the following consequence of Lemma 2.2(i).

Lemma 3.3. Under Assumption 2.1,

$$\lim_{n \to \infty} P_{\nu} \left[ \bigcap_{k=n}^{\infty} [X_1^k \in \mathcal{S}_k] \right] = 1, \quad \nu \in \mathbb{R};$$

see (1.5) for notation.

*Proof.* Let  $\nu \in \mathbb{R}$  be fixed and observe that for every  $i = 1, 2, \ldots$ ,

$$P_{\nu}[X_i \le 0] = 1 - P_{\nu}[X_i > 0] = 1 - 2\int_0^\infty f(x)G(\nu x) \, dx =: \rho_-(\nu) \in [0, 1),$$

where the inclusion stems from Lemma 2.2(i), so that for each k > 0,  $\rho_{-}(\nu)^{k} =$  $P_{\nu}[X_i \leq 0, 1 \leq i \leq k];$  similarly,  $P_{\nu}[X_i \geq 0, 1 \leq i \leq k] = \rho_{+}(\nu)^k$  for some  $\rho_{+}(\nu) \in [0, 1)$ . Since

$$[X_1^k \in \mathcal{S}_k]^c = [X_i \le 0, \ 1 \le i \le k] \cup [X_i \ge 0, \ 1 \le i \le k],$$

it follows that  $P_{\nu}\left[[X_1^k \in \mathcal{S}_k]^c\right] \leq 2\rho(\nu)^k$ , where  $\rho(\nu) \in [0,1)$  is given by  $\rho(\nu) :=$  $\max\{\rho_{-}(\nu), \rho_{+}(\nu)\}$ . Observing that (1.5) and (3.7) together yield that  $[X_{1}^{n} \in \mathcal{S}_{n}] \subset$  $[X_1^k \in \mathcal{S}_k]$  for  $k \ge n$ , it follows that

$$P_{\nu}\left[\bigcap_{k=n}^{\infty} [X_1^k \in \mathcal{S}_k]\right] \ge P_{\nu}\left[X_1^n \in \mathcal{S}_n\right] \ge 1 - 2\rho(\nu)^n$$

and the conclusion is obtained taking the limit as  $n \to \infty$ .

Notice that the last display and Lemma 3.1 together shows that, with probability increasing to 1 at a geometric rate, the function  $L_n(\cdot; X_1^n)$  has a maximizer when n is large enough.

**PROOF OF THEOREM 3.1.** Select a point  $\lambda^* \in \mathbb{R}$  and for each positive integer k define  $\lambda_k : \mathbb{R}^k \to \mathbb{R}$  as follows:  $\lambda_1(\cdot) \equiv \lambda^*$ , whereas, for  $k \geq 2$ ,  $\lambda_k(\mathbf{x}) := \lambda^*$  if  $\mathbf{x} \in \mathbb{R}^k \setminus \mathcal{S}_k$ , and  $\lambda_k(\mathbf{x}) := \lambda_k^+(\mathbf{x})$  if  $\mathbf{x} \in \mathcal{S}_k$ . Since the  $\mathcal{S}_k$ 's are open sets, Lemma 3.2(iii) implies that each function  $\lambda_k(\cdot)$  is measurable, and then  $\hat{\lambda}_k = \lambda_k(X_1^k)$  is a genuine statistic. Using (3.4)–(3.6) and (3.8), this specification yields  $[X_1^k \in S_k] \subset$  $[L_k(\hat{\lambda}_k; X_1^k) \ge L_k(\lambda; X_1^k), \lambda \in \mathbb{R}]$  for  $k \ge 2$ . Therefore, for each integer  $m \ge 2$ ,

$$\bigcap_{k=m}^{\infty} [X_1^k \in \mathcal{S}_k] \subset \bigcap_{k=m}^{\infty} [L_k(\hat{\lambda}_k; X_1^k) \ge L_k(\lambda; X_1^k), \lambda \in \mathbb{R}]$$

 $\subset \bigcup_{n=1}^{\infty} \bigcap_{k=n}^{\infty} [L_k(\hat{\lambda}_k; X_1^k) \ge L_k(\lambda; X_1^k), \lambda \in \mathbb{R}],$ 

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

a relation that yields that, for every parameter  $\nu$ ,

$$P_{\nu}\left[\bigcup_{n=1}^{\infty}\bigcap_{k=n}^{\infty}\left[L_{k}(\hat{\lambda}_{k};X_{1}^{k})\geq L_{k}(\lambda;X_{1}^{k}),\lambda\in\mathbb{R}\right]\right]\geq P_{\nu}\left[\bigcap_{k=m}^{\infty}\left[X_{1}^{k}\in\mathcal{S}_{k}\right]\right].$$

After taking the limit as  $m \to \infty$ , an application of Lemma 3.3 leads to

$$P_{\nu}\left[\bigcup_{n=1}^{\infty}\bigcap_{k=n}^{\infty}[L_{k}(\hat{\lambda}_{k};X_{1}^{k})\geq L_{k}(\lambda;X_{1}^{k}),\lambda\in\mathbb{R}]\right]=1,\quad\nu\in\mathbb{R},$$

so that, by Definition 3.1,  $\{\hat{\lambda}_n\}$  is a sequence of maximum likelihood estimators.

## 4. Consistency

The objective of this section is to show that a sequence  $\{\hat{\lambda}_n\}$  of maximum likelihood estimators of  $\lambda$  is consistent, *i.e.*, that  $\{\hat{\lambda}_n\}$  converges to the true parameter value with probability 1.

**Theorem 4.1.** Suppose that Assumption 2.1 holds and let  $\{\hat{\lambda}_n\}$  be a sequence of maximum likelihood estimators. In this context,

$$P_{\nu}\left[\lim_{n\to\infty}\hat{\lambda}_n=\nu\right]=1, \quad \nu\in\mathbb{R}.$$

The proof of this result relies on the two lemmas stated below and involves the following notation: Throughout the remainder of the section  $\{\hat{\lambda}_n\}$  is a given sequence of maximum likelihood estimators of  $\lambda$  and the event  $\Omega^*$  is given by

(4.1) 
$$\Omega^* := \bigcup_{n=1}^{\infty} \bigcap_{k=n}^{\infty} [L_k(\hat{\lambda}_k; X_1^k) \ge L_k(\lambda; X_1^k), \lambda \in \mathbb{R}].$$

Also, for each  $\nu \in \mathbb{R}$ ,

$$\Omega_{\nu}^{-} := \left[ \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} I[X_i \le -c(\nu)] = P_{\nu}[X \le -c(\nu)] \right],$$

(4.2)

$$\Omega_{\nu}^{+} := \left[ \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} I[X_i \ge c(\nu))] = P_{\nu}[X \ge c(\nu)] \right],$$

where  $c(\nu) > 0$  is as in Lemma 2.2(ii). Notice that the strong law of large numbers and Definition 3.1 yield

44 (4.3) 
$$P_{\nu}[\Omega^*] = P_{\nu}[\Omega^+] = 1 = P_{\nu}[\Omega^-] \quad \nu \in \mathbb{R}.$$
  
45

The core of the proof of Theorem 4.1 is the following boundedness property.

Lemma 4.1. Under Assumption 2.1, if  $\{\hat{\lambda}_n\}$  is a sequence of maximum likelihood estimators of  $\lambda$ , then

$$P_{\nu}[\limsup_{n \to \infty} |\lambda_n| < \infty] = 1, \quad \nu \in \mathbb{R}.$$

З

*Proof.* It will be shown, by contradiction, that the event  $[\limsup_{n \to \infty} \hat{\lambda}_n = \infty] \cap \Omega^* \cap \Omega_{\nu}^- \quad \text{is empty.}$ (4.4)З To achieve this goal, suppose that the sample trajectory  $X_1, X_2, \ldots$  is such that the above intersection occurs, and observe that along this path assertions (a)–(c) below hold: (a) Since the event  $[\limsup_{n} \hat{\lambda}_n = \infty]$  occurs when  $X_1, X_2, \ldots$  is observed, there exist a (trajectory dependent) subsequence  $\{n_k\}$  such that  $n_k \geq k$  and  $\hat{\lambda}_{n_k} \geq k$ for all positive integers k; (b) Using that  $X_1, X_2, \ldots$  is such that  $\Omega^*$  occurs, it follows that  $\hat{\lambda}_n$  maximizes  $L_n(\cdot; X_1^n)$  for n large enough, so that there exists a positive integer M such that  $L_n(\hat{\lambda}_n; X_1^n) \ge L_n(0; X_1^n) = -\log(2), \quad n \ge M;$ (c) Since the observation of  $X_1, X_2, \ldots$  implies that  $\Omega_{\nu}^-$  occurs,  $\frac{1}{n}\sum_{i=1}^{n}I[X_i \le -c(\nu)] \to P_{\nu}[X_1 \le -c(\nu)] \quad \text{as } n \to \infty.$ Notice now that for each integer  $M_1 > M$  and  $k > M_1$ , (a) and (b) together yield  $-\log(2) \le L_{n_k}(\hat{\lambda}_{n_k}; X_1^{n_k}) = \frac{1}{n_k} \sum_{i=1}^{n_k} \log(G(\hat{\lambda}_{n_k}X_i))$  $\leq \frac{1}{n_k} \sum_{i=1}^{n_k} \log(G(\hat{\lambda}_{n_k} X_i)) I[X_i \leq -c(\nu)]$  $\leq \frac{1}{n_k} \sum_{i=1}^{n_k} \log[G(-M_1 c(\nu))] I[X_i \leq -c(\nu)]$ where (3.2) with  $Q = \{i : i \leq n_k, X_i \leq -c(\nu)\}$  was used to set the second inequality, and the third one follows from the monotonicity of  $\log(G(\cdot))$ , since  $\lambda_{n_k} >$  $k > M_1$ . From this point, letting k go to  $\infty$ , (c) leads to

 $-\log(2) \le \log(G(-M_1c(\nu)))P_{\nu}[X_1 \le -c(\nu)];$ 

and, recalling that  $\lim_{x\to -\infty} \log(G(x)) = -\infty$  and that  $P_{\nu}[X_1 \leq -c(\nu)]$  and  $c(\nu)$  are both positive (by Lemma 2.2(ii)), taking the limit as  $M_1 \to \infty$ , it follows that  $-\log(2) \leq -\infty$ , which is a contradiction, establishing (4.4). Therefore,  $[\limsup_{n\to\infty} \hat{\lambda}_n = \infty] \subset (\Omega^*)^c \cup (\Omega_{\nu}^-)^c$ , inclusion that yields

$$P_{\nu}[\limsup_{n \to \infty} \hat{\lambda}_n = \infty] = 0, \quad \nu \in \mathbb{R},$$
42
43

by (4.3). Similarly, it can be established that  $[\liminf_{n\to\infty} \hat{\lambda}_n = -\infty] \subset (\Omega^*)^c \cup$  $(\Omega^+_{\nu})^c$ , so that  $P_{\nu}[\liminf_{n\to\infty} \lambda_n = -\infty] = 0$  for all  $\nu \in \mathbb{R}$ ; the conclusion follows combining this fact with the above display.

To continue, observe that since  $\log(G(\cdot)) \leq 0$ ,  $E_{\nu}[\log(G(\lambda X_1))]$  is always a well-defined non positive number, where the expectation may assume the value  $-\infty$ . Also, since  $x \mapsto x \log(x)$  is bounded on  $(0,1), E_{\nu}[\log(G(\nu x))]$  $\int_{\mathbb{R}} 2\log(G(\nu x))G(\nu x)f(x)\,dx$  is finite. 

For each  $\nu \in \mathbb{R}$ , define

**Lemma 4.2.** Suppose that Assumption 2.1 holds and let  $\nu \in \mathbb{R}$  be arbitrary but fixed. (i) [Kullback's inequality.] For each  $\lambda \in \mathbb{R} \setminus \{\nu\}$ ,  $E_{\nu}[\log(G(\lambda X_1))] < E_{\nu}[\log(G(\nu X_1))].$ (ii) Assume that  $\{r_k\}$  and  $\{s_k\}$  are two real sequences such that, for some  $\nu^* \in \mathbb{R}$ ,  $r_k \searrow \nu^*$  and  $s_k \nearrow \nu^*$  as  $k \to \infty$ . and suppose that, for every  $k = 1, 2, 3, \ldots$ , the following inequality holds:  $E_{\nu}[\log(G(\nu X))]$ (4.5) $\leq E_{\nu}[\log(G(r_k X))I[X \geq 0]] + E_{\nu}[\log(G(s_k X))I[X < 0]].$ In this case,  $\nu = \nu^*$ . *Proof.* (i) If  $\lambda \neq \nu$ , Assumption 2.1 and the strict concavity of the logarithmic function yield, via Jensen's inequality, that  $\int_{\mathbb{R}} \log\left(\frac{\rho(x;\lambda)}{\rho(x;\nu)}\right) \rho(x;\nu) \, dx \quad < \quad \log\left(\int_{\mathbb{R}} \frac{\rho(x;\lambda)}{\rho(x;\nu)} \rho(x;\nu) \, dx\right)$  $= \log\left(\int_{\mathbb{R}} \rho(x;\lambda) \, dx\right) = 0.$ Observing that  $\rho(x;\lambda)/\rho(x;\nu) = G(\lambda x)/G(\nu x)$  when  $\rho(x;\nu)$  is positive, the above inequality can be written as  $E_{\nu}[\log(G(\lambda X_1)) - \log(G(\nu X_1))] < 0$ , which yields the desired conclusion since, as already noted,  $E_{\nu}[\log(G(\nu X_1))]$  is finite. (ii) Notice that  $r_k \searrow \nu^*$  leads to  $r_k X \searrow \nu^* X$  on  $[X \ge 0]$ , whereas  $s_k \nearrow \nu^*$ implies that  $s_k X \searrow \nu^* X$  on [X < 0]. Therefore, since  $-\log(G(\cdot))$  is decreasing and nonnegative,  $0 < -\log(G(r_k X))I[X > 0] \nearrow -\log(G(\nu^* X))I[X > 0]$ and  $0 \le -\log(G(s_k X))I[X < 0] \nearrow -\log(G(\nu^* X))I[X < 0].$ Now, an application of the monotone convergence theorem yields  $E_{\nu}[\log(G(r_k X))I[X \ge 0]] \searrow E_{\nu}[\log(G(\nu^* X))I[X \ge 0]]$ and  $E_{\nu}[\log(G(s_k X))I[X < 0]] \searrow E_{\nu}[\log(G(\nu^* X))I[X < 0]],$ convergences that, after taking the limit as k goes to  $\infty$  in (4.5), lead to  $E_{\nu}[\log(G(\nu X))] < E_{\nu}[\log(G(\nu^* X))],$ and then  $\nu = \nu^*$ , by part (i). After these preliminaries, the proof of Theorem 4.1 is presented below. The argument uses the following notation, where Q stands for the set of rational numbers:

З

(

4.6) 
$$\Omega_{\nu}^{0} := \left[ \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \log(G(\nu X_{i})) = E_{\nu}[\log(G(\nu X_{1}))] \right];$$

(4.7) 
$$\Omega^{1}_{\nu} := \left[ \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \log(G(\lambda X_{i})) I[X_{i} \ge 0] \right]$$
$$= E_{\nu}[\log(G(\lambda X_{1})) I[X_{1} \ge 0]], \quad \lambda \in \mathcal{Q} ,$$

and

(4.8) 
$$\Omega_{\nu}^{2} := \left[ \lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \log(G(\lambda X_{i})) I[X_{i} < 0] \right]$$
$$= E_{\nu} [\log(G(\lambda X_{1})) I[X_{1} < 0]], \quad \lambda \in \mathcal{Q} \right].$$

Since  $\mathcal{Q}$  is denumerable, the strong law of large numbers yields  $P_{\nu}[\Omega_{\nu}^{i}] = 1$  for i = 0, 1, 2, and then, setting

(4.9) 
$$\Omega_{\nu} = \Omega_{\nu}^{0} \cap \Omega_{\nu}^{1} \cap \Omega_{\nu}^{2},$$

it follows that

$$(4.10) P_{\nu}[\Omega_{\nu}] = 1.$$

PROOF OF THEOREM 4.1. Given  $\nu \in \mathbb{R}$ , it is sufficient to show that

(4.11) 
$$\left[\limsup_{n} |\hat{\lambda}_{n}| < \infty\right] \cap \Omega^{*} \cap \Omega_{\nu} \subset \left[\lim_{n} \hat{\lambda}_{n} = \nu\right],$$

where  $\Omega^*$  and  $\Omega_{\nu}$  are specified in (4.1) and (4.9), respectively. Indeed, if this inclusion is valid, then (4.3), (4.10) and Lemma 4.1 together imply that  $P_{\nu}[\lim_{n} \hat{\lambda}_{n}]$  $\nu$ ] = 1. To establish (4.11) let  $X_1, X_2, X_3, \ldots$  be a fixed trajectory such that  $\left[\limsup_n |\hat{\lambda}_n| < \infty\right] \cap \Omega^* \cap \Omega_\nu \text{ occurs, select an arbitrary limit point } \nu^* \text{ of the }$ associated sequence  $\{\hat{\lambda}_n\}$ , and observe the following facts (a)–(c): (a)  $\nu^*$  is finite, since  $\limsup_n |\lambda_n| < \infty$  holds when  $X_1, X_2, X_3, \ldots$  is observed. Let

r and s be arbitrary rational numbers satisfying

44 (4.12) 
$$s < \nu^* < r$$

and select a sequence  $\{n_k\}$  of positive integers such that

48 (4.13) 
$$n_k > k, \quad s < \hat{\lambda}_{n_k} < r, \quad k = 1, 2, \dots$$
 48  
49 49

(b) Since the path 
$$X_1, X_2, X_3, \ldots$$
 is such that  $\Omega^*$  occurs,  $\hat{\lambda}_n$  is a maximizer of  $L_n(\cdot; X_1^n)$  when  $n$  is large enough, say  $n > M$ ; see (4.1). In particular,  $L_n(\nu; X_1^n) \le 51$ 

 $L_n(\hat{\lambda}_n; X_1^n)$  when n > M, and then, replacing n by  $n_k$  in this inequality and using the monotonicity of  $\log(G(\cdot))$  as well as (4.13), it follows that  $\frac{1}{n_k} \sum_{i=1}^{n_k} \log(G(\nu X_i)) = L_{n_k}(\nu; X_1^{n_k})$  $\leq L_{n_k}(\hat{\lambda}_{n_k}; X_1^{n_k})$  $= \frac{1}{n_k} \sum_{k=1}^{n_k} \log(G(\hat{\lambda}_{n_k} X_i))$  $\leq \frac{1}{n_k} \sum_{i=1}^{n_k} \log(G(rX_i)) I[X_i \ge 0]$ (4.14) $+\frac{1}{n_k}\sum_{i=1}^{n_k}\log(G(sX_i))I[X_i < 0], \quad k > M.$ (c) Since the trajectory  $X_1, X_2, \ldots$  is such that  $\Omega_{\nu}$  occurs, a glance to (4.6)–(4.9) immediately yields that, as  $k \to \infty$ ,  $\frac{1}{n_k} \sum_{i=1}^{n_k} \log(G(\nu X_i)) \quad \to \quad E_{\nu}[\log(G(\nu X))],$  $\frac{1}{n_k} \sum_{i=1}^{n_k} \log(G(rX_i)) I[X_i \ge 0] \quad \to \quad E_{\nu}[\log(G(rX)I[X \ge 0]], \text{ and}$  $\frac{1}{n_k} \sum_{i=1}^{n_k} \log(G(sX_i)) I[X_i < 0] \quad \rightarrow \quad E_{\nu}[\log(G(sX)I[X < 0]].$ After taking the limit as  $k \to \infty$  in (4.14), these convergences yield that  $E_{\nu}[\log(G(\nu X))] \le E_{\nu}[\log(G(rX)I[X \ge 0]] + E_{\nu}[\log(G(sX)I[X < 0]]],$ and then, since r and s are arbitrary rational numbers satisfying (4.12), from Lemma 4.2(ii) it follows that  $\nu^* = \nu$ . In short, it has been proved that along an arbitrary path  $X_1, X_2, \ldots$  for which the intersection  $\left| \limsup_n |\hat{\lambda}_n| < \infty \right| \cap \Omega^* \cap \Omega_{\nu}$  occurs, the corresponding sequence  $\{\hat{\lambda}_n\}$  has  $\nu$  as its unique limit point, so that  $\lambda_n \to \nu$  as  $n \to \infty$ . This establishes (4.11) and, as already noted, completes the proof. 5. Asymptotic Distribution

The remainder of the paper concerns the asymptotic behavior of a consistent sequence  $\{\lambda_n\}$  of maximum likelihood estimators of  $\lambda$ , whose existence is guaranteed by Assumption 2.1. As already mentioned, the large sample properties of  $\{\lambda_n\}$  will be studied under the null hypothesis  $\mathcal{H}_0: \lambda = 0$ , and the analysis below requires two properties on the densities g and f generating the family S(f,g), namely, (i) smoothness of density g outside of  $\{0\}$  and a 'good' behavior of its derivatives around  $\lambda = 0$ , and (ii) a moment-dominance condition involving both densities f and g. After a formal presentation of these assumptions, the main result is stated at the end of the section, and the corresponding proof is given after establishing the necessary technical preliminaries.

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

Assumption 5.1. For some nonnegative integer r—hereafter referred to as the critical order— the following conditions hold: (i) The symmetric density q(x) is continuous on  $\mathbb{R}$  and has derivatives up to order 2r+2 on the interval  $(0,\infty)$ ; (ii)  $D^k g(0+) = \lim_{x \to 0} D^k g(x)$  exists and is finite for k = 0, 1, 2, ..., 2r + 1, and (iii)  $D^r q(0+) \neq 0$ , whereas  $D^s q(0+) = 0$  for 0 < s < r. **Remark 5.1.** (i) Under this Assumption,  $q(\cdot)$  statisfies that  $D^rg(0+) = \lim_{x \searrow 0} r!g(x)/x^r \ge 0$ , by continuity of g if r = 0, or by L'Hopital's rule, if r > 0, so that  $D^r g(0+) > 0$ , since  $D^r g(0+)$  is no null. It follows that a density g satisfying Assumption 5.1 can be expressed as  $q(x) = |x|^r h(|x|)$  where r is a nonnegative integer (the critical order),  $h: [0,\infty) \to [0,\infty)$  is continuous with h(0) > 0, and has derivatives of order up to 2r + 2 on  $(0, \infty)$ , which are 'well-behaved' near zero so that the required lateral limits of the derivatives of q exist at x = 0. Thus, besides the smoothness requirement on the whole interval  $(0, \infty)$ , the core of Assumption 5.1 essentially concerns the local behavior of q around the origin. (ii) Under Assumption 5.1, density  $g(\cdot)$  is continuous, so that  $G(\cdot)$  has continuous derivative, and then  $\partial_{\lambda}\ell(\lambda; x) = \partial_{\lambda}\log(G(\lambda x))$  exists if  $\ell(\lambda; x)$  is finite. Suppose now that  $\hat{\lambda}_n$  maximizes  $L_n(\cdot; X_1^n)$ . In this case  $-\log(2) = L_n(0; X_1^n) \leq L_n(\hat{\lambda}_n; X_1^n)$  im-plies that  $\log(G(\hat{\lambda}_n X_i))$  is finite for every i = 1, 2, ..., n, by (3.1), so that  $L_n(\cdot; X_1^n)$ is differentiable at  $\hat{\lambda}_n$ , and then the likelihood equation holds:  $\partial_{\lambda} L_n(\hat{\lambda}_n; X_1^n) = 0$ . By symmetry, Assumption 5.1 yields that  $g(\cdot)$  also has derivatives up to order 2r+2 on the interval  $(-\infty, 0)$ ; indeed, if  $k \leq 2r+2$  then  $D^k g(x) = (-1)^k D^k g(-x)$ for  $x \neq 0$ , so that  $D^{k}g(0-) = \lim_{x \neq 0} D^{k}g(x) = (-1)^{k}D^{k}g(0+), \quad k = 0, 1, 2, \dots, 2r+1,$ (5.1)and the nullity of  $D^k g(0+)$  for  $0 \le k < r$  implies that g has null (bilateral) derivative at x = 0 of any order less that r. On the other hand, under Assumption 5.1 the cumulative distribution function G in (1.1) has derivatives up to order 2r + 3 on  $\mathbb{R} \setminus \{0\}$ , and using the relation  $D^k G(x) = D^{k-1}g(x)$  for  $x \neq 0$  and k > 0 it follows that  $D^k G(0) = 0, \quad 1 \le k \le r+1$ (5.2)and (5.3)  $D^k G(0+) = D^{k-1} g(0+), \quad D^k G(0-) = D^{k-1} g(0-), \quad r+1 \le k \le 2r+2.$ Next, define (5.4) $H(x) := \log(G(x)), \quad x \in \mathbb{R},$ and observe the equalities  $\ell(\lambda; x) = H(\lambda x), \quad x \in \mathbb{R},$ (5.5) $\partial_{\lambda}^{k}\ell(\lambda;x) = D^{k}H(\lambda x)x^{k}, \quad x \neq 0, \quad 1 < k < 2r + 3$ see (1.4). It follows that the lateral limits of  $\partial_{\lambda}^{k}\ell(\cdot;x)$  at zero are given by  $\partial_{\lambda}^{k}\ell(0+;x) = \left\{ \begin{array}{ll} D^{k}H(0+)x^{k}, & \text{if } x>0; \\ D^{k}H(0-)x^{k}, & \text{if } x<0. \end{array} \right.$ (5.6)

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

and

(5.7) 
$$\partial_{\lambda}^{k}\ell(0-;x) = \begin{cases} D^{k}H(0-)x^{k}, & \text{if } x > 0; \\ D^{k}H(0+)x^{k}, & \text{if } x < 0. \end{cases}$$

The analysis below uses (lateral) Taylor expansions of order 2r + 2 for  $L_n(\cdot; X_1^n)$ around zero, and it is necessary to have an integrable bound for the residual as well as finite second moments for the coefficients. For this reason, the following conditions will be enforced.

**Assumption 5.2.** Conditions (i) and (ii) below hold, where r is the the critical order in Assumption 5.1:

(i) 
$$E_0[X_1|^{4r+2}] = \int_{\mathbb{R}} x^{4r+2} f(x) \, dx < \infty$$

(ii) There exists a function  $W : \mathbb{R} \to [0, \infty)$  and  $\delta > 0$  such that

(5.8) 
$$|\partial_{\lambda}^{2r+3}\ell(\lambda;\cdot)| \le W(\cdot), \quad 0 < |\lambda| \le \delta, \quad and \quad \int_{\mathbb{R}} W(x)f(x) \, dx < \infty.$$

**Remark 5.2.** The moment requirement in Assumption 5.2(i) concerns only density f, whereas the dominance condition in the second part involves a relation between g and f. Using (5.4), it can be shown by induction that, for  $x \neq 0$ ,  $D^kH(x)$  is a polynomial in  $D^{s}g(x)/G(x)$ , s = 0, 1, 2, ..., k-1, and then, setting

$$M_r := \max\{|D^k g(x)|/G(x): 0 \le k \le 2r+2, \ x \ne 0\},\$$

via (5.5) it follows that there exists a constant B such that  $|\partial_{\lambda}^{2r+3}\ell(\lambda,x)| \leq |\partial_{\lambda}^{2r+3}\ell(\lambda,x)| < |\partial_{\lambda}^{2r+3}\ell(\lambda,x)| < |\partial_{\lambda}^{2r+3}\ell(\lambda,x)| < |\partial_{\lambda}^{2r+3}\ell(\lambda,x)| < |\partial_{$  $BM_r|x|^{2r+3}$  so that, if  $M_r < \infty$ , then Assumption 5.2 holds entirely if the moment condition in part (i) is valid. It is not difficult to see that  $M_r$  is finite when there exists  $x_0 > 0$  for which (a) or (b) below occur:

(a) g(x) is a rational function for  $x > x_0$ , as it is the case if  $g(\cdot)$  is a multiple of a *t*-density for x large enough;

(b)  $g(x) = p(x)e^{-\beta x}$  on  $(x_0, \infty)$ , where  $\beta$  is a positive constant and p(x) is a polynomial or, more generally, a linear combination of terms of the form  $x^{s}$ ; this occurs, for instance, when  $q(\cdot)$  is proportional to a mixture of gamma densities on  $(x_0,\infty).$ 

To state the result on the large sample distribution of maximum likelihood estimators, set

(5.9) 
$$V_{r+1} := 4 \left(\frac{D^r g(0+)}{(r+1)!}\right)^2 E_0\left[X_1^{2r+2}\right] > 0$$

which, as it will be shown later, is the variance of  $\partial_{\lambda}^{r+1}\ell(0+, X_i)/(r+1)!$ ; for the strict inequality, see Remark 5.1(i).

**Theorem 5.1.** Let  $\{\hat{\lambda}_n\}$  be a consistent sequence of maximum likelihood estimators of  $\lambda$ , and suppose that Assumptions 5.1 and 5.2 hold. In this context, under the hypothesis  $\mathcal{H}_0: \lambda = 0$ , the following convergences (i) and (ii) occur as  $n \to \infty$ , where r is the critical order in Assumption 5.1, and Z is a random variable with standard normal distribution: 

and

$$\begin{array}{c} {}^{50}\\ {}_{51} \qquad (ii) \qquad \qquad 2n[L_n(\hat{\lambda}_n;X_1^n) - L_n(0;X_1^n)] \stackrel{\mathrm{d}}{\longrightarrow} Z^2. \end{array}$$

**Remark 5.3.** (i) Suppose that the critical index r is null. In this case Theorem 5.1 (i) yields that  $(nV_1)^{1/2} \hat{\lambda}_n \xrightarrow{d} |Z| \operatorname{sign}(Z) \stackrel{d}{=} Z$ . This conclusion coincides with that obtained from the general classical results presented, for instance, in Lehmann and Casella (1998, Section 6.3), or Shao (1999, Section 4.4), where derivatives up to order 2 are required for  $g(\cdot)$  around zero; under Assumption 5.1, only the lateral limits of Dg and  $D^2g$  exist at zero. (ii) Suppose that the critical order r is positive and that  $g(\cdot)$  has (bilateral) derivative 

of order r at zero, so that  $D^r g(0+) = D^r g(0-)$ . Since  $D^r g(0+) \neq 0$  it follows from (5.1) that r is an even integer, and  $g(\cdot)$  has derivatives up to order r on the real line. Thus, setting s = r + 1, s is odd,  $\ell(\cdot; x)$  has derivatives up to order s on  $\mathbb{R}$ , and  $\partial_{\lambda}^{k}\ell(0; \cdot) = 0$  for  $1 \leq k < s$ , whereas  $\partial_{\lambda}^{s}\ell(0, x) \neq 0$  for  $x \neq 0$ . If, among other conditions,  $q(\cdot)$  has derivatives up to order 2s at zero, an application of Theorem 1 in Rotnitzky et al. (2000) yields the conclusions in Theorem 5.1; notice, however, that Assumption 5.1 only ensures the existence of the lateral limits  $D^k g(0\pm)$  for  $s < k \leq 2s$ , so that Theorem 5.1 extends Theorem 1 in Rotnitzky et al. (2000) to the framework of this work.

The rather technical proof of Theorem 5.1, requiring some explicit computations for the lateral limits  $\partial_{\lambda}^{k}\ell(0\pm;x)$  in terms of density  $g(\cdot)$ , will be given after the preliminaries established in the following section.

### 6. Technical Preliminaries

This section is dedicated to establish the basic tools that will be used to prove Theorem 5.1, namely, lateral Taylor expansions around the origin for the average likelihood  $L_n(\cdot; X_1^n)$  and its first derivative; via the absolute value function, such expansions are stated below as single equations. The following notation will be used:

(6.1) 
$$\Delta(x) := \partial_{\lambda}^{r+1} \ell(0-;x) - \partial_{\lambda}^{r+1} \ell(0+;x), \quad x \in \mathbb{R}.$$

 $\Delta$ 

**Theorem 6.1.** Suppose that Assumptions 5.1 and 5.2 hold. In this case, the assertions (i)—(iii) below occur, where r is the critical order in Assumption 5.1,  $\delta > 0$ and  $W(\cdot)$  are as in Assumption 5.2, and

$$_{n} := \frac{1}{n} \sum_{i=1}^{n} \Delta(X_{i});$$

see (6.1).

(i) For each positive integer n and  $\alpha \in (-\delta, \delta)$ ,

$$L_n(\alpha; X_1^n) - L_n(0; X_1^n)$$
45

$$= |\alpha|^{r} \alpha \begin{cases} \sum_{k=r+1}^{2r+1} \frac{\partial_{\lambda}^{k} L_{n}(0+;X_{1}^{n})}{k!} |\alpha|^{k-r-1} & 47 \\ 48 \end{cases}$$

<sup>49</sup>  
<sup>50</sup> (6.2) 
$$+ \left( \frac{\partial_{\lambda}^{2r+2} L_n(0+;X_1^n) + \Delta_n I_{(-\infty,0)}(\alpha)}{(2r+2)!} + \frac{W_n^*(\alpha)}{(2r+3)!} \alpha \right) |\alpha|^r \alpha \right\},$$
<sup>49</sup>  
<sup>50</sup> 50

$$\partial_\lambda L_n(lpha,X_1^n)$$
 3

$$= |\alpha|^r \left\{ \sum^{2r} \frac{\partial^{k+1}_{\lambda} L_n(0+;X_1^n)}{k!} |\alpha|^{k-r} \right.$$

$$\left( \partial_{\lambda}^{2r+2} L_n(0+;X_1^n) + \Delta_n I_{(-\infty,0)}(\alpha) - \tilde{W}_n(\alpha) \right) = 0$$

(6.3) 
$$+\left(\frac{\partial_{\lambda}-L_{n}(0+,X_{1})+\Delta_{n}T(-\infty,0)(\alpha)}{(2r+1)!}+\frac{W_{n}(\alpha)}{(2r+2)!}\alpha\right)|\alpha|^{r}\alpha\bigg\},$$

where the random variables  $W_n^*(\alpha)$  and  $\tilde{W}_n(\alpha)$  satisfy

(6.4) 
$$|W_n^*(\alpha)|, \ |\tilde{W}_n(\alpha)| \le W_n := \frac{1}{n} \sum_{i=1}^n W(X_i).$$

(ii) Under  $\mathcal{H}_0$ :  $\lambda = 0$  for each k = r + 1, ..., 2r + 1, the following convergences hold as  $n \to \infty$ :

(6.5) 
$$\sqrt{n} \ \partial_{\lambda}^{k} L_{n}(0+;X_{1}^{n}) \xrightarrow{\mathrm{d}} \mathcal{N}(0,v_{k}), \quad where \ v_{k} = E_{0} \left[ \left( \partial_{\lambda}^{k} \ell(0+;X_{1}) \right)^{2} \right],$$

whereas

(6.6) 
$$\Delta_n \to 0 \quad and \quad 2 \frac{\partial_{\lambda}^{2r+2} L_n(0+;X_1^n)}{(2r+2)!} \to -V_{r+1} \quad P_0\text{-}a.s.;$$

see (5.9).

The proof of this theorem relies on explicit formulas for  $\partial_{\lambda}^{k}\ell(0\pm;x)$  in terms of density  $g(\cdot)$  and, in this direction, the following lemma concerning the lateral limits at zero of the derivatives of function  $H(\cdot)$  in (5.4) will be useful.

**Lemma 6.1.** Suppose that Assumption 5.1 holds. In this case, the lateral limits at 0 of the derivatives of function  $H(\cdot)$  in (5.4) satisfy the following relations (i)–(iii): (i)  $D^k H(0+) = D^k H(0-) = 0$ ,  $1 \le k < r + 1$ ;

(*ii*) If  $r + 1 \le k < 2r + 2$ , then

$$D^k H(0+) = 2D^{k-1}g(0+)$$
 and  $D^k H(0-) = 2D^{k-1}g(0-).$ 

(*iii*) 
$$D^{2r+2}H(0+) = 2D^{2r+1}g(0+) - \frac{1}{2}\binom{2r+2}{r+1}\left(D^{r+1}H(0+)\right)^2,$$

and

$$D^{2r+2}H(0-) = 2D^{2r+1}g(0-) - \frac{1}{2}\binom{2r+2}{r+1} \left(D^{r+1}H(0-)\right)^2.$$

*Proof.* Recalling that the distribution function G(x) is continuous and has derivatives up to order 2r + 3 on  $\mathbb{R} \setminus \{0\}$ ; from (5.4) it follows that G(x)DH(x) = DG(x) and, via Leibinitz' formula,

$$G(x)D^{k}H(x) + \sum_{i=1}^{k-1} \binom{k-1}{i} D^{i}G(x)D^{k-i}H(x) = D^{k}G(x)$$

for  $x \neq 0$  and  $2 \leq k \leq 2r + 3$ ; since  $G(\cdot)$  is continuous and G(0) = 1/2, taking lateral limit as x approaches to zero these equalities lead to

$$DH(0\pm) = 2DG(0)$$

З

and, for 
$$2 \le k \le 2r+2$$
,

$$D^{k}H(0\pm) + 2\sum_{i=1}^{k-1} \binom{k-1}{i} D^{i}G(0\pm)D^{k-i}H(0\pm) = 2D^{k}G(0\pm).$$

Since  $D^k G(0\pm) = 0$  when  $1 \le k \le r$ , by (5.2), these relations yield  $D^k H(0\pm) = 0$  for  $1 \le k \le r$ , establishing part(i), as well as

(6.7) 
$$D^{r+1}H(0\pm) = 2D^{r+1}G(0\pm),$$

and

for 
$$r + 1 < k \le 2r + 2$$
,

(6.8) 
$$D^{k}H(0\pm) + 2\sum_{i=r+1}^{k-1} \binom{k-1}{i} D^{i}G(0\pm)D^{k-i}H(0\pm) = 2D^{k}G(0\pm).$$

To prove part (ii), select an integer k such that r + 1 < k < 2r + 2. In this case, if  $k > i \ge r + 1$  then  $1 \le k - i < r + 1$ , and then  $D^{k-i}H(0\pm) = 0$ , by part (i), so that the summation in the above display is null. Therefore,  $D^kH(0\pm) = 2D^kG(0\pm)$ , and combining this with (6.7) it follows that

$$D^k H(0\pm) = 2D^k G(0\pm), \quad r+1 \le k < 2r+2,$$

equalities that yield part (ii) via (5.3). To conclude, observe that if k = 2r + 2 then  $2r + 1 \ge i > r + 1$  implies that  $1 \le k - i < r + 1$ , and in this case  $D^{k-i}H(0\pm) = 0$ , by part (i), so that the terms in the summation in (6.8) with k = 2r + 2 are null when i > r + 1. Consequently,

$$D^{2r+2}H(0\pm) + 2\binom{2r+1}{r+1}D^{r+1}G(0\pm)D^{r+1}H(0\pm) = 2D^{2r+2}G(0\pm);$$

since

$$2D^{r+1}G(0\pm) = 2D^r g(0\pm) = D^{r+1}H(0\pm)$$

and  $D^{2r+2}G(0\pm) = D^{2r+1}g(0\pm)$ , by (5.3) and part (ii), respectively, the conclusion follows observing that  $\binom{2r+1}{r+1} = 2^{-1}\binom{2r+2}{r+1}$ .

The expressions in the previous lemma are used below to determine the lateral limits of  $\partial_{\lambda}^{k} \ell(\cdot; x)$  at zero in terms of density  $g(\cdot)$ .

<sup>44</sup> Lemma 6.2. Under Assumption 5.1, assertions (i)–(v) below hold:

<sup>45</sup> (i)  $\partial_{\lambda}^{k}\ell(0+;\cdot) = 0 = \partial_{\lambda}^{k}\ell(0-;\cdot)$  for  $1 \le k \le r$ .

46 (*ii*) For each  $x \in \mathbb{R}$  and  $r+1 \le k < 2r+2$ , 

$$\partial_{\lambda}^{k}\ell(0+;x) = 2D^{k-1}q(0+)|x|^{k-1}x,$$

$$\partial_{\lambda}^{k} \ell(0-;x) = 2D^{k-1}g(0-)|x|^{k-1}x.$$
49

(*iii*) 
$$\partial_{\lambda}^{k} \ell(0-;x) = (-1)^{k-1} \partial_{\lambda}^{k} \ell(0+;x)$$
 for  $r+1 \le k < 2r+2$  and  $x \in \mathbb{R}$ .

(iv) For each  $x \in \mathbb{R}$  $\partial_{\lambda}^{2r+2}\ell(0+;x) = 2D^{2r+1}g(0+)|x|^{2r+1}x$ (6.9) $-\frac{1}{2}\binom{2r+2}{r+1}\left(\partial_{\lambda}^{r+1}\ell(0+;x)\right)^{2},$ and  $\partial_{\lambda}^{2r+2}\ell(0-;x) = 2D^{2r+1}g(0-)|x|^{2r+1}x$ (6.10) $-\frac{1}{2}\binom{2r+2}{r+1}\left(\partial_{\lambda}^{r+1}\ell(0-;x)\right)^{2}.$ Consequently, (v) The difference between  $\partial_{\lambda}^{2r+2}\ell(0-;x)$  and  $\partial_{\lambda}^{2r+2}\ell(0+;x)$  is given by  $\Delta(x) = -4D^{2r+1}g(0+) |x|^{2r+1}x, \quad x \in \mathbb{R};$ (6.11)see (**6**.**1**). *Proof.* From Lemma 6.1(i), part (i) follows via (5.6) and (5.7), whereas these latter equalities and Lemma 6.1(ii) together yield that, for  $r+1 \le k \le 2r+1$ , (a) and (b) below hold: (a) For  $x \ge 0$ ,  $\partial_{\lambda}^{k}\ell(0+;x) = 2D^{k-1}g(0+)x^{k}$  and  $\partial_{\lambda}^{k}\ell(0-;x) = 2D^{k-1}g(0-)x^{k};$ (b) If x < 0,

 $\begin{array}{rcl} \partial_{\lambda}^{k}\ell(0+;x) &=& 2D^{k-1}g(0-)x^{k} \\ &=& 2D^{k-1}g(0+)(-1)^{k-1}x^{k} = 2D^{k-1}g(0+)|x|^{k-1}x, \end{array}$ 

$$\begin{aligned} \partial_{\lambda}^{\kappa}\ell(0-;x) &= 2D^{\kappa-1}g(0+)x^{\kappa} \\ &= 2D^{k-1}g(0-)(-1)^{k-1}x^{k} = 2D^{k-1}g(0-)|x|^{k-1}x, \end{aligned}$$

where (5.1) was used to set the second equalities. These facts (a) and (b) together lead to part (ii), which implies part (iii) via (5.1). To establish part (iv), notice that Lemma 6.1(iii) and (5.6) together imply that For  $x \ge 0$ 

$$\partial_{\lambda}^{2r+2}\ell(0+;x) = 2D^{2r+1}g(0+)x^{2r+2} - \frac{1}{2}\binom{2r+2}{r+1}\left(D^{r+1}H(0+)x^{r+1}\right)^2$$

$$= 2D^{2r+1}g(0+)x^{2r+2} - \frac{1}{2}\binom{2r+2}{r+1} \left(\partial_{\lambda}^{r+1}\ell(0+;x)\right)^2,$$

showing that (6.9) holds for x > 0, whereas combining Lemma 6.1(iii) with relations (5.6) and (5.1) it follows that if x < 0, then

$$\partial_{\lambda}^{2r+2}\ell(0+;$$

x)

$$= 2D^{2r+1}g(0-)x^{2r+2} - \frac{1}{2}\binom{2r+2}{r+1}\left(D^{r+1}H(0-)x^{r+1}\right)^2$$

$$=2D^{2r+1}g(0+)(-1)^{2r+1}x^{2r+1}x-\frac{1}{2}\binom{2r+2}{r+1}\left(\partial_{\lambda}^{r+1}\ell(0+;x)\right)^{2}$$

#### imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

and then (6.9) also holds for x < 0. Equality (6.10) can be established along similar lines and, finally, observing that  $\partial_{\lambda}^{r+1}\ell(0+;x)$  and  $\partial_{\lambda}^{r+1}\ell(0-;x)$  have the same absolute value, by part (iii), via part (iv), (6.11) follows immediately from (5.1) and (6.1).

Next, the above expressions will be used to write lateral Taylor expansions for  $\ell(\cdot; x)$  and  $\partial_{\lambda} \ell(\cdot; x)$  around the origin.

**Lemma 6.3.** Suppose that Assumptions 5.1 and 5.2 hold. In this case, the following expansions are valid for  $x \in \mathbb{R}$  and  $\alpha \in (-\delta, \delta) \setminus \{0\}$ :

$$\ell(lpha, x) - \ell(0; x)$$

$$= |\alpha|^r \alpha \left\{ \sum_{k=r+1}^{2r+1} \frac{\partial_\lambda^k \ell(0+;x)}{k!} |\alpha|^{k-r-1} \right.$$

$$+ \left(\frac{\partial_{\lambda}^{2r+2}\ell(0+;x) + I_{(-\infty,0)}(\alpha)\Delta(x)}{(2r+2)!} + \frac{W^{*}(\alpha,x)}{(2r+3)!}\alpha\right)|\alpha|^{r}\alpha\right\},\$$

and

$$\partial_{\lambda}\ell(\alpha,x)$$
 20  
21

$$= |\alpha|^{r} \left\{ \sum_{k=r}^{2^{r}} \frac{\partial_{\lambda}^{k+1} \ell(0+;x)}{k!} |\alpha|^{k-r} \right\}$$
22
23
24

$$+ \left(\frac{\partial_{\lambda}^{2r+2}\ell(0+;x) + I_{(-\infty,0)}(\alpha)\Delta(x)}{(2r+1)!} + \frac{\tilde{W}(\alpha,x)}{(2r+2)!}\alpha\right) |\alpha|^r \alpha \right\},$$

where  $\Delta(\cdot)$  is as in (6.1), and

(6.12) 
$$|W^*(\alpha, x)|, \ |\tilde{W}(\alpha, x)| \le W(x).$$

*Proof.* Select  $\alpha_0 \neq 0$  with the same sign as  $\alpha$  and  $|\alpha_0| < |\alpha|$ , so that the closed interval joining  $\alpha_0$  and  $\alpha$  is contained in  $(-\delta, \delta) \setminus \{0\}$ . Since  $\ell(\cdot; x)$  has derivatives up to order 2r + 3 outside 0, there exist points  $\alpha^*$  and  $\tilde{\alpha}$  between  $\alpha_0$  and  $\alpha$  such that the following Taylor expansions hold:

(6.13) 
$$\ell(\alpha, x) - \ell(\alpha_0; x) = \sum_{k=1}^{2r+2} \frac{\partial_\lambda^k \ell(\alpha_0; x)}{k!} (\alpha - \alpha_0)^k + \frac{\partial_\lambda^{2r+3} \ell(\alpha^*; x)}{(2r+3)!} (\alpha - \alpha_0)^{2r+3},$$

and

(6.14) 
$$\partial_{\lambda}\ell(\alpha, x) = \sum_{k=0}^{2r+1} \frac{\partial_{\lambda}^{k+1}\ell(\alpha_0; x)}{k!} (\alpha - \alpha_0)^k + \frac{\partial_{\lambda}^{2r+3}\ell(\tilde{\alpha}; x)}{(2r+2)!} (\alpha - \alpha_0)^{2r+2},$$

where

$$\begin{array}{l} {}^{47}_{48} & \left(6.15\right) \\ \end{array} \qquad \left|\partial_{\lambda}^{2r+3}\ell(\alpha^*;x)\right|, \left|\partial_{\lambda}^{2r+3}\ell(\tilde{\alpha};x)\right| \le W(x), \end{array}$$

by Assumption 5.2(ii). Next, the conclusions in the lemma will be obtained taking lateral limits as  $\alpha_0$  goes to zero. Recall that  $\ell(\cdot; x)$  is continuous and consider the following exhaustive cases:

З

**Case 1:**  $\alpha > 0$ . taking the limit as  $\alpha_0$  decreases to zero, the above displayed relations and Lemma 6.2(i) together yield

$$\ell(\alpha, x) - \ell(0; x)$$

$$=\sum_{k=r+1}^{2r+2} \frac{\partial_{\lambda}^{k} \ell(0+;x)}{k!} \alpha^{k} + \frac{W^{*}(\alpha;x)}{(2r+3)!} \alpha^{2r+3}$$
<sup>5</sup>
<sub>6</sub>

$$=\sum_{k=r+1}^{2r+1} \frac{\partial_{\lambda}^{k}\ell(0+;x)}{k!} \alpha^{k} + \frac{\partial_{\lambda}^{2r+2}\ell(0+;x)}{(2r+2)!} \alpha^{2r+2} + \frac{W^{*}(\alpha;x)}{(2r+3)!} \alpha^{2r+3}$$

$$\lim_{k=r+1} \kappa : (2r+2): (2r+3):$$

$$= \alpha^{r+1} \left\{ \sum_{k=r+1}^{\infty} \frac{\sigma_{\lambda} c(\sigma+x)}{k!} \alpha^{k-r-1} \right\}$$

$$+ \left(\frac{\partial_{\lambda}^{2r+2}\ell(0+;x)}{(2r+2)!} + \frac{W^*(\alpha;x)}{(2r+3)!}\alpha\right)\alpha^{r+1}\right\},\$$

and

$$= \sum_{k=r}^{2r} \frac{\partial_{\lambda}^{k+1} \ell(0+;x)}{k!} \alpha^{k} + \frac{\partial_{\lambda}^{2r+2} \ell(0+;x)}{(2r+1)!} \alpha^{2r+1} + \frac{\tilde{W}(\alpha;x)}{(2r+2)!} \alpha^{2r+2}$$

$$= \alpha^r \left\{ \sum_{k=r}^{2r} \frac{\partial_{\lambda}^{k+1} \ell(0+;x)}{k!} \alpha^{k-r} \right\}$$

where  $W^*(\alpha, x)$  is given by  $W^*(\alpha, x) := \lim_{\alpha_0 \searrow 0} \partial_{\lambda}^{2r+3} \ell(\alpha^*; x)$  and, similarly,  $\tilde{W}(\alpha,x) := \lim_{\alpha_0 \searrow 0} \tilde{\partial}_{\lambda}^{2r+3} \ell(\tilde{\alpha};x)$ , so that

$$|W^*(\alpha, x)|, |\tilde{W}(\alpha, x)| \le W(x),$$

by (6.15); since  $\alpha$  is positive, so that  $I_{(-\infty,0)}(\alpha) = 0$ , these last three displays are equivalent to (6.12) - (6.12). 

**Case 2:**  $\alpha < 0$ . In this context, taking the limit as  $\alpha_0$  increases to zero in (6.13) and (6.14), Lemma 6.2(i) yields that

$$\ell(\alpha, x) - \ell(0; x) = \sum_{k=r+1}^{2r+2} \frac{\partial_{\lambda}^{k} \ell(0-; x)}{k!} \alpha^{k} + \frac{W^{*}(\alpha; x)}{(2r+3)!} \alpha^{2r+3}$$

and

$$\partial_{\lambda}\ell(\alpha;x) = \sum_{k=r}^{2r+1} \frac{\partial_{\lambda}^{k+1}\ell(0-;x)}{k!} \alpha^k + \frac{\tilde{W}(\alpha;x)}{(2r+2)!} \alpha^{2r+2}$$

where, analogously to the previous case,  $W^*(\alpha, x) := \lim_{\alpha_0 \nearrow 0} \partial_{\lambda}^{2r+3} \ell(\alpha^*; x)$  and  $\tilde{W}(\alpha, x) := \lim_{\alpha_0 \nearrow 0} \partial_{\lambda}^{2r+3} \ell(\tilde{\alpha}; x)$  so that, again, (6.15) implies that (6.12) is valid. Observe now that Lemma 6.2(iii) allows to write

$$\sum_{k=r+1}^{2r+2} \frac{\partial_{\lambda}^{k} \ell(0-;x)}{k!} \alpha^{k} + \frac{W^{*}(\alpha;x)}{(2r+3)!} \alpha^{2r+3}$$

З

#### Cavazos-Cadena and González-Farías

$$=\sum_{\substack{k=r+1\\k=r+1}}^{2r+1} \frac{\partial_{\lambda}^{k}\ell(0+;x)}{k!} (-1)^{k-1}\alpha^{k} + \frac{\partial_{\lambda}^{2r+2}\ell(0-;x)}{(2r+2)!}\alpha^{2r+2} + \frac{W^{*}(\alpha;x)}{(2r+3)!}\alpha^{2r+3}$$

$$=\sum_{k=1}^{2r+1} \frac{\partial_{\lambda}^{k}\ell(0+;x)}{|x|} |\alpha|^{k-1} \alpha + \left(\frac{\partial_{\lambda}^{2r+2}\ell(0-;x)}{(2r+2)!} + \frac{W^{*}(\alpha;x)}{(2r+2)!} \alpha\right) (|\alpha|^{r} \alpha)^{2}$$

$$\sum_{k=r+1}^{2r+1} k! \qquad (2r+2)! \qquad (2r+3)! \quad j \leq r \leq 2r+3$$

$$= |\alpha|^r \alpha \left\{ \sum_{k=r+1}^{\infty} \frac{\partial_\lambda^{\nu}(0+;x)}{k!} |\alpha|^{k-r-1} \right\}$$

$$+\left(\frac{\partial_{\lambda}^{2r+2}\ell(0-;x)}{(2r+2)!}+\frac{W^{*}(\alpha;x)}{(2r+3)!}\alpha\right)|\alpha|^{r}\alpha\right\}$$

and

$$\sum_{k=r}^{2r+1} \frac{\partial_{\lambda}^{k+1}\ell(0-;x)}{k!} \alpha^k + \frac{\tilde{W}(\alpha;x)}{(2r+2)!} \alpha^{2r+2}$$

$$=\sum_{k=r}^{2r} \frac{\partial_{\lambda}^{k+1}\ell(0+;x)}{k!} (-1)^k \alpha^k + \frac{\partial_{\lambda}^{2r+2}\ell(0-;x)}{(2r+1)!} \alpha^{2r+1} + \frac{\tilde{W}(\alpha;x)}{(2r+2)!} \alpha^{2r+2}$$

$$=\sum_{k=r}^{n-r} \frac{\partial_{\lambda}^{k+1}\ell(0+;x)}{k!} |\alpha|^{k} + |\alpha|^{r} \left(\frac{\partial_{\lambda}^{2r+2}\ell(0-;x)}{(2r+1)!} + \frac{\tilde{W}(\alpha;x)}{(2r+2)!}\alpha\right) |\alpha|^{r}\alpha$$

$$= |\alpha|^r \left\{ \sum_{\lambda=1}^{2r} \frac{\partial_{\lambda}^{k+1} \ell(0+;x)}{k!} |\alpha|^{k-r} \right\}$$

$$\begin{pmatrix} \sum_{k=r} & k! \\ & \begin{pmatrix} \partial_{2}^{2r+2}\ell(0-;x) & \tilde{W}(\alpha;x) \end{pmatrix} \end{pmatrix} = 0$$

$$+\left(\frac{\partial_{\lambda}^{2r+2}\ell(0-;x)}{(2r+1)!}+\frac{W(\alpha;x)}{(2r+2)!}\alpha\right)|\alpha|^{r}\alpha\right\};$$

using that  $\partial_{\lambda}^{2r+2}\ell(0-;x) = \partial_{\lambda}^{2r+2}\ell(0+;x) + I_{(-\infty,0)}(\alpha)\Delta(x)$ , by (6.1), the last four displays together yield that (6.12) and (6.12) are also valid for  $\alpha < 0$ .

After the above preliminaries, the main result of this section can be established as follows.

PROOF OF THEOREM 6.1. (i) Since  $L_n(\cdot; X_1^n)$  is the average of  $\ell(\cdot; X_i)$ , i = 1, 2, ..., n, by Lemma 6.3 the two indicated expansions hold with

$$W_n^*(\alpha) = \frac{1}{n} \sum_{i=1}^n W^*(\alpha, X_i) \quad \text{and} \quad \tilde{W}_n(\alpha) = \frac{1}{n} \sum_{i=1}^n \tilde{W}(\alpha, X_i),$$

so that (6.4) is satisfied, by (6.12).

(ii) Under  $\mathcal{H}_0$ :  $\lambda = 0$ ,  $X_i$  has symmetric distribution around zero with finite moment of order 4r + 2, by Assumption 5.2(i), so that a random variable of the form  $|X_i|^{k-1}X_i$  has zero expectation and finite second moment when  $1 \le k < 2r+2$ . Thus, from Lemma 6.2, for all  $k = r + 1, \ldots 2r + 1$ ,

$$E_0[\partial_\lambda^k \ell(0+;X_i)] = 0 \quad \text{and} \quad E_0[(\partial_\lambda^k \ell(0+;X_i))^2] = v_k < \infty,$$

as well as  $E_0[\Delta(X_1)] = 0$ . From this point, (5.9) and (6.9) together yield

$$\frac{2}{(2r+2)!}E[\partial_{\lambda}^{2r+2}\ell(0+;X_i)] = -\left(\frac{1}{(r+1)!}\right)^2 E[(\partial_{\lambda}^{r+1}\ell(0+;X_i))^2]$$

$$\frac{1}{(2r+2)!} E[\partial_{\lambda} + \ell(0+;\Lambda_i)] = -\left(\frac{1}{(r+1)!}\right) E[(\partial_{\lambda} + \ell(0+;\Lambda_i))]$$

$$= -4\left(\frac{D'g(0+)}{(r+1)!}\right)^{-} E[X_{i}^{2r+2}] = -V_{r+1}.$$

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

where Lemma 6.2(ii) was used to set the second equality (notice that this shows that  $V_{r+1}$  is the variance of  $\partial_{\lambda}^{r+1}\ell(0+;X_i)/(r+1)!$ ). Now, (6.5) and (6.6) follow from the central limit theorem and the strong law of large numbers, respectively.

## 7. Proof of Theorem 5.1

After the previous preliminaries, Theorem 5.1 is finally established in this section. The core of the argument has been decoupled into two lemmas showing that, under  $\mathcal{H}_0: \lambda = 0$ , along a consistent sequence  $\{\hat{\lambda}_n\}$  of maximum likelihood estimators, (i) the expansions in Theorem 6.1 can be simplified substantially, and (ii) that  $\hat{\lambda}_n$ is no null that with probability converging to 1.

**Lemma 7.1.** Suppose that Assumptions 5.1 and 5.2 hold and let  $\{\hat{\lambda}_n\}$  be a consistent sequence of maximum likelihood estimators of  $\lambda$ . In this case, assertions (i) and (ii) below occur under the null hypothesis  $\mathcal{H}_0: \lambda = 0$ .

(i) On the event  $[|\hat{\lambda}_n| < \delta]$ , the following expressions are valid:

=

(7.1) 
$$L_n(\hat{\lambda}_n; X_1^n) - L_n(0; X_1^n)$$

$$\hat{\lambda}_{n} |\hat{\lambda}_{n}|^{r} \left[ \frac{\partial_{\lambda}^{r+1} L_{n}(0+;X_{1}^{n})}{(r+1)!} + A_{n} - \frac{B_{n}}{2} \hat{\lambda}_{n} |\hat{\lambda}_{n}|^{r} \right],$$

and

(7.2) 
$$\frac{\partial_{\lambda} L_n(\hat{\lambda}_n; X_1^n)}{r+1} = |\hat{\lambda}_n|^r \left[ \frac{\partial_{\lambda}^{r+1} L_n(0+; X_1^n)}{(r+1)!} + \tilde{A}_n - \tilde{B}_n \hat{\lambda}_n |\hat{\lambda}_n|^r \right],$$

where

$$\begin{array}{c} {}^{30}\\ {}^{31}\\ {}^{32}\\ {}^{33} \end{array}$$

$$\begin{array}{c} A_n = O_p\left(\frac{\hat{\lambda}_n}{\sqrt{n}}\right), \quad \tilde{A}_n = O_p\left(\frac{\hat{\lambda}_n}{\sqrt{n}}\right), \\ \tilde{A}_n = O_p\left(\frac{\hat{\lambda}_n}{\sqrt{n}}\right), \quad \tilde{A}_n = O_p\left(\frac{\hat{\lambda}_n}{\sqrt{n}}\right), \\ \end{array}$$

and

(7.4) 
$$\lim_{n \to \infty} B_n = \lim_{n \to \infty} \tilde{B}_n = V_{r+1} \quad P_0 \text{-}a.s.$$

(ii) Consequently,

$$2n[L_n(\hat{\lambda}_n; X_1^n) - L_n(0; X_1^n)] = 2n\hat{\lambda}_n |\hat{\lambda}_n|^r \left[ \frac{\partial_{\lambda}^{r+1} L_n(0+; X_1^n)}{(r+1)!} + A_n - \frac{B_n}{2}\hat{\lambda}_n |\hat{\lambda}_n|^r \right] + o_p(1).$$

$$A_n := \sum_{k=1}^{2r+1} \frac{\partial_{\lambda}^k L_n(0+;X_1^n)}{k!} |\hat{\lambda}_n|^{k-r-1}$$

$$n = \sum_{k=r+2}^{n} \frac{|\lambda_n|}{k!}$$

and

$$B_n := - \left( \frac{2\partial_{\lambda}^{2r+2}L_n(0+;X_1^n) + 2\Delta_n I_{(-\infty,0)}(\hat{\lambda}_n)}{(2r+2)!} + \frac{2W_n^*(\hat{\lambda}_n)}{(2r+3)!}\hat{\lambda}_n \right),$$

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

(7.1) is equivalent to (6.2) with  $\alpha = \hat{\lambda}_n$ . Similarly, defining

$$\tilde{A}_n := \frac{1}{\sqrt{2r}} \sum_{\lambda=1}^{2r} \frac{\partial_{\lambda}^{k+1} L_n(0+;X_1^n)}{\partial_{\lambda} |\hat{\lambda}_n|^{k-r}}$$

$$r+1 \sum_{k=r+1} k!$$

and

$$\tilde{B}_n := -\frac{1}{r+1} \left( \frac{\partial_{\lambda}^{2r+2} L_n(0+;X_1^n) + \Delta_n I_{(-\infty,0)}(\hat{\lambda}_n)}{(2r+1)!} + \frac{\tilde{W}_n(\hat{\lambda}_n)}{(2r+2)!} \hat{\lambda}_n \right)$$

$$= -\left(\frac{2\partial_{\lambda}^{2r+2}L_{n}(0+;X_{1}^{n})+2\Delta_{n}I_{(-\infty,0)}(\hat{\lambda}_{n})}{(2r+2)!}+\frac{\tilde{W}_{n}(\hat{\lambda}_{n})}{(r+1)(2r+2)!}\hat{\lambda}_{n}\right)$$

it follows that (7.2) is equivalent to (6.3) with  $\alpha = \hat{\lambda}_n$ . Therefore, since (6.2) and (6.3) are valid for  $|\alpha| < \delta$ , (7.1) and (7.2) hold on the event  $[|\hat{\lambda}_n| < \delta]$ . To conclude, it will be shown that (7.3) and (7.4) are satisfied. First, notice that  $A_n$  and  $\tilde{A}_n$ defined above are null for r = 0, so that (7.3) certainly occurs in this case. On the other hand, if r > 0, then convergences (6.5) established in Theorem 6.1(ii) yield that  $\partial_{\lambda}^k L_n(0+; X_1^n) = O_p(1/\sqrt{n})$  for  $r+1 \le k < 2r+1$  and then (7.3) follows, since the above expressions for  $A_n$  and  $\tilde{A}_n$  involve factors  $|\hat{\lambda}_n|^s$  with  $s \ge 1$  and

$$(7.7) P_0[\hat{\lambda}_n \to 0] = 1, 23$$

by consistency. Next, observe that

$$|W_n^*(\hat{\lambda}_n)|, \ |\tilde{W}_n(\hat{\lambda}_n)| \le W_n = \frac{1}{n} \sum_{i=1}^n W(X_i),$$
26  
27  
28  
29

by (6.4), and then the strong law of large numbers and Assumption 5.2(ii) yield that

$$\limsup_{n \to \infty} |W_n^*(\hat{\lambda}_n)|, \ \limsup_{n \to \infty} |\tilde{W}_n(\hat{\lambda}_n)| \le \int_{\mathbb{R}} W(x) f(x), \ dx < \infty \quad P_0\text{-a.s.},$$

so that

$$\lim_{n \to \infty} W_n^*(\hat{\lambda}_n) \hat{\lambda}_n = 0 = \lim_{n \to \infty} \tilde{W}_n(\hat{\lambda}_n) \hat{\lambda}_n \quad P_0\text{-a.s.},$$

by (7.7). From this point, (6.6) in Theorem 6.1(ii) and the specifications of  $B_n$  and  $\tilde{B}_n$  lead to (7.4).

(ii) Since expansion (7.1) is valid on  $[|\hat{\lambda}_n| < \delta]$ , the conclusion follows from (7.7).

**Lemma 7.2.** Suppose that Assumptions 5.1 and 5.2 are valid, let  $\{\hat{\lambda}_n\}$  be a sequence of maximum likelihood estimators of  $\lambda$ , and define

$$\Omega_n^{**} := \left[ L_n(\hat{\lambda}_n; X_1^n) \ge L_n(\lambda; X_1^n), \ \lambda \in \mathbb{R} \right] \cap \left[ \partial_{\lambda}^{r+1} L_n(0+; X_1^n) \neq 0 \right].$$

48 With this notation, assertions (i) and (ii) below occur. 49 (i)  $\hat{\lambda}_n \neq 0$  on  $\Omega_n^{**}$ . 50 Consequently, 51 (ii)  $P_{\mathbf{0}}[\hat{\lambda}_n \neq 0] \rightarrow 1$  as  $n \rightarrow \infty$ .

З

*Proof.* (i) The expansion for  $L_n(\cdot; X_1^n) - L_n(0; X_1^n)$  in Theorem 6.1(i) yields that

$$\lim_{\alpha \to 0} \frac{L_n(\alpha; X_1^n) - L_n(0; X_1^n)}{|\alpha|^r \alpha} = \partial_{\lambda}^{r+1} L_n(0+; X_1^n).$$

It follows that if  $\partial_{\lambda}^{r+1}L_n(0+;X_1^n) > 0$ , then  $L_n(\alpha;X_1^n) - L_n(0;X_1^n) > 0$  when  $\alpha$  is positive and small enough, whereas if  $\partial_{\lambda}^{r+1}L_n(0+;X_1^n) < 0$ , then  $L_n(\alpha;X_1^n) - L_n(0;X_1^n) > 0$  when  $\alpha < 0$  and  $|\alpha|$  is sufficiently small. Thus,  $\partial_{\lambda}^{r+1}L_n(0+;X_1^n) \neq 0$  implies that 0 is not a maximizer of  $L_n(\cdot;X_1^n)$ , so that, if  $\hat{\lambda}_n$  maximizes the average likelihood  $L_n(\cdot;X_1^n)$  and  $\partial_{\lambda}^{r+1}L_n(0+;X_1^n) \neq 0$  then  $\hat{\lambda}_n$  is no null, *i.e.*,  $\Omega_n^{**} \subset [\hat{\lambda}_n \neq 0]$ .

(ii) Since  $D^r g(0+) \neq 0$ , it follows that  $\partial_{\lambda}^{r+1} \ell(0+; X_i) = 2D^r g(0+)|X_i|^r X_i$  has a density, and then their average  $\partial_{\lambda}^{r+1} L_n(0+; X_1^n)$  is absolutely continuous. It follows that  $P_0[\partial_{\lambda}^{r+1} L_n(0+; X_1^n) \neq 0] = 1$ , and then  $P_0[\Omega_n^{**}] \to 1$ , by Definition 3.1 (see Remark 3.1(ii)) and, via part (i), the conclusion follows.

PROOF OF THEOREM 5.1. Suppose that Assumptions 5.1 and 5.2 hold, that the hypothesis  $\mathcal{H}_0: \lambda = 0$  occurs, and let  $\{\hat{\lambda}_n\}$  be a consistent sequence of maximum likelihood estimators. In this context, define

$$\Omega_{n,*} := \left[ L_n(\hat{\lambda}_n; X_1^n) \ge L_n(\lambda; X_1^n), \ \lambda \in \mathbb{R} \right] \cap \left[ 0 < |\hat{\lambda}_n| < \delta \right],$$

and notice that the conclusions in Lemma 7.1 occur on this event, since  $\Omega_{n,*} \subset [|\hat{\lambda}| < \delta]$ . Also, the consistency of  $\{\hat{\lambda}_n\}$ , Definition 3.1 and Lemma 7.2(ii) together imply that

(7.8) 
$$\lim_{n \to \infty} P_0[\Omega_{n,*}] = 1.$$
 26  
27

Observe now that, on the event  $\Omega_{n,*}$ , the estimator  $\hat{\lambda}_n$  is no null and maximizes  $L_n(\cdot; X_1^n)$ , so that the likelihood equation  $\partial_{\lambda}L_n(\hat{\lambda}_n; X_1^n) = 0$  holds; see Remark 5.1(ii). Via (7.2) it follows that

On 
$$\Omega_{n,*}$$
,  $\tilde{B}_n \sqrt{n} \hat{\lambda}_n |\hat{\lambda}_n|^r = \sqrt{n} \frac{\partial_{\lambda}^{r+1} L_n(0+;X_1^n)}{(r+1)!} + \sqrt{n} \tilde{A}_n$ ,

and then, from (7.8)

$$\tilde{B}_n \sqrt{n} \hat{\lambda}_n |\hat{\lambda}_n|^r = \sqrt{n} \frac{\partial_{\lambda}^{r+1} L_n(0+;X_1^n)}{(r+1)!} + \sqrt{n} \tilde{A}_n + o_p(1)$$

=

=

$$= \sqrt{n} \frac{\partial_{\lambda}^{r+1} L_n(0+;X_1^n)}{(r+1)!} + O_p(\hat{\lambda}_n) + o_p(1)$$

$$= \sqrt{n} \frac{\partial_{\lambda}^{r+1} L_n(0+;X_1^n)}{(r+1)!} + o_p(1)$$

where (7.3) was used to set the second equality, and the third one stems from  $P_0[\hat{\lambda}_n \to 0] = 1$ , by consistency; via (6.5) this yields that  $\tilde{B}_n \sqrt{n} \hat{\lambda}_n |\hat{\lambda}_n|^r = O_p(1)$ , and then (7.4) leads to

<sup>50</sup>  
<sub>51</sub> (7.9) 
$$V_{r+1}\sqrt{n} \ \hat{\lambda}_n |\hat{\lambda}_n|^r = \sqrt{n} \frac{\partial_{\lambda}^{r+1} L_n(0+;X_1^n)}{(r+1)!} + o_p(1).$$
 <sup>50</sup>  
<sub>51</sub> 51

imsart-coll ver. 2008/08/29 file: Cavazos.tex date: April 10, 2009

З
(i) Using that  $\sqrt{n} \frac{\partial_{\lambda}^{r+1} L_n(0+; X_1^n)}{(r+1)!} \xrightarrow{\mathrm{d}} \mathcal{N}(0, V_{r+1})$  (see (5.9) and (6.5)), the above display yields

(7.10) 
$$\sqrt{nV_{r+1}} \ \hat{\lambda}_n |\hat{\lambda}_n|^r \xrightarrow{\mathrm{d}} Z$$
 where Z has standard normal distribution;

since the inverse of the function  $x \mapsto x|x|^r$  is the continuous mapping  $x \mapsto |x|^{1/(r+1)}$ sign(x), it follows that

$$(nV_{r+1})^{1/(2(r+1))}\hat{\lambda}_n \xrightarrow{\mathrm{d}} |Z|^{1/(r+1)}\mathrm{sign}(Z).$$

(ii) Since  $\sqrt{n}\hat{\lambda}_n|\hat{\lambda}_n|^r = O_p(1)$ , by part (i), Lemma 7.1(ii), (7.3), (7.4) and (7.9) together yield

$$2n[L_n(\hat{\lambda}_n; X_1^n) - L_n(0; X_1^n)]$$

$$= 2n\hat{\lambda}_{n}|\hat{\lambda}_{n}|^{r} \left[ \frac{\partial_{\lambda}^{r+1}L_{n}(0+;X_{1}^{n})}{(r+1)!} + A_{n} - \frac{B_{n}}{2}\hat{\lambda}_{n}|\hat{\lambda}_{n}|^{r} \right] + o_{p}(1)$$

$$= 2\sqrt{n}\hat{\lambda}_{n}|\hat{\lambda}_{n}|^{r} \left[\sqrt{n}\frac{\partial_{\lambda}^{r+1}L_{n}(0+;X_{1}^{n})}{(r+1)!} + \sqrt{n}A_{n} - \frac{B_{n}}{2}\sqrt{n}\hat{\lambda}_{n}|\hat{\lambda}_{n}|^{r}\right] + o_{p}(1)$$

 $-\frac{V_{r+1}+o_p(1)}{2}\sqrt{n}\hat{\lambda}_n|\hat{\lambda}_n|^r\right]+o_p(1)$ 

$$= 2\sqrt{n}\hat{\lambda}_n |\hat{\lambda}_n|^r \left[ \sqrt{n} \frac{\partial_{\lambda}^{r+1} L_n(0+;X_1^n)}{(r+1)!} + O_p(\hat{\lambda}_n) \right]$$

$$= 2\sqrt{n}\hat{\lambda}_{n}|\hat{\lambda}_{n}|^{r} \left[\sqrt{n}\frac{\partial_{\lambda}^{r+1}L_{n}(0+;X_{1}^{n})}{(r+1)!} - \frac{V_{r+1}}{2}\sqrt{n}\hat{\lambda}_{n}|\hat{\lambda}_{n}|^{r}\right] + o_{p}(1)$$

$$= 2\sqrt{n}\hat{\lambda}_n |\hat{\lambda}_n|^r \left[ V_{r+1}\sqrt{n}\hat{\lambda}_n |\hat{\lambda}_n - \frac{V_{r+1}}{2}\sqrt{n}\hat{\lambda}_n |\hat{\lambda}_n|^r \right] + o_p(1)$$

$$= \left(\sqrt{nV_{r+1}} \ \hat{\lambda}_n |\hat{\lambda}_n|^r\right)^2 + o_p(1);$$

together with (7.10), this yields that  $2n[L_n(\hat{\lambda}_n; X_1^n) - L_n(0; X_1^n)] \xrightarrow{d} Z^2$ , completing the proof.

#### References

- AZZALINI, A. (1985). A class of distributions which includes the normal ones. Scand. J. Statist., 12, 171–178.
   AZZALINI, A. (1986). Further results on a class of distributions which includes the normal ones.
  - Statistica, 46, 199–208.
     [3] AZZALINI, A. and CAPITANIO, A. (1999). Statistical applications of the multivariate skew normal
  - [5] AZZALINI, A. and CAPITANIO, A. (1999). Statistical applications of the multivariate skew normal distribution. J. R. Stat. Soc., Ser. B, **61**, 579–602.
- [4] AZZALINI, A. and DALLA VALLE, A. (1996). The multivariate skew-normal distribution. *Biometrika* 83, 715–726.
- [5] AZZALINI, A. (2005). The skew-normal distribution and related multivariate families (with discussion). Scand. J. Statist., 32 159–188 (CR: 189–200).
- [6] BILLINGSLEY, P. (1995). Probability and Measure. Third Edition, Wiley, New York.
- 46 [7] CAVAZOS-CADENA R. and GONZÁLEZ-FARÍAS G. (2007). Necessary and sufficient conditions for the consistency of maximum likelihood estimators. Submitted.
- 47 [8] CHIOGNA, M. (2005). A note on the asymptotic distribution of the maximum likelihood estimator
  48 for the scalar skew-normal distribution. *Stat. Meth. & Appl.*, **14**, 331–341.
- 49 [9] DICICICIO, T. and MONTI, A. C. (2004). Inferential aspects of the skew exponential power distribution. J. Amer. Statist. Assoc., 99, 439-450.
   50 [10] CINTON M. C. (Editor) (2004). Show Elliptical Distributions and Their Applications. Chapman &
- 50515050505051Hall, London.51

З

| 1  | [11] LEHMANN E. L. and CASELLA G. (1998). Theory of Point Estimation, Second Edition. Sprin   | iger, 1         |
|----|---|-----------------|
| 2  | New York.<br>[12] NEWEY W and MCEADDEN D (1002) Estimation in large samples. In D. McEaddon and R. En   | 2               |
| 3  | (eds), Handbook of Econometrics, Vol. 4. North-Holland, Amsterdam.  | giei 3          |
| 4  | 13] PEWSEY, A. (2000). Problems of inference for Azzalini's skew-normal distribution. J. Appl. Stat   | <i>ist.</i> , 4 |
| 5  | 27.7, 859–870.<br>14 ROTNITZKY A COX D R BOTTAL M and ROBINS J (2000) Likelihood-based inference v  | 5<br>with       |
| 6  | singular information matrix. Bernoulli, 6, 243–284.   | 6               |
| 7  | 15] SARTORI, N. (2006). Bias prevention of maximum likelihood estimates for scalar skew normal  | and 7           |
| 8  | skew t distributions. J. Statist. Plann. Inference, <b>136</b> , 4259–4275.<br>[16] SHAO, J. (1999). Mathematical Statistics. Springer. New York. | 8               |
| 9  | -1,. (,   | 9               |
| 10 |   | 10              |
| 11 |   | 11              |
| 12 |   | 12              |
| 13 |   | 13              |
| 14 |   | 14              |
| 15 |   | 15              |
| 17 |   | 10              |
| 18 |   | 18              |
| 19 |   | 19              |
| 20 |   | 20              |
| 21 |   | 21              |
| 22 |   | 22              |
| 23 |   | 23              |
| 24 |   | 24              |
| 25 |   | 25              |
| 26 |   | 26              |
| 27 |   | 27              |
| 28 |   | 28              |
| 29 |   | 29              |
| 30 |   | 30              |
| 31 |   | 31              |
| 32 |   | 32              |
| 33 |   | 33              |
| 34 |   | 34              |
| 35 |   | 35              |
| 36 |   | 36              |
| 37 |   | 37              |
| 38 |   | 38              |
| 39 |   | 39              |
| 40 |   | 40              |
| 42 |   | 41              |
| 43 |   | 43              |
| 44 |   | 44              |
| 45 |   | 45              |
| 46 |   | 46              |
| 47 |   | 47              |
| 48 |   | 48              |
| 49 |   | 49              |
| 50 |   | 50              |
| 51 |   | 51              |

# Parametric Mixture Models for Estimating the Proportion of True Null Hypotheses and Adaptive Control of FDR

З

### Ajit C. Tamhane<sup>1</sup> and Jiaxiao Shi<sup>2</sup>

 $Northwestern \ University$ 

Abstract: Estimation of the proportion or the number of true null hypotheses is an important problem in multiple testing, especially when the number of hypotheses is large. Wu, Guan and Zhao (2006) found that nonparametric approaches are too conservative. We study two parametric mixture models (normal and beta) for the distributions of the test statistics or their *p*-values to address this problem. The components of the mixture are the null and alternative distributions with mixing proportions  $\pi_0$  and  $1 - \pi_0$ , respectively, where  $\pi_0$  is the unknown proportion to be estimated. The normal model assumes that the test statistics from the true null hypotheses are i.i.d. N(0,1) while those from the alternative hypotheses are i.i.d.  $N(\delta, 1)$  with  $\delta \neq 0$ . The beta model assumes that the p-values from the null hypotheses are i.i.d. U[0, 1] and those from the alternative hypotheses are i.i.d. Beta(a, b) with a < 1 < b. All parameters are assumed to be unknown. Three methods of estimation of  $\pi_0$  are developed for each model. The methods are compared via simulation with each other and with Storey's (2002) nonparametric method in terms of the bias and mean square error of the estimators of  $\pi_0$  and the achieved FDR. Robustness of the estimators to the model violations is also studied by generating data from other models. For the normal model, the parametric methods perform better compared to Storey's method with the EM method (Dempster, Laird and Rubin 1977) performing best overall when the assumed model holds; however, it is not very robust to significant model violations. For the beta model, the parametric methods do not perform as well because of the difficulties of estimation of parameters, and Storey's nonparametric method turns out to be the winner in many cases. Therefore the beta model is not recommended for use in practice. An example is given to illustrate the methods. Contents  $\mathbf{2}$ 2.1Test Statistics (TS) Method ..... 308 2.22.33.13.2

<sup>1</sup>Department of IEMS, Northwestern University, Evanston, IL 60208

<sup>2</sup>Department of Statistics, Northwestern University, Evanston, IL 60208

AMS 2000 subject classifications: Primary 62F10; secondary 62F12

*Keywords and phrases:* Beta model; Bias-correction; EM method; Mixture model; False discovery rate (FDR); Least squares method; Maximum likelihood method; Normal model; *p*-values. 51

imsart-coll ver. 2008/08/29 file: Tamhane.tex date: April 10, 2009

| 1  | <b>3.3 EM Method</b>                    | 1  |
|----|---|----|
| 2  | 4 Adaptive Control of FDR               | 2  |
| 3  | 5 Simulation Results                    | 3  |
| 4  | 5.1 Simulation Results for Normal Model | 4  |
| 5  | 5.2 Simulation Results for Beta Model   | 5  |
| 6  | 6 Example                               | 6  |
| 7  | 7 Concluding Remarks                    | 7  |
| 8  | Acknowledgment                          | 8  |
| 9  | Appendix A                              | 9  |
| 10 | <b>References</b>                       | 10 |
| 11 |   | 11 |

# 1. Introduction

Suppose that m null hypotheses,  $H_{01}, \ldots, H_{0m}$ , are to be tested against alternatives,  $H_{11}, \ldots, H_{1m}$ . Let  $X_1, \ldots, X_m$  be the test statistics and  $p_1, \ldots, p_m$  their p-values. Throughout we assume that the  $X_i$ 's and hence the  $p_i$ 's are mutually independent. Suppose that some unknown number  $m_0$  of the hypotheses are true and  $m_1 =$  $m-m_0$  are false. We wish to estimate  $m_0$  or equivalently the proportion  $\pi_0 = m_0/m$ of the true hypotheses based on the  $X_i$ 's or equivalently the  $p_i$ 's. The estimate  $\hat{m}_0$  is useful for devising more powerful adaptive multiple comparison procedures (MCPs) to control an appropriate type I error rate, e.g., the familywise error rate (FWE) (Hochberg and Tamhane 1987) in the Bonferroni procedure or the false discovery rate (FDR) in the Benjamini and Hochberg (1995) procedure. These procedures normally use the total number m as a conservative upper bound on the number of true hypotheses. Adaptive procedures based on  $\hat{m}_0$  are especially useful in largescale multiplicity testing problems arising in microarray data involving m of the order of several thousands.

A number of methods have been proposed for estimating  $m_0$  starting with Schweder and Spjøtvoll (1982); see, e.g., Hochberg and Benjamini (1990), Ben-jamini and Hochberg (2000), Turkheimer, Smith and Schmidt (2001), Storey (2002), Storey et al. (2004), Jiang and Doerge (2005) and Langaas et al. (2005). Many of these methods reject the p-values that differ significantly from the null U[0, 1] distri-bution as non-null and exclude them from the estimation process. Different formal or graphical tests are used for this purpose. For example, consider Storey's (2002) method with a fixed  $\lambda$ -level test for a sufficiently large  $\lambda$  (e.g.,  $\lambda = 0.5$ ) to reject any p-value  $\leq \lambda$  as non-null. (It should be noted that  $\lambda$  is not fixed but is a tuning parameter whose value is determined from the data to minimize the mean square error of the estimate of  $\pi_0$  using bootstrap.) Let  $N_r(\lambda) = \sharp(p_i \leq \lambda)$  denote the number of rejected hypotheses and  $N_a(\lambda) = \sharp(p_i > \lambda)$  the number of accepted hypotheses at level  $\lambda \in (0, 1)$ . If type II errors are ignored for a sufficiently large  $\lambda$ then 

(1.1) 
$$E[N_a(\lambda)] \approx m_0(1-\lambda).$$

Storey's (ST) estimator is given by 

(1.2) 
$$\widehat{\pi}_0(\lambda) = \frac{N_a(\lambda)}{m(1-\lambda)} \text{ or } \widehat{m}_0(\lambda) = \frac{N_a(\lambda)}{1-\lambda}.$$

Schweder and Spjøtvoll's (1982) method visually fits a straight line through the origin to the plot of  $N_a(p_{(i)}) = m - i$  vs.  $1 - p_{(i)}$   $(1 \le i \le m)$  for large values

З

of the  $p_{(i)}$ . The slope of the fitted line is taken as an estimate of  $m_0$  according to Equation (1.1). Because these estimators attribute all nonsignificant *p*-values to the true null hypotheses (type II errors are ignored) and do not explicitly model the non-null *p*-values, they tend to be positively biased which results in conservative adaptive control of any type I error rate.

To get a handle on type II errors, so that both the null and non-null *p*-values can be utilized to estimate  $\pi_0$ , the mixture model approach has been proposed by several authors. The mixture model differs from the setup given in the first paragraph in that the number of true hypotheses is a random variable (r.v.) and  $m_0$  is its expected value. Specifically, let  $Z_i$  be a Bernoulli r.v. which equals 1 with probability  $\pi_0$  if  $H_{0i}$  is true and 0 with probability  $\pi_1 = 1 - \pi_0$  if  $H_{0i}$  is false. Assume that the  $Z_i \ (1 \le i \le m)$  are i.i.d. Then the number of false hypotheses,  $M_0 = \sum_{i=1}^m Z_i$ , is a binomial r.v. with parameters m and  $\pi_0$ , and  $E(M_0) = m_0 = m\pi_0$ . 

A parametric mixture model was considered by Hsueh, Chen, and Kodell (2003) (HCK). They assumed the following simple hypothesis testing setup. Suppose that all m hypotheses pertain to the means of the normal distributions with  $H_{0i}: \mu_i = 0$ versus  $H_{1i}$ :  $\mu_i > 0$ . (HCK considered a two-sided alternative, but that is not germane to their method.) Conditional on  $Z_i$ , the test statistic  $X_i \sim N(\delta_i, 1)$ , where  $\delta_i$  is the standardized  $\mu_i$  with  $\delta_i = 0$  if  $Z_i = 1$  and  $\delta_i = \delta > 0$  if  $Z_i = 0$ where HCK assumed that  $\delta$  is known. We refer to this model as the normal model, which was also used by Black (2004) to study the bias of Storey's (2002) estimator. An expression for the expected number of  $X_i$ 's that are greater than any specified threshold can be derived using this setup. By plotting the corresponding observed number of  $X_i$ 's against the threshold,  $m_0$  could be estimated as the slope of the straight line through the origin using least squares (LS) regression. 

The normal model is the topic of Section 2. We first extend the HCK estimation method to the unknown  $\delta$  case, which is a nonlinear least squares (NLS) regression problem. Next we note that the HCK method makes use of only the number of  $X_i$ 's that are greater than a specified threshold; it does not make use of the *magnitudes* of the  $X_i$ 's. Therefore we propose two alternative methods of estimation which utilize the magnitudes of the  $X_i$ 's in an attempt to obtain a better estimate of  $\delta$  and thereby a better estimate of  $m_0$ . The first of these alternative methods is similar to the LS method of HCK, but uses the sample mean (instead of the number) of the  $X_i$ 's that are greater than a specified threshold. We refer to it as the test statistics (TS) method. The second method is the EM method of Dempster, Laird and Rubin (1977) which finds the maximum likelihood estimators (MLEs) of the mixture distribution of the  $X_i$ 's.

This normal mixture model approach in conjunction with the EM algorithm was proposed by Guan, Wu and Zhao (2004) and most recently by Iyer and Sarkar (2007). So, although the use of the EM algorithm for estimation in the context of the present problem is not new, we perform a comprehensive comparison of it with the other two new methods, and find that it performs best when the assumed model is correct, but is not robust to model violations.

In the second approach discussed in Section 3, the non-null *p*-values are modeled by a beta distribution with unknown parameters a and b (denoted by Beta(a, b)). We refer to this model as the *beta model*. Here we restrict to estimation methods based on *p*-values since the  $X_i$ 's can have different null distributions. All three estimators (HCK, TS and EM) are also derived for the beta model.

We stress that both the normal and beta models are simply "working" models
intended to get a handle on type II errors. We do not pretend that these models are
strictly true. Therefore robustness of the estimators to the model assumptions is an

З

important issue. In the simulation comparisons for the normal model, robustness of the fixed  $\delta$  assumption is tested by generating different  $\delta_i$ 's for the false hypotheses from a normal distribution. Robustness of the normal model assumption is tested by generating  $p_i$ 's for the false hypotheses using the beta model and transforming them to the  $X_i$ 's using the inverse normal transformation. Similarly, the robustness of the beta model is tested by generating  $X_i$ 's using the normal model and transforming them to  $p_i$ 's.

Adaptive control of FDR using different estimators of  $m_0$  is the topic of Section 4. The ST, HCK, TS and EM estimators are compared in a large simulation study in Section 5. The performance measures used in the simulation study are the biases and mean square errors of the estimators of  $\pi_0$  and FDR. An example illustrating application of the proposed methods is given in Section 6. Conclusions are summarized in Section 7. Proofs of some technical results are given in the Appendix.

## 2. Normal Model

The normal mixture model can be expressed as

(2.1) 
$$f(x_i) = \pi_0 \phi(x_i) + \pi_1 \phi(x_i - \delta),$$

where  $f(x_i)$  is the p.d.f. of  $X_i$  and  $\phi(\cdot)$  is the p.d.f. of the standard normal distribution. Although  $\delta$  will need to be estimated, we are not too concerned about its estimation accuracy since, after all, it is a parameter of a working model.

#### 2.1. Hsueh, Chen, and Kodell (HCK) Method

Let

(2.2) 
$$\beta(\delta,\lambda) = P_{H_{1i}}\{p_i > \lambda\} = P_{H_{1i}}\{X_i < z_\lambda\} = \Phi\left(z_\lambda - \delta\right)$$

denote the type II error probability of a test performed at level  $\lambda$  where  $\Phi(\cdot)$  is the standard normal c.d.f. and  $z_{\lambda} = \Phi^{-1}(1-\lambda)$ . Then  $E[N_r(\lambda)] = m_0\lambda + (m-m_0)[1-\beta(\delta,\lambda)]$ , and hence

(2.3) 
$$E[N_r(\lambda)] - m\Phi\left(-z_{\lambda} + \delta\right) = m_0[\lambda - \Phi\left(-z_{\lambda} + \delta\right)].$$

For  $\lambda = p_{(i)}$ , i = 1, 2, ..., m, the term inside the square brackets in the R.H.S. of the above equation is

(2.4) 
$$x_i = p_{(i)} - \Phi \left( -z_{p_{(i)}} + \delta \right)$$

and the L.H.S. can be estimated by

<sup>44</sup> (2.5) 
$$y_i = i - m\Phi \left( -z_{p_{(i)}} + \delta \right).$$

If  $\delta$  is assumed to be known then we can calculate  $(x_i, y_i)$ , i = 1, 2, ..., m, and using (2.3) fit an LS straight line through the origin by minimizing  $\sum_{i=1}^{m} (y_i - m_0 x_i)^2$  with respect to (w.r.t.)  $m_0$ . The LS estimator of  $m_0$  is given by

$$\widehat{m}_{0} = \frac{\sum_{i=1}^{m} x_{i} y_{i}}{\sum_{i=1}^{m} x_{i}^{2}}.$$

imsart-coll ver. 2008/08/29 file: Tamhane.tex date: April 10, 2009

Tamhane and Shi

We first extend the HCK estimator to the unknown  $\delta$  case by incorporating estimation of  $\delta$  as part of the NLS problem of minimizing  $\sum_{i=1}^{m} (y_i - m_0 x_i)^2$  w.r.t. 2  $m_0$  and  $\delta$ . The iterative algorithm for this purpose is given below. The initial values for this algorithm as well as the algorithms for the TS and EM estimators were determined by solving the following two moment equations for  $m_0$  and  $\delta$ : 

(2.7) 
$$\sum_{i=1}^{m} X_i = (m - m_0)\delta \text{ and } \sum_{i=1}^{m} X_i^2 = m_0 + (m - m_0)(\delta^2 + 1).$$

## HCK Algorithm

<sup>12</sup> Step 0: Compute initial estimates  $\widehat{m}_0$  and  $\widehat{\delta}$  by solving (2.7). Let  $\widehat{\pi}_0 = \widehat{m}_0/m$ . <sup>13</sup> Step 1: Set  $\delta = \widehat{\delta}$  and compute  $(x_i, y_i)$ , i = 1, 2, ..., m, using (2.4) and (2.5). <sup>14</sup> Step 2: Compute  $\widehat{m}_0$  using (2.6) and  $\widehat{\pi}_0 = \widehat{m}_0/m$ . <sup>15</sup> Step 3: Find  $\widehat{\delta}$  to minimize  $\sum_{i=1}^{m} (y_i - m_0 x_i)^2$ .

**Step 4:** Return to Step 1 until convergence.

**Remark:** One could use weighted least squares to take into account the heteroscedasticity of the  $y_i$ 's. We tried this, but the resulting NLS problem was computationally much more intensive without a collateral gain in the efficiency of the estimators.

#### 

#### 2.2. Test Statistics (TS) Method

As noted in Section 1, we hope to improve upon the HCK estimator by utilizing the information in the magnitudes of the  $X_i$ 's. Toward this end we first propose an estimator analogous to the HCK estimator except that it uses the sample mean (rather than the number) of the  $X_i$ 's that are significant at a specified level  $\lambda$ . Define

$$S_a(\lambda) = \{i : p_i > \lambda\} = \{i : X_i < z_\lambda\} \text{ and } S_r(\lambda) = \{i : p_i \le \lambda\} = \{i : X_i \ge z_\lambda\}.$$

Then  $N_a(\lambda) = |S_a(\lambda)|$  and  $N_r(\lambda) = |S_r(\lambda)|$ . Finally define

$$\overline{X}_a(\lambda) = \frac{1}{N_a(\lambda)} \sum_{i \in S_a(\lambda)} X_i \text{ and } \overline{X}_r(\lambda) = \frac{1}{N_r(\lambda)} \sum_{i \in S_r(\lambda)} X_i.$$

To derive the expected values of these sample means the following lemma is useful.

### Lemma 1. Define

$$c_{0a}(\lambda) = E_{H_{0i}}\left(X_i | X_i < z_{\lambda}\right), c_{0r}(\lambda) = E_{H_{0i}}\left(X_i | X_i \ge z_{\lambda}\right),$$

and

$$c_{1a}(\delta,\lambda) = E_{H_{1i}}\left(X_i|X_i < z_{\lambda}\right), c_{1r}(\delta,\lambda) = E_{H_{1i}}\left(X_i|X_i \ge z_{\lambda}\right).$$

Then

$$c_{0a}(\lambda) = -rac{\phi(z_{\lambda})}{1-\lambda}, c_{0r}(\lambda) = rac{\phi(z_{\lambda})}{\lambda}$$

49 and

$$c_{1a}(\delta,\lambda) = \delta - \frac{\phi(z_{\lambda} - \delta)}{\Phi(z_{\lambda} - \delta)}, c_{1r}(\delta,\lambda) = \delta + \frac{\phi(\delta - z_{\lambda})}{\Phi(\delta - z_{\lambda})}.$$

*Proof.* The proof follows from the following expressions for the conditional expectations of  $X \sim N(\mu, 1)$ :  $E(X|X \le x) = \mu - \frac{\phi(x-\mu)}{\Phi(x-\mu)}$  and  $E(X|X > x) = \mu + \frac{\phi(\mu-x)}{\Phi(\mu-x)}$ . The desired expected values of  $\overline{X}_a(\lambda)$  and  $\overline{X}_r(\lambda)$  are then given by the following lemma. Lemma 2. Let  $g(\pi_0, \delta, \lambda) = P\left\{Z_i = 1 | X_i < z_\lambda\right\} = \frac{\pi_0(1-\lambda)}{\pi_0(1-\lambda) + \pi_1 \Phi\left(z_\lambda - \delta\right)}$ (2.8)and  $h(\pi_0, \delta, \lambda) = P\left\{Z_i = 1 | X_i \ge z_\lambda\right\} = \frac{\pi_0 \lambda}{\pi_0 \lambda + \pi_1 \Phi\left(-z_\lambda + \delta\right)}.$ (2.9)Then $E[\overline{X}_a(\lambda)] = g(\pi_0, \delta, \lambda)c_{0a}(\lambda) + [1 - g(\pi_0, \delta, \lambda)]c_{1a}(\delta, \lambda)$ (2.10)and $E[\overline{X}_r(\lambda)] = h(\pi_0, \delta, \lambda)c_{0r}(\lambda) + [1 - h(\pi_0, \delta, \lambda)]c_{1r}(\delta, \lambda),$ (2.11)where  $c_{0a}(\lambda), c_{0r}(\lambda), c_{1a}(\delta, \lambda)$  and  $c_{1r}(\delta, \lambda)$  are as given in Lemma 1. *Proof.* Given in the Appendix. To develop an estimation method analogous to the HCK method note that from (2.3) and (2.11) we get (2.12) $E[N_r(\lambda)]E[\overline{X}_r(\lambda)] - m\Phi\left(-z_{\lambda} + \delta\right)c_{1r}(\delta, \lambda) = m_0\left[\lambda c_{0r}(\lambda) - \Phi\left(-z_{\lambda} + \delta\right)c_{1r}(\delta, \lambda)\right].$ For  $\lambda = p_{(i)}$ , i = 1, 2, ..., m, the term inside the square brackets in the R.H.S. of the above equation is  $x_{i} = p_{(i)}c_{0r}(p_{(i)}) - \Phi\left(-z_{p_{(i)}} + \delta\right)c_{1r}(\delta, p_{(i)})$ (2.13)and the L.H.S. can be estimated by  $y_i = i\overline{X}_r(p_{(i)}) - m\Phi\left(-z_{p_{(i)}} + \delta\right)c_{1r}(\delta, p_{(i)})$  $= \sum_{i=m-i+1}^{m} X_{(j)} - m\Phi\left(-z_{p_{(i)}} + \delta\right) c_{1r}(\delta, p_{(i)}).$ (2.14)Then from (2.12) we see that a regression line of  $y_i$  versus  $x_i$  can be fitted through the origin with slope  $m_0$  by minimizing  $\sum_{i=1}^{m} (y_i - m_0 x_i)^2$  w.r.t.  $m_0$  and  $\delta$ . The algorithm to solve this NLS regression problem is exactly analogous to the HCK algorithm.

З

2.3. EM Method

Whereas the HCK and TS methods compute the LS estimators of  $\pi_0$  and  $\delta$  (for two different regression models), the EM method computes their MLEs. For these MLEs to exist, it is necessary that  $\pi_0$  be bounded away from 0 and 1. The steps in the EM algorithm are as follows.

## EM Algorithm

**Step 0:** Compute initial estimates  $\hat{m}_0$  and  $\hat{\delta}$  by solving (2.7). Let  $\hat{\pi}_0 = \hat{m}_0/m$ . **Step 1 (E-step):** Calculate the posterior probabilities:

$$\widehat{\pi}_0(X_i) = \frac{\widehat{\pi}_0 \phi(X_i)}{\widehat{\alpha}_0 \phi(X_i) - \widehat{\alpha}_0 \phi(X_i)}$$

$$\widehat{\pi}_0(X_i) = \widehat{\pi}_0 \phi(X_i) + \widehat{\pi}_1 \phi(X_i - \widehat{\delta})$$
<sup>13</sup>
<sup>14</sup>

and  $\hat{\pi}_1(X_i) = 1 - \hat{\pi}_0(X_i), \ i = 1, 2, \dots, m.$ 

Step 2 (M-step): Calculate new estimates:

$$\widehat{\pi}_0 = \frac{\sum_{i=1}^m \widehat{\pi}_0(X_i)}{m} \text{ and } \widehat{\delta} = \frac{\sum_{i=1}^m \widehat{\pi}_1(X_i)X_i}{\sum_{i=1}^m \widehat{\pi}_1(X_i)}.$$

**Step 3:** Return to Step 1 until convergence.

## 3. Beta Model

In many applications the normal model may be inappropriate because the test statistics may not be normally distributed or different types of test statistics (e.g., normal, t, chi-square, Wlicoxon, log-rank) may be used to test different hypotheses or only the *p*-values of the test statistics may be available. In these cases we use the *p*-values to estimate  $\pi_0$ .

We propose to model the non-null *p*-values by a Beta(a, b) distribution given by

$$g(p|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{a-1} (1-p)^{b-1}$$

with unknown parameters a and b with a < 1 and b > 1. This restriction is imposed in order to ensure that g(p|a, b) is decreasing in p. It is well-known that the nonnull distribution of the p-values must be right-skewed and generally decreasing in shape (see Hung, O'Neill, Bauer and Kohne 1997). Langaas et al. (2005) imposed the same restriction in their nonparametric estimate of the non-null distribution of p-values.

Of course, the null distribution of a *p*-value is Beta(1, 1), i.e., the U[0, 1] distribution. As in the case of the normal model, the beta model can be represented as a mixture model for the distribution of the  $p_i$ :

45 (3.1) 
$$f(p_i) = \pi_0 \times 1 + \pi_1 g(p_i | a, b).$$

The parameters a and b must be estimated along with  $\pi_0$ . This problem is analogous to that encountered for the normal model with the difference that in addition to  $\pi_0$ , we have to estimate two parameters, a and b, instead of a single parameter  $\delta$ . We first extend the HCK method for the normal model discussed in Section 2.1 to this beta model. З

3.1. Hsueh, Chen, and Kodell (HCK) Method

Denote the type II error probability of a test performed at level  $\lambda$  by

(3.2) 
$$\beta(a,b,\lambda) = P_{H_{1i}}\{p_i > \lambda\} = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{\lambda}^{1} p^{a-1}(1-p)^{b-1}dp = 1 - I_{\lambda}(a,b),$$

where  $I_{\lambda}(a, b)$  is the incomplete beta function. Put

(3.3) 
$$x_i = p_{(i)} - I_{p_{(i)}}(a, b) \text{ and } y_i = i - mI_{p_{(i)}}(a, b).$$
 10

Then the HCK method amounts to solving the NLS problem of minimizing  $\sum_{i=1}^{m} (y_i - m_0 x_i)^2$  w.r.t.  $m_0$  and (a, b) (subject to a < 1 < b). Gauss-Newton method (Gill et al. 1981) was used to perform minimization w.r.t. (a, b). The initial starting values for this algorithm as well as the algorithms for the TS and EM estimators described below were determined by solving the following three moment equations for  $m_0$  and (a, b):

$$\sum_{i=1}^{m} p_i = \frac{1}{2}m_0 + \frac{a}{a+b}m_1,$$
<sup>18</sup>
<sup>19</sup>
<sub>20</sub>

<sup>21</sup>  
<sup>22</sup>  
<sup>23</sup> (3.4) 
$$\sum_{i=1}^{m} p_i^2 = \frac{1}{3}m_0 + \frac{a(a+1)}{(a+b)(a+b+1)}m_1,$$
<sup>21</sup>  
<sup>22</sup>  
<sup>23</sup> 23

$$\sum_{i=1}^{m} p_i^3 = \frac{1}{4}m_0 + \frac{a(a+1)(a+2)}{(a+b)(a+b+1)(a+b+2)}m_1.$$

### 3.2. Test Statistics (TS) Method

Here the TS estimator is based on the average of the "accepted" or "rejected" *p*-values defined as

$$\overline{p}_a(\lambda) = \frac{1}{N_a(\lambda)} \sum_{i \in S_a(\lambda)} p_i \text{ and } \overline{p}_r(\lambda) = \frac{1}{N_r(\lambda)} \sum_{i \in S_r(\lambda)} p_i.$$

Analogous to Lemma 1, we have the following lemma.

Lemma 3. Define

$$d_{0a}(\lambda) = E_{H_{0i}}(p_i | p_i > \lambda), d_{0r}(\lambda) = E_{H_{0i}}(p_i | p_i \le \lambda),$$

and

$$d_{1a}(a,b,\lambda) = E_{H_{1i}}(p_i | p_i > \lambda), d_{1r}(a,b,\lambda) = E_{H_{1i}}(p_i | p_i \le \lambda).$$

Then we have

$$d_{0a}(\lambda) = \frac{\lambda+1}{2}, d_{0r}(\lambda) = \frac{\lambda}{2}$$

and

$$d_{1a}(a,b,\lambda) = \frac{1-I_{\lambda}(a+1,b)}{1-I_{\lambda}(a,b)} \cdot \frac{a}{a+b}, d_{1r}(a,b,\lambda) = \frac{I_{\lambda}(a+1,b)}{I_{\lambda}(a,b)} \cdot \frac{a}{a+b}.$$

*Proof.* Straightforward. The next lemma gives  $E[\overline{p}_a(\lambda)]$  and  $E[\overline{p}_r(\lambda)]$ ; its proof is exactly analogous to that of Lemma 2. Lemma 4. Let  $g(\pi_0, a, b, \lambda) = P\{Z_i = 1 | p_i > \lambda\} = \frac{\pi_0(1 - \lambda)}{\pi_0(1 - \lambda) + \pi_1[1 - I_\lambda(a, b)]}$ and $h(\pi_0, a, b, \lambda) = P\left\{Z_i = 1 | p_i \le \lambda\right\} = \frac{\pi_0 \lambda}{\pi_0 \lambda + \pi_1 I_\lambda(a, b)}.$ Then  $E[\overline{p}_a(\lambda)] = q(\pi_0, a, b, \lambda)d_{0a}(\lambda) + [1 - q(\pi_0, a, b, \lambda)]d_{1a}(a, b, \lambda)$ (3.5)and  $E[\overline{p}_r(\lambda)] = h(\pi_0, a, b, \lambda) d_{0r}(\lambda) + [1 - h(\pi_0, a, b, \lambda)] d_{1r}(a, b, \lambda),$ (3.6)where  $d_{0a}(\lambda), d_{0r}(\lambda), d_{1a}(a, b, \lambda)$  and  $d_{1r}(a, b, \lambda)$  are as given in Lemma 3. The equations for the TS estimator are derived as follows. Analogous to (2.12), we obtain  $E[N_r(\lambda)]E[\overline{p}_r(\lambda)] - mI_{\lambda}(a,b)d_{1r}(a,b,\lambda) = m_0[\lambda d_{0r}(\lambda) - I_{\lambda}(a,b)d_{1r}(a,b,\lambda)].$ For  $\lambda = p_{(i)}$ , i = 1, 2, ..., m, the term in the square brackets in the R.H.S. of the above equation equals  $x_i = \frac{p_{(i)}^2}{2} - \frac{a}{a+b} I_{p_{(i)}}(a+1,b)$ and the L.H.S. can be estimated by

$$y_i = \sum_{j=1}^{a} p_{(j)} - \frac{a}{a+b} I_{p_{(i)}}(a+1,b).$$

The TS algorithm for the normal model can be modified to minimize  $\sum_{i=1}^{m} (y_i - m_0 x_i)^2$  by replacing the minimization with respect to  $\delta$  by minimization with respect to (a, b).

#### 3.3. EM Method

The steps in the EM algorithm, which gives the MLEs of  $\pi_0$  and (a, b), are as follows. As in the case of the normal model, for these MLEs to exist, it is necessary that  $\pi_0$  be bounded away from 0 and 1.

47 Step 0: Initialize  $\hat{m}_0$  and  $(\hat{a}, \hat{b})$  by solving (3.5). Let  $\hat{\pi}_0 = \hat{m}_0/m$ . 48 Step 1 (E-Step): Calculate the posterior probabilities: 

$$\widehat{\pi}_0(p_i) = \frac{\widehat{\pi}_0}{\widehat{\pi}_0 + \widehat{\pi}_1 g(p_i | \widehat{a}, \widehat{b})}$$

imsart-coll ver. 2008/08/29 file: Tamhane.tex date: April 10, 2009

and  $\hat{\pi}_1(p_i) = 1 - \hat{\pi}_0(p_i)$ , i = 1, 2, ..., m. **Step 2** (M-Step): Calculate  $\hat{a}$  and  $\hat{b}$  as solutions of the equations (see Equations (21.1) and (21.2) in Johnson and Kotz 1970):

$$\psi(a) - \psi(a+b) = \frac{\sum_{i=1}^{m} \hat{\pi}_1(p_i) \ln p_i}{\sum_{i=1}^{m} \hat{\pi}_1(p_i)},$$

$$\psi(b) - \psi(a+b) = \frac{\sum_{i=1}^{m} \widehat{\pi}_1(p_i) \ln(1-p_i)}{\sum_{i=1}^{m} \widehat{\pi}_1(p_i)},$$

where  $\psi(\cdot)$  is the digamma function (i.e., the derivative of the natural logarithm of the gamma function). Also calculate

$$\widehat{\pi}_0 = \frac{\sum_{i=1}^m \widehat{\pi}_0(p_i)}{m}.$$

**Step 3:** Return to Step 1 until convergence.

# 4. Adaptive Control of FDR

We now discuss the use of the estimate  $\hat{m}_0$  for adaptively controlling the FDR. The control is assumed to be strong control (Hochberg and Tamhane 1987), i.e., FDR  $\leq \alpha$  for some specified  $\alpha < 1$  for all possible combinations of true and false null hypotheses and the respective parameter values. Let R be the total number of rejected hypotheses and let V be the number of true hypotheses that are rejected. Benjamini and Hochberg (1995) introduced the definition

$$FDR = E\left[\frac{V}{R}\right] = E\left[\frac{V}{R}\middle| R > 0\right]P(R > 0),$$

where 0/0 is defined as 0. Benjamini and Hochberg (1995) gave a step-up (SU) procedure that controls FDR  $\leq \alpha$ .

Storey (2002) considered a single-step (SS) procedure (which he referred to as the fixed rejection region method) that rejects  $H_{0i}$  if  $p_i \leq \gamma$  for some common fixed threshold  $\gamma$ . His focus was on estimating the FDR. He proposed the following nonparametric estimator:

(4.1) 
$$\widehat{\mathrm{FDR}}_{\lambda}(\gamma) = \frac{\widehat{\pi}_0(\lambda)\gamma}{\{N_r(\lambda) \lor 1\}/m},$$
36  
37  
38

where  $\hat{\pi}_0(\lambda)$  is given by (1.2). The solution  $\hat{\gamma}$  to the equation  $\widehat{FDR}_{\lambda}(\gamma) = \alpha$  can be used in an MCP that tests each hypothesis at the  $\hat{\gamma}$ -level. Storey, Taylor and Siegmund (2004, Theorem 3) have shown that this heuristic procedure (which uses a slightly modified estimator of  $\pi_0$ ) controls the FDR. The heuristic procedures proposed below along the same lines (which use parametric estimators of the FDR) have not been rigorously shown to control the FDR.

We propose the following parametric estimator of the FDR:

(4.2) 
$$\widehat{\text{FDR}}(\gamma) = \frac{\widehat{\pi}_0 \gamma}{\widehat{\pi}_0 \gamma + \widehat{\pi}_1 [1 - \beta(\cdot, \gamma)]},$$
47
48

where  $\beta(\cdot, \gamma)$  is either  $\beta(\hat{\delta}, \gamma)$  given by (2.2) for the normal model or  $\beta(\hat{a}, \hat{b}, \gamma)$  given by (3.2) for the beta model. To adaptively control the FDR at level  $\alpha$ , we use 51

imsart-coll ver. 2008/08/29 file: Tamhane.tex date: April 10, 2009

the same heuristic procedure as above except that  $\hat{\gamma}$  is obtained by setting this parametric estimator equal to  $\alpha$ .

We may confine attention to  $\alpha \leq \pi_0$  since if  $\alpha > \pi_0$  then one can choose  $\hat{\gamma} = 1$ , and reject all hypotheses while still controlling the FDR =  $\pi_0 < \alpha$ . Existence and uniqueness of  $\hat{\gamma}$  for  $\alpha \in (0, \pi_0]$  is proved in the following two lemmas for the normal and beta models, respectively.

**Lemma 5.** For the normal model, the solution  $\widehat{\gamma}$  to the equation  $\widehat{FDR}(\gamma) = \alpha$ , where  $\widehat{FDR}(\gamma)$  and  $\beta(\widehat{\delta}, \gamma)$  are given by (4.2) and (2.2), respectively, exists and is unique for  $\alpha \in (0, \pi_0]$ .

*Proof.* Given in the Appendix.

**Lemma 6.** For the beta model, assuming  $0 < \hat{a} < 1 < \hat{b}$ , the solution  $\hat{\gamma}$  to the equation  $\widehat{FDR}(\gamma) = \alpha$ , where  $\widehat{FDR}(\gamma)$  and  $\beta(\hat{\delta}, \gamma)$  are given by (4.2) and (3.2), respectively, exists and is unique for  $\alpha \in (0, \pi_0]$ .

*Proof.* Given in the Appendix.

To develop an adaptive FDR-controlling procedure for the normal mixture model, Iyer and Sarkar (2007) took a somewhat different approach via the following asymptotic result of Genovese and Wasserman (2002): Assume that the  $p_i$  are independent U[0, 1] when the  $H_{0i}$  are true and have a common distribution F when the  $H_{0i}$  are false. Then the nominal  $\alpha$ -level Benjamini and Hochberg SU procedure is asymptotically (as  $m \to \infty$ ) equivalent to Storey's SS procedure that rejects  $H_{0i}$  if  $p_i \leq \hat{\gamma}$ where  $\hat{\gamma}$  is the solution to the equation

$$F(\gamma) = \rho \gamma \text{ and } \rho = \frac{1 - \alpha \pi_0}{\langle 1 - \alpha \rangle}.$$
 26

$$\alpha(1-\pi_0)$$
  
10 procedure actually controls the FDR conservatively at 29

Furthermore, since the SU procedure actually controls the FDR conservatively at  $\pi_0 \alpha$  level, exact control at level  $\alpha$  is achieved by replacing  $\alpha$  in the expression for  $\rho$  by  $\alpha/\pi_0$ , which results in the following equation for  $\gamma$ :

(4.3) 
$$F(\gamma) = \rho \gamma \text{ and } \rho = \frac{\pi_0(1-\alpha)}{\alpha(1-\pi_0)}.$$

By writing  $F(\gamma) = 1 - \beta(\cdot, \gamma)$ , we see that  $\widehat{\text{FDR}}(\gamma) = \alpha$  and (4.3) are identical if  $\pi_0$  is replaced by  $\widehat{\pi}_0$  in (4.3). If and Sarkar (2007) used the solution  $\widehat{\gamma}$  from (4.3) in Storey's SS procedure with  $F(\gamma) = \Phi(\delta - z_{\gamma})$ , and  $\delta$  and  $\pi_0$  replaced by their estimates  $\widehat{\delta}$  and  $\widehat{\pi}_0$  obtained from the EM method, which results in an adaptive FDR-controlling procedure, which is identical to the one proposed before.

#### 

#### 5. Simulation Results

We compared different estimators by conducting an extensive simulation study. The ST estimator was used with  $\lambda = 0.5$  throughout. The estimators were compared in terms of their accuracy of estimation of  $\pi_0$  and control of FDR at a nominal  $\alpha = 0.10$  using the SS procedure. The bias and mean square error (MSE) of the estimators were used as the performance measures. The results for the normal model are presented in Section 5.1 and for the beta model in Section 5.2. Throughout we used m = 1000 and the number of replications was also set equal to 1000. We varied

 $\pi_0$  from 0.1 to 0.9 in steps of 0.1. The values  $\pi_0 = 0$  and 1 were excluded because  $\hat{\pi}_0$  exhibits erratic results in these extreme cases; also FDR = 0 when  $\pi_0 = 0$ . In each simulation run, first the random indexes of the true and false hypotheses were generated by generating Bernoulli r.v.'s  $Z_i$ . Then the respective  $X_i$ 's or the  $p_i$ 's were generated using the appropriate null or alternative distributions. The bias of each  $\hat{\pi}_0$  estimator was estimated as the difference between the average of the  $\hat{\pi}_0$  values over 1000 replicates and the true value of  $\pi_0$ . In the case of FDR, the bias was estimated as the difference between the FDR values over 1000 replicates and the nominal  $\alpha = 0.10$ . The MSE was computed as the sum of the square of the bias and the variance of the  $\hat{\pi}_0$  (or FDR) values averaged over 1000 replicates. The detailed numerical results are given in Shi (2006); here we only present graphical plots to save space.

#### 5.1. Simulation Results for Normal Model

З

Simulations were conducted in three parts. In the first part, the true model for the non-null hypotheses was set to be the same as the assumed model by generating the  $X_i$ 's for the false hypotheses from a  $N(\delta, \sigma^2)$  distribution with a fixed  $\delta = 2$  and  $\sigma = 1$ . In the other two parts of simulations, robustness of the assumed model was tested by generating the  $X_i$ 's for the false hypotheses from different distributions than the assumed one. In the second part, the  $X_i$ 's for the false hypotheses were generated from  $N(\delta_i, \sigma^2)$  distributions where the  $\delta_i$ 's were themselves drawn from a  $N(\delta_0, \sigma_0^2)$  distribution with  $\delta_0 = 2$  and  $\sigma_0 = 0.25$  corresponding to an approximate range of [1,3] for the  $\delta_i$ . In the third part, the  $p_i$ 's for the false hypotheses were generated from a Beta(a, b) distribution with a = 0.5 and b = 2, and the  $X_i$ 's were computed using the inverse normal transformation  $X_i = \Phi^{-1}(1 - p_i)$ .

**Results for Fixed**  $\delta$ : The bias and the square root of the mean square error  $(\sqrt{\text{MSE}})$  of  $\hat{\pi}_0$  for ST, HCK, TS and EM estimators are plotted in Fig. 1. Note from Equation (2.3) that the bias of the ST estimator is given by

(5.1) 
$$\operatorname{Bias}[\widehat{\pi}_0(\lambda)] = \frac{1 - \pi_0}{1 - \lambda} \Phi(z_{\lambda} - \delta).$$

Also, using the fact that  $N_a(\lambda)$  has a binomial distribution with number of trials m and success probability,

 $p = P\{p_i > \lambda\} = \pi_0(1-\lambda) + (1-\pi_0)\Phi(z_\lambda - \delta),$ 

and using Equation (1.2) for  $\hat{\pi}_0(\lambda)$ , we have

(5.2) 
$$\operatorname{Var}[\hat{\pi}_0(\lambda)] = \frac{p(1-p)}{m(1-\lambda)^2}.$$
 41  
42

These formulae were used to compute the bias and MSE of the ST estimator instead of estimating them by simulation. Note that the bias of the ST estimator decreases linearly in  $\pi_0$ . The  $\sqrt{\text{MSE}}$  plot for the ST estimator is also approximately linear because the bias is the dominant term in MSE. This is true whenever the alternative is fixed for all false null hypotheses.

<sup>49</sup> The TS estimator does not offer an improvement over the HCK estimator, as we <sup>50</sup> had hoped, and in fact performs slightly worse in terms of MSE for  $\pi_0 \leq 0.5$ . We sus-<sup>51</sup> pect that this result is due to the bias introduced when the term  $E[N_r(\lambda)]E[\overline{X}_r(\lambda)]$  <sup>51</sup>

#### Tamhane and Shi

in Equation (2.12) is estimated by  $i\overline{X}_r(p_{(i)})$  for  $\lambda = p_{(i)}$  because of the fact that the product of the expected values does not equal the expected value of the product of two dependent r.v.'s. The EM estimator has consistently the lowest bias and the lowest MSE.



FIG 1. Bias and  $\sqrt{MSE}$  of  $\hat{\pi}_0$  for ST, HCK, TS and EM Estimators for Normal Model (Data Generated by Normal Model with Fixed  $\delta$ )

The bias and  $\sqrt{\text{MSE}}$  of  $\widehat{\text{FDR}}$  for ST, HCK, TS and EM estimators are plotted in Fig. 2. We see that the ST estimator leads to a large negative bias which means that, on the average,  $\widehat{\text{FDR}}$  is less than the nominal  $\alpha = 0.10$  resulting in conservative control of FDR. The other three estimators yield approximately the same level of control. Surprisingly, there is very little difference in the MSEs of the four estimators. The EM estimator still is the best choice with the lowest bias and the lowest MSE throughout the entire range of  $\pi_0$  values.

36
37
38
39
40
41
42
43
44
45
46
47
48
49



FIG 2. Bias and  $\sqrt{MSE}$  of  $\widehat{FDR}$  for ST, HCK, TS and EM Estimators for Normal Model (Data Generated by Normal Model with Fixed  $\delta$ )

**Results for Random**  $\delta$ : The bias and  $\sqrt{\text{MSE}}$  of  $\hat{\pi}_0$  and of FDR for ST, HCK, TS and EM estimators are plotted in Figs. 3 and 4, respectively. By comparing these results with those for fixed  $\delta = 2$ , we see that, as one would expect, there is a slight degradation in the performance of every estimator because the assumed model does not hold. The comparisons between the four estimators here are similar to those for fixed  $\delta$  with the estimators ranked as EM > HCK > TS > ST.



FIG 3. Bias and  $\sqrt{MSE}$  of  $\hat{\pi}_0$  for ST, HCK, TS and EM Estimators for Normal Model (Data Generated by Normal Model with Random  $\delta$ )



FIG 4. Bias and  $\sqrt{MSE}$  of  $\widehat{FDR}$  for ST, HCK, TS and EM Estimators for Normal Model (Data Generated by Normal Model with Random  $\delta$ )

**Robustness Results for Data Generated by Beta Model:** The bias and  $\sqrt{\text{MSE}}$  of  $\hat{\pi}_0$  and of  $\widehat{\text{FDR}}$  for ST, HCK, TS and EM estimators are plotted in Figs. 5 and 6, respectively. Looking at Fig. 5 first, we see that the biases and MSEs of all four estimators are an order of magnitude higher compared to the normal model data which reflects lack of robustness.

Tamhane and Shi

З

It is interesting to note that the EM estimator is no longer uniformly best for estimating  $\pi_0$ . In fact, the HCK estimator has a lower bias and MSE for  $0.2 \le \pi_0 \le 0.7$ . The lack of robustness of the EM estimator is likely due to the strong dependence of the likelihood methods on distributional assumptions. On the other hand, for the least squares methods, the dependence on the assumed distribution is only through its first moment and hence is less strong. As far as control of FDR is concerned, there are not large differences between the proposed estimators. However, when  $\pi_0 = 0.9$  the proposed estimators exceed the nominal FDR by as much as 0.05, while the ST estimator still controls FDR conservatively. In conclusion, the HCK estimator performs best for the middle range of  $\pi_0$  values.



FIG 5. Bias and  $\sqrt{MSE}$  of  $\hat{\pi}_0$  for ST, HCK, TS and EM Estimators for Normal Model (Data Generated by Beta Model)



Generated by Beta Model)

# 5.2. Simulation Results for Beta Model

**Results for Beta**(0.5, 2) **Data:** In this case the non-null *p*-values were generated from a Beta(a, b) distribution with a = 0.5, b = 2.0 and the null p-values were generated from the U[0,1] distribution. As before, the bias and variance of the ST estimator were not estimated from simulations, but were computed using Equations (5.1) and (5.2) with  $\Phi(z_{\lambda} - \delta)$  replaced by  $1 - I_{\lambda}(a, b)$ . Note that the bias of the ST estimator decreases linearly in  $\pi_0$  in this case as well and  $\sqrt{\text{MSE}}$  decreases approximately linearly. From Fig. 7 we see that all estimators of  $\pi_0$ , except ST, have significant negative biases particularly over the interval [0.2, 0.5] and for  $\pi_0 \geq 0.7$ , resulting in the achieved FDR significantly exceeding the nominal value of  $\alpha = 0.10$  over the corresponding ranges of  $\pi_0$  as can be seen from Fig. 8. Comparing the results here with those for the normal model with the fixed  $\delta$  case, we see that the biases and MSEs of all estimators are an order of magnitude higher in the present case. The reason behind this poor performance of the beta model probably lies in the difficulty of estimating the parameters a, b of the beta distribution. Only the ST estimator controls FDR conservatively and has the smallest MSE for  $0.2 \leq \pi_0 \leq 0.7$ . Thus the ST estimator has the best performance since it is a nonparametric estimator (and the performance would be even better if  $\lambda$  is not fixed, but is used as a tuning parameter). In other words, the benefits of using a parametric model are far outweighed by the difficulty of estimating the parameters of the model resulting in less efficient estimators.



FIG 7. Bias and  $\sqrt{MSE}$  of  $\hat{\pi}_0$  for ST, HCK, TS and EM Estimators for Beta Model (Data Generated by Beta Model)

Robustness Results for Data Generated by Normal Model: In this case we generated the data by the normal model with  $N(2, 1^2)$  as the alternative distribution. The *p*-values were then computed and all four methods of estimation were applied. The results are plotted in Figs. 9 and 10. From these figures we see that none of the proposed estimators exhibit consistent negative bias as they did when the data were generated according to the beta model. This is somewhat surprising since one would expect these estimators to perform more poorly when the assumed model does not hold as in the present case. We also see that the EM estimator performs worse than other estimators. Thus lack of robustness of the EM estimator

imsart-coll ver. 2008/08/29 file: Tamhane.tex date: April 10, 2009

З

З



FIG 8. Bias and  $\sqrt{MSE}$  of  $\widehat{FDR}$  for ST, HCK, TS and EM Estimators for Beta Model (Data Generated by Beta Model)

to the model assumptions is demonstrated again, and for the same reason. The TS estimator generally has the lowest bias for estimating  $\pi_0$  and its achieved FDR is closest to the nominal  $\alpha$ ; the ST estimator has the second best performance.



FIG 9. Bias and  $\sqrt{MSE}$  of  $\hat{\pi}_0$  for ST, HCK, TS and EM Estimators for Beta Model (Data Generated by Normal Model)

# 6. Example

We consider the National Assessment of Educational Progress (NAEP) data analyzed by Benjamini and Hochberg (2000). The data pertain to the changes in the
average eighth-grade mathematics achievement scores for the 34 states that participated in both the 1990 and 1992 NAEP Trial State Assessment. The raw *p*-values
for the 34 states are listed in the increasing order in Table 2. The FWE controlling



FIG 10. Bias and  $\sqrt{MSE}$  of  $\widehat{FDR}$  for ST, HCK, TS and EM Estimators for Beta Model (Data Generated by Normal Model)

Bonferroni procedure and the Hochberg (1988) procedure both identified only 4 significant results (those with *p*-values  $\leq p_{(4)} = 0.0002$ ) Application of the FDR controlling non-adaptive Benjamini-Hochberg SU procedure resulted in 11 significant results. By applying their method they estimated  $\hat{m}_0 = 7$  ( $\hat{\pi}_0 = 0.2059$ ); using this value in the adaptive version of their procedure yielded 24 significant results.

We applied the three methods of estimation considered in this paper to these data under both the normal and beta models. The estimates  $\hat{\pi}_0$  and the associated  $\hat{\delta}$  or  $(\hat{a}, \hat{b})$  values are given in Table 1. We see that for both models, the HCK and EM methods give smaller estimates of  $\pi_0$  than does the TS method. The  $\hat{\gamma}$ -values obtained by solving the equation  $\widehat{\text{FDR}}(\gamma) = \alpha$  for  $\alpha = 0.05$  are inversely ordered.



|                    | No     | ormal Mo | del                    | Beta Model |        |                        |  |  |
|--------------------|--------|----------|------------------------|------------|--------|------------------------|--|--|
|                    | HCK    | TS       | $\mathbf{E}\mathbf{M}$ | HCK        | TS     | $\mathbf{E}\mathbf{M}$ |  |  |
| $\widehat{\pi}_0$  | 0.1317 | 0.3233   | 0.1407                 | 0.0096     | 0.1307 | 0.0160                 |  |  |
| $\widehat{\gamma}$ | 0.3163 | 0.0918   | 0.2946                 | 1.0000     | 0.3092 | 1.0000                 |  |  |
| $\widehat{\delta}$ | 1.8285 | 2.2657   | 1.9221                 | -          | -      | -                      |  |  |
| $\widehat{a}$      | -      | —        | -                      | 0.3291     | 0.4474 | 0.3210                 |  |  |
| $\widehat{b}$      | _      | -        | -                      | 2.0764     | 3.2842 | 1.9313                 |  |  |
| $N_r$              | 28     | 21       | 27                     | 34         | 27     | 34                     |  |  |

|  | $N_r =$ | Number | of rejected | hypotheses |
|--|---------|--------|-------------|------------|
|--|---------|--------|-------------|------------|

The *p*-values  $\leq \hat{\gamma}$  are declared significant. From Table 2, we see that the number of significant *p*-values for HCK, TS and EM for the normal model are 28, 21 and 27, respectively. Thus, HCK and EM methods give more rejections than Benjamini and Hochberg's (2000) adaptive SU procedure.

Before fitting the beta mixture model, it is useful to plot a histogram of the *p*-values. This histogram is shown in Fig. 11. It has a decreasing shape, and assuming that the majority of the *p*-values are non-null, it corresponds to a < 1 and b > 1.

З



FIG 11. Histogram of the p-Values for the NAEP Data

HCK and EM methods yield  $\hat{\pi}_0 < \alpha = 0.05$ , hence  $\hat{\gamma} = 1$  which means that all 34 hypotheses are rejected. This evidently liberal result is likely due to underestimation of  $\pi_0$  using the beta model as noted in Section 5.2. The TS method yields  $\hat{\pi}_0 = 0.1307$  and  $\hat{\gamma} = 0.3092$ , which are close to the estimates produced by the HCK and EM methods for the normal model and it rejects the same 27 hypotheses.

Rejections of hypotheses with large *p*-values will justifiably raise many eyebrows. This appears to be a problem with FDR-controlling procedures when there are many hypotheses that are clearly false (with *p*-values close to zero) which lowers the bar for rejection for other hypotheses. Shaffer (2005) has discussed this problem and has suggested imposing additional error controlling requirements in order to limit such dubious rejections. This is a topic for further research.

# 7. Concluding Remarks

In this paper we offered two different mixture models for estimating the number of true null hypotheses by modeling the non-null *p*-values. For each model (the normal and beta), three methods of estimation were developed: HCK, TS and EM. Generally speaking, these parametric estimators outperform (in terms of the accuracy of the estimate of  $\pi_0$  and control of the FDR) the nonparametric ST estimator for the normal model but not for the beta model. The reason for this is that the normal model is easier to estimate and so the benefits of the parametric estimators are not significantly compromised by the errors of estimation. On the other hand, the beta model is difficult to estimate and so the benefits of the parametric estimators are lost. Therefore we do not recommend the use of the beta model in practice. 

For normally distributed test statistics, the EM estimator generally performs best
followed by the HCK and TS estimators. However, the EM estimator is not robust
to the violation of the model assumptions. If the EM estimator for the normal model
is applied to the data generated from the beta model or vice versa, its performance
is often worse than that of the HCK estimator, and sometimes even that of the ST

З

| C.     | toto | n roluo         | UCK      | TC  | БМ         | State | n voluo         | UCK      | TC  | БМ       |
|--------|------|-----------------|----------|-----|------------|-------|-----------------|----------|-----|----------|
| د<br>٦ | tate | <i>p</i> -value | пск<br>* | 15  | E/IVI<br>* | State | <i>p</i> -value | пск<br>* | 15  | E/IVI    |
| R      |      | 0.00000         | -1-<br>- | -1- | -1-<br>-   | NY    | 0.05802         |          | -1- | -1-      |
| N      | 1N   | 0.00002         | , î      | *   | *          | OH    | 0.06590         | ^        | *   | <b>*</b> |
| H      | n    | 0.00002         | *        | *   | *          | CA    | 0.07912         | *        | *   | *        |
| N      | IC   | 0.00002         | *        | *   | *          | MD    | 0.08226         | *        | *   | *        |
| N      | IH   | 0.00180         | *        | *   | *          | WV    | 0.10026         | *        |     | *        |
| IA     | A    | 0.00200         | *        | *   | *          | VA    | 0.14374         | *        |     | *        |
| C      | o    | 0.00282         | *        | *   | *          | WI    | 0.15872         | *        |     | *        |
| Т      | X    | 0.00404         | *        | *   | *          | IN    | 0.19388         | *        |     | *        |
| II     | D    | 0.00748         | *        | *   | *          | LA    | 0.20964         | *        |     | *        |
| А      | Z    | 0.00904         | *        | *   | *          | MI    | 0.23522         | *        |     | *        |
| K      | XY   | 0.00964         | *        | *   | *          | DE    | 0.31162         | *        |     |          |
| 0      | )K   | 0.02036         | *        | *   | *          | ND    | 0.36890         |          |     |          |
| C      | т    | 0.04104         | *        | *   | *          | NE    | 0.38640         |          |     |          |
| Ν      | IM   | 0.04650         | *        | *   | *          | NJ    | 0.41998         |          |     |          |
| W      | VY   | 0.04678         | *        | *   | *          | AL    | 0.44008         |          |     |          |
| F      | 'L   | 0.05490         | *        | *   | *          | AR    | 0.60282         |          |     |          |
| P      | PA   | 0.05572         | *        | *   | *          | GA    | 0.85628         |          |     |          |



estimator. The TS estimator did not improve on the HCK estimator in all cases as we had hoped. Thus our final recommendation is to use the normal model with the EM method if the test statistics follow approximately normal distributions and the HCK method otherwise. If only the p-values calculated from various types of test statistics are available then the ST method is recommended; alternatively the p-values may be transformed using the inverse normal transform and then the HCK method may be applied.

#### Acknowledgment

This research was partially supported by the National Heart, Lung and Blood Institute Grant 1 R01 HL082725-01A1 and the National Security Agency Grant H98230-07-1-0068. The authors are extremely grateful especially to one of the two referees who pointed out some crucial errors in the earlier version of the paper.

#### Appendix A

*Proof of Lemma 2.* We have

$$E[\overline{X}_{a}(\lambda)] = E\left\{\frac{1}{N_{a}}\sum_{i\in S_{a}(\lambda)}X_{i}\right\}$$

$$= E\left\{ E\left[ \left. \frac{1}{n_a} \sum_{i \in s_a} X_i \right| S_a(\lambda) = s_a, N_a(\lambda) = n_a \right] \right\}$$

 $g(\pi_0, \delta, \lambda)c_{0a}(\lambda) + [1 - g(\pi_0, \delta, \lambda)]c_{1a}(\delta, \lambda).$ 

$$= E\left\{\frac{1}{n_a} \cdot n_a[g(\pi_0, \delta, \lambda)c_{0a}(\lambda) + [1 - g(\pi_0, \delta, \lambda)]c_{1a}(\delta, \lambda)]\right\}$$

З

Tamhane and Shi

In the penultimate step above, we have used the fact that conditionally on  $X_i \leq z_{\lambda}$ , the probability that  $Z_i = 1$  is  $g(\pi_0, \delta, \lambda)$  and the probability that  $Z_i = 0$  is  $1 - g(\pi_0, \delta, \lambda)$ . Furthermore, the conditional expectation of  $X_i$  in the first case is  $c_{0a}(\lambda)$  and in the second case it is  $c_{1a}(\delta,\lambda)$ . The expression for  $E[\overline{X}_r(\lambda)]$  follows similarly.

*Proof of Lemma 5.* By substituting for  $\beta(\cdot, \gamma)$  from (2.2) and dropping carets on  $\widehat{\mathrm{FDR}}(\gamma), \widehat{\pi}_0, \widehat{\pi}_1$  and  $\widehat{\delta}$  for notational convenience, the equation to be solved is

$$FDR(\gamma) = \frac{\pi_0}{\pi_0 + \pi_1 \Phi(\delta - z\gamma)/\gamma} = \alpha.$$

It is easy to check that FDR(0) = 0 and  $FDR(1) = \pi_0$ . We shall show that  $FDR(\gamma)$ is an increasing function of  $\gamma$  which will prove the lemma. Thus we need to show that  $u(\delta, \gamma) = \Phi(\delta - z_{\gamma})/\gamma$  is decreasing in  $\gamma$ . By implicit differentiation of the equation  $\Phi(z_{\gamma}) = 1 - \gamma$ , we get

$$\frac{dz\gamma}{d\gamma} = -\frac{1}{\phi(z\gamma)}.$$
<sup>16</sup>
<sup>17</sup>
<sup>18</sup>
<sup>19</sup>

Hence,

$$\frac{du(\delta,\gamma)}{d\gamma} = \frac{\gamma\phi(\delta-z\gamma) - \phi(z\gamma)\Phi(\delta-z\gamma)}{\gamma^2\phi(z\gamma)}.$$

Therefore we need to show that

$$v(\delta,\gamma) = \phi(z\gamma) \Phi(\delta-z\gamma) - \gamma \phi(\delta-z\gamma) > 0 \ \forall \ \delta > 0.$$

Now  $v(0, \gamma) = 0$ . Therefore we must show that

$$\frac{dv(\delta,\gamma)}{d\delta} = \phi(\delta - z\gamma)[\phi(z\gamma) + \gamma(\delta - z\gamma)] > 0,$$

which reduces to the condition:  $w(\delta, \gamma) = \phi(z_{\gamma}) + \gamma(\delta - z_{\gamma}) > 0$ . Since  $w(\delta, \gamma)$  is increasing in  $\delta$ , it suffices to show that

$$w(0,\gamma) = \phi(z\gamma) - \gamma z\gamma > 0.$$

By putting  $x = z_{\gamma}$  and hence  $\gamma = \Phi(-x)$  the above inequality becomes

which is the Mills' ratio inequality (Johnson and Kotz 1970, p. 279). This completes the proof of the lemma.

*Proof of Lemma* 6. By substituting for  $\beta(\cdot, \gamma)$  from (3.2) and dropping carets on  $\widehat{\mathrm{FDR}}(\gamma), \widehat{\pi}_0, \widehat{\pi}_1, \widehat{a}$  and  $\widehat{b}$  for notational convenience, the equation to be solved is

(A.1) 
$$FDR(\gamma) = \frac{\pi_0}{\pi_0 + \pi_1 I_{\gamma}(a, b)/\gamma} = \alpha.$$
 43  
44  
45

Note that FDR(0) = 0 and  $FDR(1) = \pi_0$ . To show that  $FDR(\gamma)$  is an increasing function of  $\gamma$  we need to show that  $I_{\gamma}(a,b)/\gamma$  decreases in  $\gamma$ . To see this, note that the derivative of  $I_{\gamma}(a,b)/\gamma$  w.r.t.  $\gamma$  is proportional to  $\gamma g(\gamma | a, b) - I_{\gamma}(a, b)$ , which is negative since the beta p.d.f.  $g(\gamma|a, b)$  is strictly decreasing in  $\gamma$  for a < 1 and b > 1, and so  $\gamma q(\gamma | a, b) < I_{\gamma}(a, b)$ . It follows therefore that the equation FDR( $\gamma$ ) =  $\alpha$  has a unique solution in  $\gamma \in (0, 1)$  for  $\alpha \in (0, \pi_0]$ . 

| <ol> <li>BENJAMIN, Y. and HOCHBERG, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. <i>Journal of the Royal Statistical Society, Scr. B</i>, 57, 289-300.</li> <li>BENJAMIN, Y. and HOCHBERG, Y. (2000). On the adaptive control of Like discovery rates in multiple testing with independent statistics. <i>Journal of Educational Statistica</i>, 28, 00-83.</li> <li>BEAGE, M. &amp; (2004). A nucleon strend of the statistics of the testing of the testing of the statistics. <i>Journal of Royal Statistical Society, Scr. B</i>, 30, 1-38.</li> <li>FINNER, H. and ROTERS, M. (2001). On the false discovery rate and expected type I error. <i>Biometrical Journal</i>, 8, 985-1005.</li> <li>GENVERE, C. and WASEBMAN, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. <i>Journal of Royal Statistical Society, Scr. B</i>, 64, 490-517.</li> <li>Chu, P. E., MCRAN, W. and WURT, M. H. (1981). <i>Practical Optimization</i>. Academic Press, Londen and New York.</li> <li>GLAZ, Z. WU, B. and ZMAO, H. (2004). Model-based approach to FDR estimation. <i>Tech. Report 2004-016</i>. Division of Biostatistics, Univ. of Minnesota, Minnespoli, M.N.</li> <li>HOCHBRA, TS, 800-803.</li> <li>HOCHBRA, TG, AND AWARNA, A. C. (1987). Multiple Comparison Procedures John Wiley and Stons. New York.</li> <li>HISTEH, H., CHES, J. J. and KODEL, R. L. (2003). Comparison of methods for estimating the number of true null hypotheses in multiple testing. <i>Journal of Biophymacetal Statistics</i>, 30, 675-689.</li> <li>HUNG, H. M., O'NELL, R. T., BAUER, P. and KONE</li></ol>   | 1  | $\mathbf{Ref}$ | erences  | 1  |
|--|----|----------------|--|----|
| <ol> <li>DEALMAN, Y. and ROUMBAR, Y. (1996). Controlling the last statistical Society, Soc. B, 57, 289–4</li> <li>BERLAND, Y. and HOURBER, Y. (2000). On the adaptive control of the false discovery rate in multiple testing with independent statistical. Journal of Educational Statistica, 50, 60-83.</li> <li>DEALE, M. A. (2004). A note on the adaptive control of false discovery rates. Journal of Royal Statistical Society, Soc. B, 69, 69, 87-804.</li> <li>DEMENTER, A. P., LARD, N. M. and RURN, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. Journal of Royal Statistical Society, Soc. B, 98, 1-93.</li> <li>PENNER, H. and ROTERS, M. (2001). On the false discovery rate and expected type I error. Biometric of Journal 8, 985-1005.</li> <li>GENOTESE, C. and WASSERMAN, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. Journal of Royal Statistical Society, Soc. B, 94, 1-99-517.</li> <li>GLIL, P. E., MURKAY, W. and WRUIT, M. H. (1981). Practical Optimization. Academic Press, London and New York.</li> <li>GUAN, Z., WU, B. and ZIMO, H. (2004). Model-based approach to FDR estimation. Tech. Report Bg004-016. Division of Biostatistics, Univ of Minesota, Minnegolis, MN.</li> <li>HOCHIBER, Y. (1988). A sharper Bonferoni procedure for multiple tests of significance Biometryka, 75, 809-803.</li> <li>HOCHIBER, Y. (1988). A sharper Bonferoni procedure for multiple testing. Statistics in Medicine, 9, 811-813.</li> <li>HOCHIBER, Y. (1989). An alter statistic Journal of Biopharmacentical Statistics, 13, 677-689.</li> <li>HOCHIBER, Y. (1990). More powerful procedures for multiple significance testing. Statistics in Medicine, 9, 811-813.</li> <li>HOCHIBER, Y. (1990). An adaptive single-step FDR procedures with applications to DNA microarray analysis. Biometrical Journal of Biopharmacentical Statistics, 13, 677-689.</li> <li>HOCHIBER, T. and Spirotycolucy. E. (1997). The behavior of the rue null hypotheses for multiple</li></ol>   | 2  | [1]            | Drugger V and Hearppea V (1007) Claster line the file line and A spectral and  | 2  |
| <ol> <li>DENAMEN, Y. and HOCHEREG, Y. (2000). On the adaptive control of the false discovery rate in multiple testing with independent statistics. Journal of Educational Statistics, 25, 60–83.</li> <li>BLACK, M. A. (2001). A note on the adaptive control of false discovery rates. Journal of Royal Statistical Society, Ser. B, 60, 297–304.</li> <li>DENETER, A. P., LARIN, N. M. and RUBN, D. E. (1977). Maximum likelihood from incomplete data via the EM algorithm. Journal of Royal Statistical Society, Ser. B, 39, 1–38.</li> <li>FINNAR, H. and ROTRES, M. (2001). On the false discovery rate and expected type I error. Biometric real Journal, 6, 985-1005.</li> <li>GENOTES, C. ANDREN, W. and WAGRET, M. H. (1981). Fractical Optimization. Academic Press, 13</li> <li>GIOM, Z., WU, B. and Ziako, H. (2004). Model-based approach to FDR estimation. Tech. Report 14 (2004). 400-407. Distain of Biotastistics, Univ of Minesoda, Minnegrolis, MN.</li> <li>HOCHERG, Y. (1988). A sharper Bonferroni procedure for multiple tests of significance testing. Statistics in Medicine, 9, 811–818.</li> <li>HOCHERG, Y. and DENAMEN, Y. (1990). More powerful procedures for suthiple significance testing. Statistics in Medicine, 9, 811–818.</li> <li>HOCHERG, Y. and ENAMEN, Y. (1990). More powerful procedures for suthiple significance testing. Statistics in Medicine, 9, 811–818.</li> <li>HOCHERG, Y. and ENAMEN, Y. (1990). More powerful procedures of nethods for estimating the number of true null hypothesis in untriple testing. Journal of Biopharmotor of the p-value when the alternative hypothesis in Mintple Journal, 91, 21–718.</li> <li>HUNG, H. M., O'NELL, R. T., EXCER, P. and KOREK, K. (1997). The behavior of the p-value when the alternative hypothesis in control of Devian Janual, 91, 21–718.</li> <li>HUNG, H. M., ONDEL, R. W. (2003). Continuous Univariate Distributions, 1. John Wiley &amp; Songer Proceedures, Statistical Society, Sci. 76, 90, 197–108.</li> <li>HUNG, H. M., ONDELE, R. W. (2003). Continuou</li></ol>   | 3  | [1]            | Denjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: A practical and<br>powerful approach to multiple testing. <i>Journal of the Royal Statistical Society. Ser. B.</i> 57, 289–         | 3  |
| <ol> <li>EEMARDS, Y. and HOGIBEG, Y. (2000). On the adaptive control of the false discovery rate in 6 false discovery rates. Journal of Royal Statistical Society, Ser. B, 69, 207–304.</li> <li>DEMEYER, A. P., LAND, N. M. and RUBN, D. B. (1977). Maximum likelihood from incomplete data via the FM algorithm. Journal of Royal Statistical Society, Ser. B, 39, 1–38.</li> <li>FINNER, H. and ROTERS, M. (2001). On the false discovery rate and expected type I error. Biometrical Journal, 8, 985-1005.</li> <li>GENONER, C. and WASEBLAN, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. Journal of Royal Statistical Society, Ser. B, 64, 499-517.</li> <li>GEL, F. E., MURAN, W. and WHART, M. H. (1981). Practical Optimization. Academic Press, 120, 2007.</li> <li>GEL, P. E., MURAN, W. and WHART, M. H. (1981). Practical Optimization. Academic Press, 120, 2007.</li> <li>GONZ, Z., WY, D. 202, 2010. H. (2004). Model-based approach to TDR estimation. Tech. Report 42, 2007. 2007.</li> <li>GONZ, Z., WY, U1985). A sharper Bonferroni procedure for multiple tests of significance. 16 Journator, 70, 1985). A sharper Bonferroni procedure for multiple statistics, 136, 767-689.</li> <li>Horcumator, Y. and TAMIANG, A. C. (1987). Multiple Comparison Procedures. John Wiley and Sons, New York.</li> <li>Hersten, H., CURNI, J. J. and KONEL, R. L. (2003). Comparison of methods for estimating the number of true null hypotheses in multiple testing. Journal of Bopharmaceutical Statistics, 13, 67-689.</li> <li>Horcumator, Y. and TAMIANG, A. C. (1987). Multiple Comparison of the bryalue when the alternative hypothesis is true. Biometrics, 53, 11-22.</li> <li>Hynt, V. and Status, S. (2007). Continuous Univariate Distributions, I. John Wiley and Sons, New York.</li> <li>Hatas, H., CURNI, R. T., Batter, P. and KONEK, K. (1997). The behavior of the p-value when the alternative hypothesis is true. Biometrics, 53, 11-22.</li> <li>Hynt, V. and Status, S. (2007). An ad</li></ol>   | 4  |                | 300.   | 4  |
| <ul> <li>a BLACK, M. A. (2004). A note on the adaptive control of fast discovery rates. Journal of Royal Statistical Society, Ser. B, 66, 297–304.</li> <li>DEARM, M. M. (2004). A note on the adaptive control of fast discovery rates. Journal of Royal Statistical Society, Ser. B, 80, 1–38.</li> <li>Prexan, H. and FURES, M. (2001). On the fast discovery rate and expected type I error. Diometer resource of the fast of Royal Statistical Society, Ser. B, 80, 1–38.</li> <li>Frexan, H. and FURES, M. (2001). On the fast discovery rate and expected type I error. Diometer resource of the fast of Royal Statistical Society, Ser. B, 64, 09–517.</li> <li>Giux, P. E., MURRAY, W. and WHGHY, M. H. (1981). Practical Optimization. Academic Press, Landon and New York.</li> <li>Giux, Z. M. W. B. and Zuko, H. (2004). Model-based approach to FDR semimation. Tech. Report 42, 2004-016. Division of Biostatistics, Univ. of Minesota, Minnesota, Min</li></ul>   | 5  | [2]            | BENJAMINI, Y. and HOCHBERG, Y. (2000). On the adaptive control of the false discovery rate in  | 5  |
| <ul> <li>Statistical Society, Ser. B, 66, 207–304.</li> <li>Deversen, A. P., LARD, N. M. and RUBN, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. <i>Journal of Royal Statistical Society, Ser. B</i>, 39, 1–38.</li> <li>Persen, H. and Korzen, M. (2002). On the false discovery rate and expected type I error. <i>Biometrical Journal</i>, 8, 985-1005.</li> <li>Gesoverse, C. and Wassenava, N. L. (2002). Operating characteristics and extensions of the false discovery rate procedure. <i>Journal of Royal Statistical Society, Ser. B</i>, 46, 409–517.</li> <li>G. Lu, P. E., MURAY, W. and WAGRY, M. H. (1981). <i>Practical Optimization</i>. Academic Press, 12, 2004.</li> <li>Gorousz, C. Nu, R. and Zitao, H. (2004). Model-based approach to FDR estimation. <i>Tech. Report 2004-016</i>. Division of Biostatistics, Univ of Minnesota, Minneapolis, MN.</li> <li>Horomense, Y. and TakutaNY, A. (1990). More powerful procedure for multiple tests of significance. <i>Biometrika</i>, 75, 800-803.</li> <li>Horomense, Y. and TakutaNY, A. (1990). Multiple Comparison Procedures. John Wiley and Sons, New York.</li> <li>Hassiman, Y. (1990). More powerful procedures for multiple significance testing. <i>Statistics in Medicine</i>, 9, 811-818.</li> <li>Horomense, Y. and TakutaNY, A. C. (1987). <i>Multiple Comparison Procedures</i>, John Wiley and Sons, New York.</li> <li>Hexen, H. C. UNN, J. J. and KODEL, R. L. (2003). Comparison of methods for estimating the number of true null hypotheses in multiple testing. <i>Journal of Biopharmaceutical Statistics</i>, 13, 67-689.</li> <li>Huey, Lu AN, SANA, S. (2007). An adaptive single-step FDR procedure with applications to DNA microarray analysis. <i>Biometrical Journal</i>, 49, 127-135.</li> <li>Jukawa, H. and Denze, F. W. (2006). Estimating the proportion of the neu null hypotheses for multiple comparisons. Preprint.</li> <li>Horometay, A. (9, 49-49).</li> <li>Huwa, K. and Sanka, S. (2007). <i>Continuous Univariate Distributions</i>, J. John Wiley &amp; Song.</li> <li>Huwa, K. and Sanka, S. (2007). Continuous Univariate Distributio</li></ul>  | 6  | [3]            | multiple testing with independent statistics. <i>Journal of Educational Statistics</i> , <b>25</b> , 60–83.<br>BLACK M A (2004) A note on the adaptive control of false discovery rates. <i>Journal of Boual</i> | 6  |
| <ul> <li>[4] DEMFETRE, A. P., LARD, N. M. and RUBN, D. B. (1977). Maximum likelihood from incomplete data via the EM Algorithm. Journal of Royal Statistical Society, Ser. B, 39, 1-38.</li> <li>[5] FINNER, H. and ROTERS, M. (2001). On the false discovery rate and expected type I error. Biometrical Journal, 8, 985-1005.</li> <li>[6] GENOVESE, C. and WASSEMAN, L. (2002). Operating characteristics and extensions of the false lideovery rate procedure. Journal of Royal Statistical Society, Ser. B, 64, 499-517.</li> <li>[7] CHL, P. E., MUHRAY, W. and WHUEHT, M. H. (1981). Practical Optimization. Academic Press, London and New York.</li> <li>[8] GUAN, Z., WU, B. and ZINO, H. (2004). Model-based approach to FDR estimation. Tech. Report 12,004-016. Division of Biostatintiscs, Univ. of Minneota, Minneapolis, MN.</li> <li>[9] HOCHBERG, Y. (1988). A sharper Bonferroni procedure for multiple tests of significance. Biometrika, 75, 800-803.</li> <li>[10] HOCHBERG, Y. and BENJAMINI, Y. (1990). More powerful procedures for multiple significance testing. Statistics in Modifue, 9, 811-818.</li> <li>[11] HOCHBERG, Y. and KONELL, R. L. (2003). Comparison of methods for estimating the number of true mult hypotheses in multiple testing. Journal of Biophammaceutical Statistics, 13, 675-689.</li> <li>[12] HENGH, H., CHEN, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number of true mult hypotheses in multiple testing. Journal of Biophammaceutical Statistics, 13, 675-672.</li> <li>[13] HUNG, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparison. Procedures, Many 19, 197-135.</li> <li>[14] VIER, V. and SMAKAR, S. (2007). An adaptive single-step PDR procedures. Paper presented at the first multiple comparison. Procedures, 19, 40, 655-672.</li> <li>[15] JANG, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparison. Procedures, Paper presented at the first maphicati</li></ul>   | 7  | [0]            | Statistical Society, Ser. B, 66, 297–304.  | 7  |
| <ol> <li>via the EM algorithm. Journal of Neural of Neural Statistical Society, Ser. B, 39, 1–38.</li> <li>[5] FINRE, H. and ROTTES, M. (2001). On the fails discovery rate and expected type I error. Biometrical Journal, 8, 965-1005.</li> <li>[6] G. KENWES, C. and WASSEMAN, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. Journal of Royal Statistical Society, Ser. B, 64, 499–517.</li> <li>[7] G. TL, P. E., MURAY, W. and WKREYT, M. H. (1981). Practical Optimization. Academic Press, London and New York.</li> <li>[8] GUAX, Z., WU, B. and ZIAO, H. (2004). Model-based approach to FDR estimation. Tech. Report 2004-016. Division of Biostatistics, Univ. of Minneotola, Minneapolia, MN.</li> <li>[10] HOCHBERG, Y. And BENJAMIN, Y. (1990). More powerful procedures for multiple tests of significance. Biometrika, 75, 800-803.</li> <li>[11] HOCHBERG, Y. and TAMIANER, A. C. (1987). Multiple Comparison Procedures. John Wiley and Sons, New York.</li> <li>[12] BUERH, H., CHER, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number of true mult hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, 13, 075-699.</li> <li>[13] HUNG, H. M., ONELL, R. T., BUERE, P. and KONEX, K. (1997). The behavior of the p-value when the alternative hypothesis is true. Biometrics, 53, 11-22.</li> <li>[14] VIRE, V. and SAUKAN, S. (2007). Continuous Univariate Distributions, I. John Wiley &amp; Sons, The., New York.</li> <li>[15] JANG, H. and DORBERG, R. W. (2005). Estimating the proportion of the run mult hypotheses for multiple comparisons. Preprint.</li> <li>[16] JONSSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons, Inc., New York.</li> <li>[17] LANGAM, M., LINDQUETE, S. H. And FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses for multiple comparisons. Preprint.</li> <li>[18] JUNGER, J. A. and KOTZ, S. (1970). Continuous Univariate Distributions, Chan</li></ol>   | 8  | [4]            | DEMPSTER, A. P., LAIRD, N. M. and RUBIN, D. B. (1977). Maximum likelihood from incomplete data   | 8  |
| <ol> <li>Fried Journal, 8, 985-1005.</li> <li>GROVERS, C. and WASEBANA, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. Journal of Royal Statistical Society, Scr. B, 64, 499-517.</li> <li>GUL, P. E., MURAN, W. and WRIGHY, M. H. (1981). Practical Optimization. Academic Press, London and New York.</li> <li>GUAN, Z., WU, B. and ZHO, H. (2004). Model-based approach to FDR estimation. Tech. Report 42004-016. Division of Biostatistics, Univ. of Minnesota, Minneapolis, MN.</li> <li>HOCHBER, Y. (1988). A sharper Bonferroni procedure for multiple tests of significance. Biometrika, 75 800-803.</li> <li>HOCHBER, Y. and BENJAMINI, Y. (1990). More powerful procedures for multiple significance testing. Statistics in Medicine, 9, 811-818.</li> <li>HOCHBER, Y. and BENJAMINI, Y. (1990). More powerful procedures for multiple significance testing. Statistics in Medicine, 9, 811-818.</li> <li>HOCHBER, Y. and KOELL, R. L. (2003). Comparison of methods for estimating the number of three mult hypotheses in multiple testing. Journal of Biopharmaccuical Statistics, 13, 675-692.</li> <li>HUNG, H. M., O'NELL, R. T., BAUER, P. and KONEX, K. (1997). The behavior of the p-value when the alternative hypothesis in three. Biometrics, 53, 11-22.</li> <li>IVRE, V. and SAKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA mitple comparison. Properint.</li> <li>JANG, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons. Preprint.</li> <li>JANG, H. and DOERGE, R. W. (2005). Estimating the proportion of true null hypotheses for multiple comparisons. Preprint.</li> <li>JANG, H. and DOERGE, R. W. (2005). Estimating the proportion of true null hypotheses with Applications to DNA microarray data. Journal of Royal Statistical Society, Ser. B, 67, 555-572.</li> <li>SouwEDER, T. and SPIOTOLE, E. (1982). Flots of p-values to evaluate many tests simultaneously. Biometrika, 69, 493</li></ol>   | 9  | [5]            | via the EM algorithm. Journal of Royal Statistical Society, Ser. B, <b>39</b> , 1–38.<br>FINNER H and ROTERS M (2001) On the false discovery rate and expected type I error <i>Biomet</i> -                      | 9  |
| <ol> <li>GENOYES, C. and WASERMAN, L. (2002). Operating characteristics and extensions of the false discovery rate procedure. Journal of Royal Statistical Society, Ser. B, 64, 499–517.</li> <li>GILL, P. E., MURRAY, W. and WIGUT, M. H. (1981). Practical Optimization. Academic Press, London and New York.</li> <li>GCAN, Z., WU, B. and ZIAO, H. (2004). Model-based approach to FDR estimation. Tech. Report 2004-016. Division of Biotatistics, UN: v. of Minnesota, Minnespolis, MN.</li> <li>HOCHBERG, Y. (1985). A sharper Bonferroni procedure for multiple tests of significance. Biometrika, 75, 800-803.</li> <li>HOCHBERG, Y. and BENAMIN, Y. (1990). More powerful procedures for multiple significance testing. Statistics in Medicine, 9, 811–818.</li> <li>HOCHBERG, Y. and TAMIANE, A. C. (1987). Multiple Comparison of methods for estimating the number of three mult hypotheses in true. Biometrika, 75, 601–803.</li> <li>HSUR, H., CURN, J. J. and KODEL, R. L. (2003). Comparison of methods for estimating the number of three multipotheses in true. Biometrika, 51, 40, 127–129.</li> <li>HSURA, H., O'NELL, R. T., BAUB, P. and KOINS, K. (1997). The behavior of the p-value when the alternative hypothesis is true. Biometrika, 64, 127–123.</li> <li>HUKE, H. M., O'NELL, R. T., BAUB, P. and KOINS, J. (4), 127–135.</li> <li>JANG, H. and DOEGGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons. Preprint.</li> <li>JANG, H. and DOEGGE, R. W. (2005). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser. B, 66, 9403–502.</li> <li>SCHWEDER, T. and SPJOTVOLT, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 60, 6493–502.</li> <li>SCHWEDER, T. and SPJOTVOLT, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 60, 6493–502.</li> <li>SCHWEDER, T. and SPJOTVOLT, E. (1982). Plots of p-values to evaluate many tests sim</li></ol>   | 10 | [0]            | rical Journal, 8, 985-1005.  | 10 |
| <ul> <li>discovery rate procedure. Journal of Noyal Statistical Society, Ser. B, 64, 499–517.</li> <li>[7] CILL, P. E., MURAY, W. and WHGHT, M. H. (1981). Protectical Optimization. Academic Press,<br/>London and New York.</li> <li>[8] GUAN, Z., WU, B. and ZHAO, H. (2004). Model-based approach to FDR estimation. Tech. Report<br/>2004-016. Division of Biostatistics, Univ. of Minnesota, Minnesople, MN.</li> <li>[9] HOCHBERG, Y. (1988). A Sharper Bonferroni procedure for multiple tests of significance.<br/>Biometrika, 75, 800-803.</li> <li>[10] HOCHBERG, Y. (1988). A Sharper Bonferroni procedures for multiple significance testing.<br/>Statistics in Medicine, 9, 811-818.</li> <li>[11] HOCHBERG, Y. and TAMANK, A. C. (1987). Multiple Comparison Procedures. John Wiley and<br/>Sons, New York.</li> <li>[12] HENT, H. C. RUN, J. J. and KODEL, R. L. (2003). Comparison of methods for estimating the number<br/>of true null hypothesis is true. Biometrics, 53, 11-22.</li> <li>[14] HENG, H. M. (NNEL, R. T., BAUR, P. and KONINK, K. (1997). The behavior of the p-value when<br/>the alternative hypothesis is true. Biometrics, 53, 11-22.</li> <li>[15] JIANG, H. and DARCE, R. W. (2003). Estimating the proportion of the true null hypotheses for<br/>multiple comparisons. Preprint.</li> <li>[16] JOINSON, N. L. and KORZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,<br/>Inc., New York.</li> <li>[17] LANGAAS, M., LINQUER, B. H. and FERINKSTAD, E. (2004). Estimating the proportion of true null<br/>hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.<br/>B, 67, 555–572.</li> <li>[18] STOREY, J. (2005). Multiple requirements for multiple test procedures. Paper presented at the<br/>IVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>[29] STOREY, J. (2005). Multiple requirements for multiple test procedures. Paper presented at the<br/>IVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>[20] STOREY, J. (2002). A direct approach to false discovery rates.</li></ul>  | 11 | [6]            | GENOVESE, C. and WASSERMAN, L. (2002). Operating characteristics and extensions of the false   | 11 |
| <ol> <li>I. Dondon and New York.</li> <li>London and New York.</li> <li>Conx, Z., WU, B. and ZHAO, H. (2004). Model-based approach to FDR estimation. Tech. Report 2004-016. Division of Biostatistics, Univ. of Minnesota, Minneapolis, NN.</li> <li>B. HOCHIBERG, Y. (1988). A sharper Bonferroni procedure for multiple tests of significance. Biometrika, 75, 800-803.</li> <li>HOCHIBERG, Y. and BENJAMIN, Y. (1990). More powerful procedures for multiple significance testing. Statistics in Medicine, 9, 811-818.</li> <li>HOCHIBERG, Y. and RENJAMIN, A. C. (1987). Multiple Comparison Procedures. John Wiley and Sons, New York.</li> <li>HSUEH, H., CHEN, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number of true null hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, 13, 675-689.</li> <li>HENKE, H. M., O'NELL, R. T., BATER, P. and KOINE, K. (1997). The behavior of the p-value when the alternative hypothesis is true. Biometrics, 53, 11-22.</li> <li>HENK, H. M., O'NELL, R. Y., BATER, P. and KOINE, K. (1997). The behavior of the p-value when the alternative hypothesis. Biometrical Journal, 49, 127-135.</li> <li>JIANG, H. and DORGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons, Preprint.</li> <li>JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons, Inc., New York.</li> <li>T. LANGAAS, M., LINDQUET, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrica, 69, 493-502.</li> <li>SUAPEER, T. and SPIOTVOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrica, 69, 493-502.</li> <li>SUAPEER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the TVth International Conference on Multiple Comparison Procedures. Paper presented at the TVth International Conference on Multiple Comparison procedures. Paper presented at the TVth International Conference on Multiple Co</li></ol>   | 12 | [7]            | discovery rate procedure. Journal of Royal Statistical Society, Ser. B, 64, 499–517.   | 12 |
| <ol> <li>[8] GUAN, Z., WU, B. and ZHAO, H. (2004). Model-based approach to FDR estimation. Tech. Report 2004-016. Division of Biostatistics, Ulv. vol Minnesoda, Minneapolis, NN.</li> <li>[9] HOCHBREG, Y. (1988). A sharper Bonferroni procedure for multiple tests of significance. Biometrika, <b>75</b>, 800-803.</li> <li>[10] HOCHBREG, Y. and BENAMIN, Y. (1990). More powerful procedures for multiple significance testing.</li> <li>[11] Katistice in Medicine, <b>9</b>, 811-818.</li> <li>[12] HUNG, H. M., ONELL, R. L. (2003). Comparison of methods for estimating the number of true null hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, <b>13</b>, 675-689.</li> <li>[13] HUNG, H. M., ONELL, R. T. BAUR, P. and KONEK, K. (1997). The behavior of the p-value when the alternative hypothesis is true. Biometrics, <b>53</b>, 11-22.</li> <li>[14] VIE, V. and SLKKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA microarray analysis. Biometrical Journal, <b>49</b>, 127-135.</li> <li>[15] JUNG, H. and DOEKE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons. Preprint.</li> <li>[16] JOINSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons, Inc., New York.</li> <li>[17] LANCAAS, M., LENQUEY, B. H. and FERKINSEND, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser. B, <b>64</b>, 243-502.</li> <li>[18] SCHWEDR, T. and SPERVOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, <b>69</b>, 403-502.</li> <li>[19] Statistical Society, Ser. B, <b>66</b>, 187-205.</li> <li>[20] SUM, J. (2005). Multiple requirements for multiple test procedures. Paper presented at the IVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>[21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society, Ser. B, <b>66</b>, 1</li></ol>  | 13 | [•]            | London and New York.   | 13 |
| <ul> <li>2004-016. Division of Biostatistics, Univ. of Minnesota, Minneso</li></ul>   | 14 | [8]            | GUAN, Z., WU, B. and ZHAO, H. (2004). Model-based approach to FDR estimation. Tech. Report   | 14 |
| <ol> <li>F. Biometrika, <b>75</b>, 800–803.</li> <li>Biometrika, <b>75</b>, 800–803.</li> <li>Hoenneac, Y. and BEXAMNN, Y. (1990). More powerful procedures for multiple significance testing.</li> <li>Biometrika, <b>76</b>, 800–803.</li> <li>Hoenneac, Y. and TAMHANE, A. C. (1987). Multiple Comparison Procedures. John Wiley and Sons, New York.</li> <li>HSUEH, H., CHEN, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number of true mult hypotheses in multiple testing. Journal of Biopharmaccurical Statistics, <b>18</b>, 675–689.</li> <li>HUKG, H. M., O'NELL, R. T., BAUER, P. and KONNE, K. (1997). The behavior of the p-value when the alternative hypothesis is true. Biometrics, <b>53</b>, 11–22.</li> <li>I-YER, V. and SARKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA microarray analysis. Biometrical Journal, <b>49</b>, 127–135.</li> <li>JANKO, H. and DOEKGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons. Preprint.</li> <li>JOINSSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons, Inc., New York.</li> <li>LANGAAS, M., LINQUIST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Boyal Statistical Society, Ser. B, <b>64</b>, 595–572.</li> <li>SULAFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the TV: th International Conference on Multiple Comparison. Proceedures. Shanghai, China.</li> <li>SULAFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the TV: th International Conference on Multiple Comparison Proceedures. Shanghai, China.</li> <li>SULAFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the TV: th International Conference on Multiple Comparison Proceedures. Shanghai, China.</li> <li>SURAFER, J. P. (2005). Multiple requirements for multiple test p</li></ol>  | 15 | [9]            | 2004-016. Division of Biostatistics, Univ. of Minnesota, Minneapolis, MN.<br>HOCHBERG, V (1988) A sharper Bonferroni procedure for multiple tests of significance  | 15 |
| <ol> <li>HOCIBERG, Y. and BENNANN, Y. (1990). More powerful procedures for multiple significance testing.</li> <li>Statistics in Medicine, 9, 811–818.</li> <li>HOCIBERG, Y. and TAMBARE, A. C. (1987). Multiple Comparison Procedures. John Wiley and<br/>Sons, New York.</li> <li>HSUER, H., CHEN, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number<br/>of true null hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, 13, 675–689.</li> <li>HUNG, H. M., O'NELL, R. T., BAUER, P. and KONEK, K. (1997). The behavior of the p-value when<br/>the alternative hypothesis is true. Biometrics, 53, 11–22.</li> <li>HYR, V. and DOERGE, R. W. (2007). An adaptive single-step FDR procedure with applications to DNA<br/>microarray analysis. Biometrical Journal, 49, 127–135.</li> <li>JANG, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for<br/>multiple comparisons. Preprint.</li> <li>JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,<br/>Inc., New York.</li> <li>LANGAAS, M., LINDQUST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null<br/>hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.<br/>B, 67, 555–572.</li> <li>SUMEDER, T. and STSOTYOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously.<br/>Biometrika, 69, 493–502.</li> <li>SINJFER, J. C2002). Multiple requirements for multiple test procedures. Phaper presented at the<br/>Vh International Conference on Multiple Comparison Procedures. Shanghia, China.</li> <li>SUM, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br/>to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br/>Strong J. (2002). A direct approach to false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187–205.</li> <li>TUKHUEMERER, F. E., SMITH, C. B. and SCHMUN, D. (2004). Strong contr</li></ol>   | 16 | [9]            | Biometrika, 75, 800–803.   | 16 |
| <ul> <li>Statistics in Medicine, 9, 811-818.</li> <li>[14] HOCHERG, V. and TAMINE, A. C. (1987). Multiple Comparison Procedures. John Wiley and<br/>Sons, New York.</li> <li>[12] HSUER, H., CHUN, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number of<br/>of true null hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, 13, 675-689.</li> <li>[13] HUNG, H. M., O'NELL, R. T., BAUER, P. and KOINE, K. (1997). The behavior of the p-value when<br/>the alternative hypothesis is true. Biometrics, 53, 11-22.</li> <li>[14] IYER, V. and SARKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA<br/>microarray analysis. Biometrical Journal 49, 127-135.</li> <li>[15] JLANG, H. and DORGE, R. W. (2005). Estimating the proportion of the true null hypotheses for<br/>multiple comparisons. Preprint.</li> <li>[16] JOINSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,<br/>Inc., New York.</li> <li>[17] LANGAS, M., LENDQUST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null<br/>hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.<br/>B, 67, 555-572.</li> <li>[18] SCHWEDER, T. and SPJøTVOL, E. (1982). Plots of p-values to evaluate many tests simultaneously.<br/>Biometrika, 69, 493-502.</li> <li>[19] SHAFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the<br/>IVth International Conference on Multiple Comparison Procedures. Shanghai, China.</li> <li>[20] Stin, J. (2000). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br/>to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br/>Northwestern University, Evanston, IL.</li> <li>[21] STOREY, J. (2002). A direct approach to false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187-205.</li> <li>[23] TUKKIMEMER, F. E., SMITH, C. B. and SCHMIT, K. (2001). Estimation of the number of true' null<br/>hypotheses in multivar</li></ul>   | 17 | [10]           | HOCHBERG, Y. and BENJAMINI, Y. (1990). More powerful procedures for multiple significance testing.   | 17 |
| <ol> <li>F. M. KARLER, H. CHEN, R. L. (2003). Comparison of methods for estimating the number of true null hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, 31, 675-689.</li> <li>HUNG, H. M., O'NELL, R. T., BLUER, P. and KONER, K. (1997). The behavior of the p-value when the alternative hypothesis is itrue. Biometrics, 53, 11-22.</li> <li>HUNG, H. M., O'NELL, R. T., BLUER, P. and KONER, K. (1997). The behavior of the p-value when the alternative hypothesis is itrue. Biometrics, 50, 11-22.</li> <li>IYER, V. and SARKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA microarray analysis. Biometrical Journal, 49, 127-135.</li> <li>JIANG, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons. Preprint.</li> <li>JOINSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons, Inc., New York.</li> <li>LANCAAS, M., LINNQUER, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser. B, 67, 55-572.</li> <li>SCHWEDER, T. and STIPTOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 69, 493-502.</li> <li>SCHWEDER, T. 2005). Multiple requirements for multiple test procedures. Paper presented at the TVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>STOREY, J. (2002). Adirect approach to false discovery rates. Journal of Royal Statistical Society, Ser. B, 64, 297-304.</li> <li>STOREY, J. (2002). Adirect approach to false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 64, 297-304.</li> <li>STOREY, J. (2002). Adirect approach to false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 66, 187-205.</li> <li>STOREY, J. (2002). Adirect approach to false discovery rates: a unified approach. Journal of</li></ol>  | 18 | [11]           | Statistics in Medicine, 9, 811–818.<br>HOCHBERG V and TAMHANE A C (1987) Multiple Comparison Procedures John Wiley and   | 18 |
| <ol> <li>HSUER, H., CHEN, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number<br/>of true null hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, 13, 675–689.</li> <li>HUNG, H. M., O'NELL, R. T., BAUER, P. and KONNE, K. (1997). The behavior of the p-value when<br/>the alternative hypothesis is true. Biometrics, 53, 11–22.</li> <li>HYER, V. and SARKA, S. (2007). An adaptive single-step FDR procedure with applications to DNA<br/>microarray analysis. Biometrical Journal, 49, 127–135.</li> <li>JAING, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for<br/>multiple comparisons. Preprint.</li> <li>JOINSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,<br/>Inc., New York.</li> <li>LANGAAS, M., LINQOUST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null<br/>hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.<br/>B, 67, 555–572.</li> <li>SCHWEDER, T. and SF197VOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously.<br/>Biometrika, 69, 493–502.</li> <li>SIR, J. (2006). Multiple requirements for multiple test procedures. Paper presented at the<br/>TVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>SIN, J. (2006). Adirect approach to false discovery rates. Journal of Royal Statistical Society,<br/>Scr. B, 64, 297–304.</li> <li>STOREY, J., TAVLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation<br/>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187–205.</li> <li>STOREY, J., TAVLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point. Similation<br/>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187–205.</li> <li>TURKIEMERER, F. E., SUTH, C. B. and SIEGMUND, D</li></ol>   | 19 | [++]           | Sons, New York.  | 19 |
| 21       of true null hypotheses in multiple testing. Journal of Biopharmacettical Statistics, 13, 673–689.       21         22       [14] IVER, U., ONELL, R. T., BALER, P. and KOINE, K. (1997). The behavior of the p-value when the alternative hypothesis is true. Biometrics, 53, 11–22.       [21]         23       [14] IVER, V. and SARKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA microarray analysis. Biometrical Journal, 49, 127–135.       23         24       [15] JJANG, H. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley & Sons, Inc., New York.       26         26       [16] JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley & Sons, Inc., New York.       27         27       LANGAAS, M., LINDQUIST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser. 9, 67, 555–572.       29         28       SCHWEDER, T. and SPISTYOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 69, 493–502.       29         29       Sit, J. (2006). Multiple requirements for multiple test procedures. Paper presented at the IVth International Conference on Multiple Comparison Procedures, Snaphai, China.       30         30       Sit, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society, Ser. 8, 64, 297–304.       35         31       Strokey, J. TAYLOR, J. E. and SEGMUND, D. (2004). Strong control, conservative point estimation and simultaneo   | 20 | [12]           | HSUEH, H., CHEN, J. J. and KODELL, R. L. (2003). Comparison of methods for estimating the number   | 20 |
| <ul> <li>[15] FIGHT MER, P. MARKER, S. (2007). An adaptive single-step FDR procedure with applications to DNA the adit of the alternative hypothesis is true. Biometrica, 53, 11–22.</li> <li>[14] YPER, V. and SARKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA discontrol of the true null hypotheses for multiple comparisons. Preprint.</li> <li>[15] JANKG, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons. Preprint.</li> <li>[16] JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons, Inc., New York.</li> <li>[17] LANGAAS, M., LINDQUIST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser. B, 64, 555–572.</li> <li>[18] SCHWEDER, T. and SPISTVOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 69, 493–502.</li> <li>[19] SIAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Phanghai, China.</li> <li>[20] SUL, J. (2000). Improved Estimation of the Proportion of True Null Hypotheses with Applications at to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics, Northwestern University, Evanston, IL.</li> <li>[21] STOREY, J. CAUOS). A length of the Statistical Society, Ser. B, 64, 297–304.</li> <li>[22] STOREY, J. TAVLOR, J. E. and SEGMUND, D. (2004). Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 64, 297–304.</li> <li>[23] STOREY, J. TAVLOR, J. E. and SEGMUND, D. (2004). Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 66, 187–205.</li> <li>[24] WU, B., CUAN, Z. and ZIAO, H. (2006). Parametric and nonparametric FDR estimation revis</li></ul>   | 21 | [13]           | of true null hypotheses in multiple testing. Journal of Biopharmaceutical Statistics, <b>13</b> , 675–689.<br>HUNG H M O'NEUL B T BAUER P and KOHNE K (1997) The behavior of the p-value when                    | 21 |
| <ul> <li>[14] IVER, V. and SARKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA microarray analysis. Biometrical Journal, 49, 127–135.</li> <li>[15] JANG, H. and DORKE, R. W. (2005). Estimating the proportion of the true null hypotheses for multiple comparisons. Preprint.</li> <li>[16] JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons, Inc., New York.</li> <li>[17] LANGAAS, M., LINDQUST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser. B, 67, 555–572.</li> <li>[18] SCHWEDER, T. and SPIOTVOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 69, 493–502.</li> <li>[19] SIAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the TVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>[20] SIR, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications 33 to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics, Northwestern University, Evanston, IL.</li> <li>[21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society, Ser. B, 64, 297–304.</li> <li>[22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 66, 187–205.</li> <li>[23] TURKNEIMER, F. E., SMITH, C. B. and SIEGMUND, K. (2001). Estimation of the number of 'true' null hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.</li> <li>[24] WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited. Biometrics, 62, 735–744.</li> <li>[24] WU, B., GUAN, Z. and ZHAO, H. (2006).</li></ul>  | 22 | [10]           | the alternative hypothesis is true. <i>Biometrics</i> , <b>53</b> , 11–22.   | 22 |
| <ul> <li>microarray analysis. Biometrical Journal, 49, 127–135.</li> <li>[15] JIANC, H. and DOERGE, R. W. (2005). Estimating the proportion of the true null hypotheses for<br/>multiple comparisons. Preprint.</li> <li>[16] JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,<br/>Inc., New York.</li> <li>[17] LANCAAS, M., LINDQUEN, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null<br/>hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.<br/>B, 67, 555–572.</li> <li>[18] SCINWEDER, T. and SPJØTVOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously.<br/>Biometrika, 69, 493–502.</li> <li>[19] SHAFEER, J. P. (2005). Multiple requirements for multiple test procedures. Phaper presented at the<br/>IVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>[20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br/>to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br/>Northwestern University, Evanston, IL.</li> <li>[21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,<br/>Ser. B, 64, 297–304.</li> <li>[22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation<br/>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. R, 66, 187–205.</li> <li>[23] TURKHEMER, F. E., SMITH, C. B. and SIEGMUND, K. (2001). Estimation of the number of 'true' null<br/>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.</li> <li>[24] Wu, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.<br/>Biometrics, 62, 735–744.</li> <li>[25]</li> </ul>   | 23 | [14]           | IYER, V. and SARKAR, S. (2007). An adaptive single-step FDR procedure with applications to DNA   | 23 |
| <ul> <li>[16] JUNKON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,</li> <li>[16] JUNNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,</li> <li>[17] LANGAAS, M., LINDQUIST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.</li> <li>[18] SCIWEDER, T. and SPJØTVUL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 69, 493–502.</li> <li>[19] SHAFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the TVH International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>[20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistical Society, Ser. B, 64, 297–304.</li> <li>[21] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 64, 297–304.</li> <li>[23] TURKHEIMER, F. E., SMITH, C. B. and SIEGMUND, D. (2001). Estimation of the number of 'true' null hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.</li> <li>[24] WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited. Biometrics, 62, 735–744.</li> <li>[34] 44</li> <li>[35] 44</li> <li>[36] 44</li> <li>[37] 44</li> <li>[38] 44</li> <li>[39] 44</li> <li>[30] 44</li> <li>[31] 44</li> <li>[32] 55</li> <li>[33] 56</li> <li>[34] 56</li> <li>[35] 56</li> <li>[35] 56</li> </ul>  | 24 | [15]           | microarray analysis. Biometrical Journal, <b>49</b> , 127–135.<br>JUANG H and DOEBGE R W (2005) Estimating the proportion of the true null hypotheses for  | 24 |
| <ul> <li>[16] JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley &amp; Sons,<br/>Inc., New York.</li> <li>[17] LANGAAS, M., LINDQUIST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null<br/>hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.<br/>B, 67, 555-572.</li> <li>[18] SCHWEDER, T. and SPJØTVOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously.<br/>Biometrika, 69, 493-502.</li> <li>[19] SHAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the<br/>YVH International Conference on Multiple Comparison Procedures. Shanghai, China.</li> <li>[20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br/>to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br/>Northwestern University, Evanston, IL.</li> <li>[21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,<br/>Ser. B, 64, 297-304.</li> <li>[22] STOREY, J. (2002). A direct approach to false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187-205.</li> <li>[23] TURKHEIMER, F. E., SMITH, C. B. and SEGMUND, D. (2004). Strong control, conservative point estimation<br/>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187-205.</li> <li>[24] WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.<br/>Biometrics, 62, 735-744.</li> <li>[25] TURKHEIMER, F. E., SMITH, C. B. and SCHMUT, K. (2001). Estimation of the number of 'true' null<br/>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920-930.</li> <li>[24] WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.<br/>Biometrics, 62, 735-744.</li> <li>[26] Songer, Song</li></ul>   | 25 | [10]           | multiple comparisons. Preprint.  | 25 |
| 27       Inc., New York.       27         10       LANGAS, M., LINDQUIST, B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser. B, 67, 555–572.       28         29       [18] SCHWEDER, T. and SPJØTVOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously. Biometrika, 69, 493–502.       30         31       [19] SHAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the VIV International Conference on Multiple Comparison Procedures, Shanghai, China.       32         32       [20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics, Northwestern University, Evanston, IL.       34         34       Northwestern University, Evanston, IL.       36         35       [21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society, Ser. B, 64, 297–304.       36         37       STOREY, J., TAYLOR, J. E. and SEGMUND, D. (2004). Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 66, 187–205.       38         39       [23] TURKHEMER, F. E., SMITH, C. B. and SEGMUND, K. (2001). Estimation of the number of 'true' null hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         41       Biome   | 26 | [16]           | JOHNSON, N. L. and KOTZ, S. (1970). Continuous Univariate Distributions, I. John Wiley & Sons,   | 26 |
| 28       [11] Indication, in Integent in DNA microarray data. Journal of Royal Statistical Society, Ser.<br>B, 67, 555–572.       29         30       [18] SCHWEDER, T. and SPJøTYOLL, E. (1982). Plots of p-values to evaluate many tests simultaneously.<br>Biometrika, 69, 493–502.       30         31       [19] SIAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the<br>IVth International Conference on Multiple Comparison Procedures, Shanghai, China.       32         32       [20] Stu, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br>to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br>Northwestern University, Evanston, IL.       34         33       [21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,<br>Ser. B, 64, 297–304.       36         34       Ser. B, 64, 297–304.       36         35       [22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation<br>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br>Royal Statistical Society, Ser. B, 66, 187–205.       38         36       [23] TURKHEMER, F. E., SMITH, C. B. and SCHMDT, K. (2001). Estimation of the number of 'true' null<br>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         41       Biometrics, 62, 735–744.       41         42       43       44         44       44       44     <  | 27 | [17]           | Inc., New York.<br>LANGAAS M. LINDOUIST B. H. and FERKINGSTAD, E. (2004). Estimating the proportion of true null.  | 27 |
| <ul> <li>B, 67, 555-572.</li> <li>[18] SCHWEDER, T. and SPIøTVOL, E. (1982). Plots of p-values to evaluate many tests simultaneously.<br/>Biometrika, 69, 493-502.</li> <li>[19] SHAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the<br/>IVth International Conference on Multiple Comparison Procedures. Shanghai, China.</li> <li>[20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br/>to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br/>Northwestern University, Evanston, IL.</li> <li>[21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,<br/>Ser. B, 64, 297-304.</li> <li>[22] STOREY, J. (2002). A direct approach to false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187-205.</li> <li>[23] TURKHEMER, F. E., SMITH, C. B. and SCIMMDT, K. (2001). Estimation of the number of 'true' null<br/>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920-930.</li> <li>[24] WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.<br/>Biometrics, 62, 735-744.</li> <li>[24] Wu, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.</li> <li>[25] 44</li> <li>[26] 45</li> <li>[27] 55, 54, 54</li> <li>[28] 50</li> <li>[29] 50</li> <li>[29] 50</li> <li>[20] 50</li> <li>[20] 50</li> <li>[21] 50</li> <li>[22] 50</li> <li>[23] 50</li> <li>[24] 50</li> <li>[25] 50</li> <li>[24] 50</li> <li>[25] 50</li> <li>[26] 50</li> <li>[27] 50</li> <li>[28] 50</li> <li>[29] 50</li> <li>[29] 50</li> <li>[29] 50</li> <li>[29] 50</li> <li>[20] 50</li> <li>[20] 50</li> <li>[20] 50</li> <li>[21] 50</li> <li>[22] 50</li> <li>[23] 50</li> <li>[24] 50</li> <li>[25] 50</li> <li>[25] 50</li> <li>[26] 50</li> <li>[27] 50</li> <li>[28] 50</li> <li>[29] 50</li> <li>[29] 50</li> <li>[29] 50</li> <li>[20] 50</li> <li>[20] 50</li> <li>[20] 50</li> <li>[20] 50</li> <li>[20] 50<td>28</td><td>[1]</td><td>hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.</td><td>28</td></li></ul> | 28 | [1]            | hypotheses, with application to DNA microarray data. Journal of Royal Statistical Society, Ser.  | 28 |
| <ul> <li>[18] SCHWEDER, T. and SPJøTVOLL, E. (1982). Plots of <i>p</i>-values to evaluate many tests simultaneously.<br/>Biometrika, 69, 493-502.</li> <li>[19] SHAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the<br/>IVth International Conference on Multiple Comparison Procedures, Shanghai, China.</li> <li>[20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br/>to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br/>Northwestern University, Evanston, IL.</li> <li>[21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,<br/>Ser. B, 64, 297-304.</li> <li>[22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation<br/>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br/>Royal Statistical Society, Ser. B, 66, 187-205.</li> <li>[23] TURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null<br/>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920-930.</li> <li>[24] WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.<br/>Biometrics, 62, 735-744.</li> <li>[25] Multiple Comparison Proceeding and the statistic and nonparametric FDR estimation revisited.</li> <li>[26] TURKHEIMER, F. 2, STATE, A.</li> <li>[27] TURKHEIMER, F. 2, STATE, S.</li> <li>[28] TURKHEIMER, F. 2, STATE, S.</li> <li>[29] TURKHEIMER, F. 3, STATE, S.</li> <li>[29] TURKHEIMER, F. 3, STATE, S.</li> <li>[29] TURKHEIMER, F. 4, STATE, S.</li> <li>[29] TURKHEIMER, F. 4, STATE, S.</li> <li>[20] TURKHEIMER, F. 4, STATE, S.</li> <li>[20] TURKHEIMER, F. 4, STATE, S.&lt;</li></ul>   | 29 | [10]           | <i>B</i> , <b>67</b> , 555–572.  | 29 |
| 31       [19] SHAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the       31         32       [19] SHAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Shanghai, China.       32         33       [20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics, Northwestern University, Evanston, IL.       34         35       [21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society, Ser. B, 64, 297–304.       36         36       Royal Statistical Society, Ser. B, 66, 187–205.       37         37       [23] TURKNEIMER, F. E., SMITH, C. B. and SICMINDT, K. (2001). Estimation of the number of 'true' null hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         41       Biometrics, 62, 735–744.       41         42       43       44         44       45       46         47       48       49         50       50       50         51       51       51  | 30 | [18]           | SCHWEDER, T. and SPJøTVOLL, E. (1982). Plots of <i>p</i> -values to evaluate many tests simultaneously.<br><i>Biometrika</i> , <b>69</b> , 493–502   | 30 |
| 32       IVth International Conference on Multiple Comparison Procedures, Shanghai, China.       32         33       [20] SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics, Northwestern University, Evanston, IL.       34         35       [21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society, Ser. B, 64, 297–304.       36         37       [22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of Royal Statistical Society, Ser. B, 66, 187–205.       38         39       [23] TURKHEIMER, F. E., SHITH, C. B. and SCHMDT, K. (2001). Estimation of the number of 'true' null hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         41       Biometrics, 62, 735–744.       41         42       43       44         44       44       44         45       46       46         46       47       47         48       49       49         50       50       50  | 31 | [19]           | SHAFFER, J. P. (2005). Multiple requirements for multiple test procedures. Paper presented at the  | 31 |
| 33       [20] Shi, J. (2006). Improved Estimation of the Proportion of The Null Hypotheses with Applications       33         34       to Adaptive Control of FDR and Drug Screening. Doctoral Dissertation. Department of Statistics,<br>Northwestern University, Evanston, IL.       34         35       [21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,<br>Ser. B, 64, 297–304.       36         37       [22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation<br>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br>Royal Statistical Society, Ser. B, 66, 187–205.       38         39       [23] TURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null<br>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         41       Biometrics, 62, 735–744.       41         42       43       44         44       44       44         45       46       46         47       48       48         49       49       50         50       50       50   | 32 | [20]           | IVth International Conference on Multiple Comparison Procedures, Shanghai, China.  | 32 |
| 34       Northwestern University, Evanston, IL.       34         35       [21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,       35         36       Ser. B, 64, 297–304.       36         37       [22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation       37         38       monthwestern University, Evanston, V. (2004). Strong control, conservative point estimation       37         39       [23] TURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null       39         40       hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         41       Biometrics, 62, 735–744.       41         42       42       42         43       44       44         44       44       44         45       46       46         46       47       47         47       48       49         50       50       50         51       50       50   | 33 | [20]           | SHI, J. (2006). Improved Estimation of the Proportion of True Null Hypotheses with Applications<br>to Adaptive Control of FDB and Drug Screening. <i>Doctoral Dissertation</i> . Department of Statistics        | 33 |
| 35       [21] STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,<br>Ser. B, 64, 297–304.       36         37       [22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation<br>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br>Royal Statistical Society, Ser. B, 66, 187–205.       38         39       [23] TURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null<br>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         41       Biometrics, 62, 735–744.       41         42       42         43       44         44       45         45       46         46       47         47       48         49       49         50       50   | 34 |                | Northwestern University, Evanston, IL.   | 34 |
| 36     Ser. B, 64, 297–304.     36       37     [22] STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation<br>and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br>Royal Statistical Society, Ser. B, 66, 187–205.     38       39     [23] TURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null<br>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.     40       41     Biometrics, 62, 735–744.     41       42     42       43     44       44     44       45     46       46     47       47     48       49     50       51     51  | 35 | [21]           | STOREY, J. (2002). A direct approach to false discovery rates. Journal of Royal Statistical Society,   | 35 |
| 37       11       11       11       37       37         38       and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of<br>Royal Statistical Society, Ser. B, 66, 187–205.       38         39       [23]       TURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null<br>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.       40         10       [24]       WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.       41         41       Biometrics, 62, 735–744.       41         42       42       42         43       44       44         44       45       46         45       46       47         46       47       48         49       49       49         50       50       50   | 36 | [22]           | Ser. B, 04, 291-304.<br>STOREY, J., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control. conservative point estimation   | 36 |
| 38Royal Statistical Society, Ser. B, 66, 187–205.3839[23] TURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null<br>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.4040[24] WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.<br>Biometrics, 62, 735–744.41424243444445454646474748495050505151   | 37 | []             | and simultaneous conservative consistency of false discovery rates: a unified approach. Journal of   | 37 |
| 39[25]FURKHEIMER, F. E., SMITH, C. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null<br>hypotheses in multivariate analysis of neuroimaging data. NeuroImage, 13, 920–930.3040[24]WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.<br>Biometrics, 62, 735–744.41424344444545454646474849505050515151  | 38 | [00]           | Royal Statistical Society, Ser. B, 66, 187–205.  | 38 |
| 40       [24] Wu, B., Guan, Z. and Zhao, H. (2006). Parametric and nonparametric FDR estimation revisited.       40         41       Biometrics, 62, 735–744.       42         43       43       43         44       45       45         46       47       48         49       50       50         51       51       51  | 39 | [23]           | IURKHEIMER, F. E., SMITH, U. B. and SCHMIDT, K. (2001). Estimation of the number of 'true' null hypotheses in multivariate analysis of neuroimaging data. <i>NeuroImage</i> . <b>13</b> , 920–930                | 39 |
| 41     Biometrics, 62, 735–744.     41       42     42       43     43       44     43       45     46       46     46       47     48       48     49       50     50       51     51   | 40 | [24]           | WU, B., GUAN, Z. and ZHAO, H. (2006). Parametric and nonparametric FDR estimation revisited.   | 40 |
| 42     42       43     43       44     44       45     45       46     46       47     47       48     49       50     50       51     51  | 41 |                | <i>Biometrics</i> , <b>62</b> , 735–744.   | 41 |
| 43     43       44     44       45     45       46     46       47     47       48     49       50     50       51     51  | 42 |                |  | 42 |
| 44     44       45     45       46     46       47     47       48     49       50     50       51     51  | 43 |                |  | 43 |
| 45     45       46     46       47     47       48     49       50     50       51     51  | 44 |                |  | 44 |
| 46     46       47     47       48     48       49     49       50     50       51     51  | 45 |                |  | 45 |
| 47     47       48     48       49     49       50     50       51     51  | 46 |                |  | 46 |
| 48     48       49     49       50     50       51     51  | 47 |                |  | 47 |
| 49     49       50     50       51     51  | 48 |                |  | 48 |
| 50         50           51         51  | 49 |                |  | 49 |
| 51 51  | 50 |                |  | 50 |
|  | 51 |                |  | 51 |

# Bayesian Decision Theory for Multiple Comparisons

З

# Charles Lewis<sup>1</sup> and Dorothy T. Thayer<sup>2</sup>

Fordham University and Educational Testing Service

Abstract: Applying a decision theoretic approach to multiple comparisons very similar to that described by Lehmann (1950, 1957a, 1957b), we introduce a loss function based on the concept of the false discovery rate (FDR). We derive a Bayes rule for this loss function and show that it is very closely related to a Bayesian version of the original multiple comparisons procedure proposed by Benjamini and Hochberg (1995) to control the sampling theory FDR. We provide the results of a Monte Carlo simulation that illustrates the very similar sampling behavior of our Bayes rule and Benjamini and Hochberg's procedure when applied to making all pair-wise comparisons in a one-way fixed effects analysis of variance setup with 10 and with 20 means.

#### Contents

| L | Introduction   |
|---|--|
| 2 | Setup  |
| 3 | Bayes Decision Rule                                    |
| 4 | A Bayesian Version of Benjamini & Hochberg's Procedure |
| 5 | Simulation Results                                     |
| 6 | Conclusions  |

### 1. Introduction

A previous paper by the authors (Lewis & Thayer, 2004) considered the application of Bayesian decision theory to the multiple comparisons problem for random effects designs, following the earlier work of Shaffer (1999), Duncan (1965), and Waller and Duncan (1969). In our paper, we demonstrated that the Bayes rule for a per-comparison "0-1" loss function controls a random effects version of the false discovery rate (FDR), thus supporting and extending Shaffer's (1999) results.

A recent paper by Sarkar and Zhou (2008) adopts a random effects setup very similar to that of our earlier paper. Rather than considering Bayes rules, they introduce a procedure that controls the random effects FDR discussed by us while maximizing the random effects per-comparison power rate that we had considered. This approach produces substantial power gains over other procedures (including ours), but it "declares even small differences significant when  $\tau$  [the between-groups standard deviation] is large, thereby achieving [even] greater power than the unadjusted (per-comparison) procedure for large values of  $\tau$ " (Sarkar & Zhou, 2008, p. 692). We view this as a weakness, rather than a strength, of their method, as it seems to ignore the basic principle behind multiple comparisons procedures, namely

| <sup>1</sup> Fordham | University | and | Educational | Testing | Service |
|----------------------|------------|-----|-------------|---------|---------|
|----------------------|------------|-----|-------------|---------|---------|

<sup>2</sup>Educational Testing Service

that making multiple inferences calls for increased conservatism relative to making
a single inference.

The present study considers a more general setting for making multiple comparisons and introduces a new loss function that is more directly tied to the *FDR*. We derive a Bayes rule for this loss function and show that it is very closely related to a Bayesian version of the original multiple comparisons procedure proposed by Benjamini and Hochberg (1995) to control the sampling theory *FDR*. We provide the results of a Monte Carlo simulation that illustrates the very similar sampling behavior of our Bayes rule and Benjamini and Hochberg's procedure when applied to testing all pairwise comparisons for a one-way fixed effects analysis of variance setup with 10 and with 20 means.

#### 2. Setup

We start with a general likelihood  $p(\mathbf{y}|\theta)$ , prior  $p(\theta)$ , and resulting posterior  $p(\theta|\mathbf{y})$ . Let  $\psi = \mathbf{f}(\theta)$  be a vector of m "contrasts" among the elements of  $\theta$ . Suppose our goal is to identify the sign of each of the elements of  $\psi$ , given  $\mathbf{y}$ . In the language of decision theory, for each  $\psi_i$ ,  $i = 1, \dots, m$ , we will take action  $a_i$ , with  $a_i = +1$  used to indicate that we declare  $\psi_i$  to be positive,  $a_i = -1$  indicating that we declare  $\psi_i$  to be negative, and  $a_i = 0$  used to indicate that we are unable to determine sign of  $\psi_i$ . Although directly inspired by Williams, Jones and Tukey (1999), and Jones and Tukey (2000), this approach to (multiple) hypothesis testing has its origins in the much earlier work of Lehmann (1950, 1957a, 1957b).

To continue, we introduce two component loss functions:  $L_1(\psi_i, a_i) = 1$  if the signs of  $\psi_i$  and  $a_i$  disagree, and  $L_1 = 0$  otherwise (used to count the number of incorrect sign declarations);  $L_2(\psi_i, a_i) = 1$  if  $a_i = 0$ , and  $L_2 = 0$  otherwise (used to count the signs not declared). These actions and losses are very similar to those given by Lehmann (1957b, p. 549). They differ from conventional treatments of hypothesis testing in the sense that they focus on identifying the sign of each contrast and do not formally consider the possibility that the value of the contrast could be (exactly) 0. The reasonableness of this approach, compared with conventional point hypothesis testing is emphasized by Jones and Tukey (2000), among others. We now propose a loss function that combines  $L_1$  and  $L_2$  as follows:

(1) 
$$L_{DFDR}(\psi, \mathbf{a}) = \frac{\sum_{i=1}^{m} L_1(\psi_i, a_i)}{\max\left\{1, m - \sum_{i=1}^{m} L_2(\psi_i, a_i)\right\}} + \left(\frac{\alpha}{2}\right) \frac{\sum_{i=1}^{m} L_2(\psi_i, a_i)}{m},$$

for a fixed choice of  $0 < \alpha < 1$  (such as  $\alpha = 0.05$ ). Here, *DFDR* (in notation introduced by Shaffer, 2002) stands for Directional False Discovery Rate. The first term in Equation 1 is the sample value of the *DFDR* for a given  $\psi$  and a vector of actions **a**, namely the number of incorrect sign declarations, divided by the total number of signs declared (or divided by 1 if no signs are declared by **a**).

The second term in Equation 1 is  $\alpha/2$  times the sample proportion of signs not declared. This term may be interpreted as a sample per-comparison Type II error rate, weighted by a relative importance factor of  $\alpha/2$ . Using a per-comparison formulation, as well as assigning this loss component a small weight, serves to emphasize that failure to declare a sign is considered to be much less serious than declaring that sign incorrectly. This emphasis is in keeping with the concern about controlling Type I errors (at the expense of making Type II errors) in traditional

imsart-coll ver. 2008/08/29 file: Thayer.tex date: April 10, 2009

З

treatments of the multiple comparisons problem. The Bayes decision rule for the complete loss function in Equation 1 minimizes its posterior expected value, in this sense balancing the two types of losses against each other, with the major focus being on reducing the *DFDR*.

#### 3. Bayes Decision Rule

To identify the actions that minimize the posterior expected value of the loss function given in Equation 1, we begin by introducing some notation. If  $\Pr(\psi_i > 0 | \mathbf{y}) > 0.5$ , define  $a_i^* = +1$  and  $p_i = \Pr(\psi_i < 0 | \mathbf{y})$ ; if  $\Pr(\psi_i > 0 | \mathbf{y}) \leq 0$ 0.5, define  $a_i^* = -1$  and  $p_i = \Pr(\psi_i > 0 | \mathbf{y})$ . Note that  $a_i^*$  and  $p_i$  are related by the following result:  $E_{\theta|\mathbf{y}}[L_1(\psi_i, a_i^*)|\mathbf{y}] = p_i$ . Now order the  $p_i$  so that  $p_{(1)} \leq p_i$  $\dots \leqslant p_{(m)}. \text{ Define } \mathbf{a}^{(k)} \text{ for } k = 1, \dots, m \text{ as } a^{(k)}_{(i)} = a^*_{(i)}, \text{ for } i = 1, \dots, k, \text{ and } a^{(k)}_{(i)} = 0, \text{ for } i = k+1, \dots, m. \text{ For } k = 0, \text{ take } a^{(0)}_{(i)} = 0, \text{ for } i = 1, \dots, m.$ The posterior expected loss for  $\mathbf{a}^{(k)}$  is given by

З

(2) 
$$E_{\theta|\mathbf{y}}\left[L_{DFDR}\left(\psi, \mathbf{a}^{(k)}\right)|\mathbf{y}\right] = \frac{\sum_{i=1}^{n} p_{(i)}}{\max\left\{1, k\right\}} + \left(\frac{\alpha}{2}\right)\left(1 - \frac{k}{m}\right).$$

ı

Clearly,  $\mathbf{a}^{(k)}$  minimizes the posterior expected loss among all action vectors  $\mathbf{a}$  that declare exactly k signs. Let  $k_{DFDR}$  be the value of k for which the posterior expected loss given in Equation 2 is minimized. (The value of  $k_{DFDR}$  in a given setting would normally be determined by an exhaustive search over all values of  $k = 0, \dots, m$ .) The Bayes decision rule for this problem is given by  $\delta_{DFDR}(\mathbf{y}) = \mathbf{a}^{(k_{DFDR})}$ , and the corresponding Bayes risk is

$$r\left(\delta_{DFDR}\right) = E_{\theta,\mathbf{y}}\left[L_{DFDR}\left(\psi,\delta_{DFDR}\left(\mathbf{y}\right)\right)\right]$$

$$T(0DFDR) = D\theta_{\mathbf{y}} \left[ DFDR(\psi, 0DFDR(\mathbf{y})) \right]$$
$$\begin{bmatrix} k_{DFDR} \\ k_{DFDR} \end{bmatrix}$$

 $= E_{\mathbf{y}} \left[ \frac{\sum_{i=1}^{m_{DFDR}} p_{(i)}}{\max\{1, k_{DFDR}\}} + \left(\frac{\alpha}{2}\right) \left(1 - \frac{k_{DFDR}}{m}\right) \right].$ 

The latter expectation in Equation 3 is taken with respect to the predictive distribution of y, and it should be noted that the  $p_{(i)}$  and, consequently,  $k_{DFDR}$  all depend on y.

Since  $E_{\theta|\mathbf{y}}\left[L_{DFDR}\left(\psi, \mathbf{a}^{(0)}\right)|\mathbf{y}\right] = \alpha/2$  for all  $\mathbf{y}$ , it follows that  $r\left(\delta_{DFDR}\right) \leqslant \alpha/2$ . Consequently,

$$\sum_{i=1}^{m} L_1\left(\psi_i, \delta_{DFDR,i}\left(\mathbf{y}\right)\right)$$

(4) 
$$E_{\theta,\mathbf{y}}\left|\frac{\sum\limits_{i=1}^{L}L_1\left(\psi_i,\delta_{DFDR,i}\left(\mathbf{y}\right)\right)}{\max\left\{1,m-\sum\limits_{i=1}^{m}L_2\left(\psi_i,\delta_{DFDR,i}\left(\mathbf{y}\right)\right)\right\}}\right| \leqslant \frac{\alpha}{2}.$$

Γ

$$E_{\theta,\mathbf{y}}\left[\frac{\frac{1}{i=1}}{\max\left\{1,m-\sum_{i=1}^{m}L_{2}\left(\psi_{i},\delta_{DFDR,i}\left(\mathbf{y}\right)\right)\right\}}\right] \leqslant \frac{\alpha}{2}.$$

Equation 4 says that a Bayesian version of the *DFDR* is bounded by  $\alpha/2$  when the Bayes rule  $\delta_{DFDR}$  is used. It may be worth observing that these results apply to a very general class of multiple comparison problems. Essentially the only restriction is that the set of contrasts be finite. Indeed, these do not even have to be contrasts in the usual sense of that term. They could also, for example, be a set of independent parameters that formed a family of interest. 

# 4. A Bayesian Version of Benjamini & Hochberg's Procedure

Next, we consider the multiple comparisons procedure proposed by Benjamini and Hochberg (1995) and modified for directional testing by Williams, Jones and Tukey (1999).

However, we will translate the procedure into our Bayesian framework. Define  $k_{DB \& H}$  to be the largest value of  $k = 1, \cdots, m$ , such that

$$p_{(k)} \leqslant \left(\frac{\alpha}{2}\right) \left(\frac{k}{m}\right),$$

with  $k_{DB\&H} = 0$  if no such value of k exists. Define  $\delta_{DB\&H}(\mathbf{y}) = \mathbf{a}^{(k_{DB\&H})}$ . If  $k_{DB\&H} = 0$ , then the posterior expected loss for  $\mathbf{a}^{(k_{DB\&H})}$  is equal to  $\alpha/2$ . If  $k_{DB\&H} > 0$ , the posterior expected loss for  $\mathbf{a}^{(k_{DB\&H})}$  is given by Equation 2 as

$$\sum_{n \in \mathcal{N}} \sum_{n \in \mathcal{N}} n$$

(5) 
$$E_{\theta|\mathbf{y}}\left[L_{DFDR}\left(\psi, \mathbf{a}^{(k_{DB \& H})}\right)|\mathbf{y}\right] = \frac{\sum\limits_{i=1}^{N} P(i)}{k_{DB \& H}} + \left(\frac{\alpha}{2}\right)\left(1 - \frac{k_{DB \& H}}{m}\right).$$

From the definition of  $k_{DB \& H}$ , it follows that

 $k_{DFDR}$ 

(6) 
$$p_{(i)} \leqslant \left(\frac{\alpha}{2}\right) \left(\frac{k_{DB \& H}}{m}\right) \text{ for } i = 1, \cdots, k_{DB \& H}.$$

Consequently, applying Inequality 6 to Equation 5, it follows that

$$E_{\theta|\mathbf{y}}\left[L_{DFDR}\left(\psi, \mathbf{a}^{(k_{DB \& H})}\right)|\mathbf{y}\right] \leqslant \left(\frac{\alpha}{2}\right) \left(\frac{k_{DB \& H}}{m}\right) + \left(\frac{\alpha}{2}\right) \left(1 - \frac{k_{DB \& H}}{m}\right) = \frac{\alpha}{2}.$$

Since this inequality holds for all y, it implies that  $r(\delta_{DB\& H}) \leq \alpha/2$ , and so, just as with  $\delta_{DFDR}$ ,

$$E_{\theta,\mathbf{y}}\left[\frac{\sum_{i=1}^{m} L_1\left(\psi_i, \delta_{DB \& H, i}\left(\mathbf{y}\right)\right)}{\sum_{i=1}^{m} L_1\left(\psi_i, \delta_{DB \& H, i}\left(\mathbf{y}\right)\right)}\right] \leqslant \frac{\alpha}{2}$$

(7) 
$$E_{\theta,\mathbf{y}}\left[\frac{i=1}{\max\left\{1,m-\sum_{i=1}^{m}L_{2}\left(\psi_{i},\delta_{DB\&H,i}\left(\mathbf{y}\right)\right)\right\}}\right] \leqslant \frac{\alpha}{2}.$$

Equation 7 says that  $\delta_{DB \& H}$  also controls our Bayesian *DFDR*.

This seems like an appropriate place to note that what has just been established (namely the fact that  $\delta_{DB\& H}$  controls the DFDR for an arbitrary set of contrasts) is a Bayesian, rather than a sampling theory result. Indeed, Benjamini and Hochberg's (1995) sampling theory procedure has only been shown to control the sampling theory FDR in special circumstances, such as the case of independent tests. In particular, it has not been shown to provide sampling theory control of the FDR when making all pairwise comparisons among a set of means in a one-way, fixed effects analysis of variance setup.

Since  $\delta_{DFDR}$  is a Bayes decision rule, it must be the case that  $r(\delta_{DFDR}) \leq$  $r(\delta_{DB\&H})$ . Moreover, it is also possible to show that  $k_{DFDR} \ge k_{DB\&H}$  for all **y**. To see this, suppose the contrary:  $k_{DB \& H} = k_{DFDR} + d$  with d > 0. By the definition of  $k_{DFDR}$ , we must have

$$\sum_{i=1}^{49} \frac{\sum_{i=1}^{k_{DFDR}} p_{(i)}}{\max\{1, k_{DFDR}\}} + \left(\frac{\alpha}{2}\right) \left(1 - \frac{k_{DFDR}}{m}\right) < \frac{\sum_{i=1}^{k_{DB\&H}} p_{(i)}}{k_{DB\&H}} + \left(\frac{\alpha}{2}\right) \left(1 - \frac{k_{DB\&H}}{m}\right)$$

imsart-coll ver. 2008/08/29 file: Thayer.tex date: April 10, 2009

 $k_{DB \& H}$ 

З

or

$$\sum_{i=1}^{k_{DFDR}} p_{(i)} \qquad \qquad \sum_{i=1}^{k_{DFDR}} p_{(i)} + \sum_{i=k_{DFDR}+1}^{k_{DFDR}} p_{(i)}$$

З

$$\frac{1}{\max\left\{1, k_{DFDR}\right\}} + \left(\frac{1}{2}\right) \left(1 - \frac{DFDR}{m}\right) < \frac{1}{k_{DFDR} + d}$$

(8) 
$$+\left(\frac{\alpha}{2}\right)\left(1-\frac{k_{DFDR}+d}{m}\right).$$

Note that a strict inequality has been used here, implying that, in case of ties,  $k_{DFDR}$  would be chosen to be the largest value of k that minimizes the posterior expected loss. To continue, using the definition of  $k_{DB\& H}$ ,

(9) 
$$\sum_{i=k_{DFDR}+1}^{k_{DFDR}+d} p_{(i)} \leq d\left(\frac{\alpha}{2}\right) \left(\frac{k_{DFDR}+d}{m}\right).$$

Combining Inequalities 8 and 9 gives

$$\max\{1, k_{DFDR}\} \quad (2) \quad (m) \quad k_{DFDR} + d \\ + \left(\frac{\alpha}{2}\right) \left(\frac{d}{m}\right) + \left(\frac{\alpha}{2}\right) \left(1 - \frac{k_{DFDR} + d}{m}\right)$$

or

$$\frac{\sum_{i=1}^{k_{DFDR}} p_{(i)}}{\sum_{i=1}^{k_{DFDR}} p_{(i)}} < \frac{\sum_{i=1}^{k_{DFDR}} p_{(i)}}{\sum_{i=1}^{k_{DFDR}} p_{(i)}},$$

(10) 
$$\frac{i=1}{\max\{1, k_{DFDR}\}} < \frac{i=1}{k_{DFDR}+d},$$

If  $k_{DFDR} > 0$ , Inequality 10 implies that d = 0, contrary to our initial assumption. If  $k_{DFDR} = 0$ , both sides of Inequality 10 would be 0, contradicting the strict inequality. Thus, we have demonstrated that  $k_{DFDR} \ge k_{DB\&H}$  for all **y**. In other words, the Bayes rule  $\delta_{DFDR}$  will always declare at least as many signs as  $\delta_{DB\&H}$ .

## 5. Simulation Results

It is important to recall that the procedure actually proposed by Benjamini and Hochberg (1995) uses sampling theory p-values (one-tailed values in Williams, Jones and Tukey's 1999 version), rather than posterior tail probabilities and controls the sampling theory version of the FDR (or DFDR in Willams et al.'s version). Now consider a standard multiple comparisons problem: the one-way, fixed effects analysis of variance setup, with the  $\psi_i$  chosen to be all pair-wise differences among the group means. In this case, the relevant sampling theory and Bayesian (based on a vague prior for all parameters) p-values are identical tail probabilities from the appropriate Student's t-distribution (see, for instance, Box & Tiao, 1973, p. 140).

Tables 1 and 2 give the results of sampling theory simulations (based on 25,000 replications for each condition) of one-way, fixed effects ANOVA setups, considering all pair-wise differences for 10 evenly spaced means, and for 25 evenly spaced means. In these simulations, the within-group variance was set at 3.0 with n = 3 sample observations per group (so the sampling variances of the sample means are all equal to 1.0 and the within-group degrees of freedom equals 20 in the first case and 50 in 

| the second case),<br>the sampling the<br>there is no theory.<br>The <i>DFDR</i> va-<br>used as the first<br>means, and a rar<br>across the two set<br>$\tau$ ) to index the sp<br>pling theory aver<br>All four quantities<br>replications at ear | In addition<br>every <i>DFDR</i><br>y to suppor<br>lues in Tab<br>term of ou-<br>age of spacin<br>tups, we us<br>pread of the<br>rages of the<br>es for a given<br>ach spread of | n, we choose $\alpha$ at $\alpha/2$<br>t that co-<br>le 1 are s<br>r loss fun-<br>ngs amor-<br>ed the po-<br>e means. '<br>e sample<br>n number<br>of the me | bese $\alpha = 0.0$<br>= 0.025,<br>ontrol for p<br>sampling the<br>nction for<br>ng the mea<br>opulation s<br>The Avera<br>per-compa-<br>r of means<br>eans. Note | 05, with t<br>although<br>pair-wise of<br>neory aver<br>the two r<br>ns. To ma<br>tandard of<br>ge Power<br>arison cor<br>are comp<br>that the | he intention<br>we empha-<br>comparison<br>rages of the<br>rules with<br>ake the rest<br>leviation (<br>values in 7<br>rect sign of<br>uted from<br>spread is | on of controlling<br>asize again that<br>is.<br>e sample <i>DFDR</i><br>two numbers of<br>sults comparable<br>denoted here by<br>Table 2 are sam-<br>declaration rate.<br>the same 25,000<br>effectively given |
|---|--|--|---|--|---|--|
| in units equal to   | the standa   | rd errors  | of the sau  | nple mea   | ns For the  | e spread labeled   |
| " $0.00+$ " all popu  | lation mea   | n values   | were set e  | ual. and   | an arbitra  | arv ordering was   |
| chosen to evaluat   | te the "wro  | ng sign"   | errors. Bot   | h proced   | ures conse  | rvatively control  |
| the $DFDR$ for all  | l conditions   | consider   | ed. The Ba  | aves rule r  | orocedure   | provides slightly  |
| greater per-com   | parison pow  | er than t  | hat of Ber  | niamini a  | nd Hochbe   | erg in these con-  |
| ditions, but the a  | actual differ  | ences are  | e trivial.  | -J   |   |  |
| Table 1. Sampling   | Theory DFL<br>and I  | OR for All<br>Benjamini  | Pair-wise C<br>& Hochberg   | omparisons<br>'s Procedui  | Made Usin   | g Our Bayes Rule   |
|   | $\tau$ : Spread  | 10 N   | Aeans   | 25 N   | leans   |  |
|   | of Means   | $\delta_{DFDR}$  | $\delta_{DB \& H}$  | $\delta_{DFDR}$  | δ <sub>DB &amp; H</sub>   |  |
| -   | 0.00+  | 0.0204   | 0.0171  | 0.0206   | 0.0176  |  |
|   | 0.721  | 0.0062   | 0.0044  | 0.0067   | 0.0046  |  |
|   | 3.606  | 0.0005   | 0.0005  | 0.0013   | 0.0012  |  |
|   | 5.408  | 0.0001   | 0.0001  | 0.0006   | 0.0006  |  |
|   | 7.211  | 0.0000   | 0.0000  | 0.0003   | 0.0003  |  |
| -   | 14.422   | 0.0000   | 0.0000  | 0.0000   | 0.0000  |  |
| Table 2. Sampling   | Theory Avera<br>Rule an  | ge Power f<br>id Benjam  | for All Pair-<br>ini & Hochb  | wise Compa<br>erg's Proce  | arisons Mad<br>dure.  | e Using Our Bayes  |
|   | $\tau$ : Spread  | 10 N   | Ieans   | 25 N   | feans   |  |
|   | of Means   | δηέρρ  | δρρ ε. Π  |  | δDR & H   |  |

| $\tau$ : Spread | 10 Means        |                    | 25 Means        |                    |  |
|-----------------|-----------------|--------------------|-----------------|--------------------|--|
| of Means        | $\delta_{DFDR}$ | $\delta_{DB \& H}$ | $\delta_{DFDR}$ | $\delta_{DB \& H}$ |  |
| 0.00+           | 0.002           | 0.002              | 0.001           | 0.000              |  |
| 0.721           | 0.022           | 0.016              | 0.012           | 0.007              |  |
| 3.606           | 0.634           | 0.621              | 0.604           | 0.594              |  |
| 5.408           | 0.783           | 0.778              | 0.741           | 0.737              |  |
| 7.211           | 0.860           | 0.857              | 0.813           | 0.811              |  |
| 14.422          | 0.984           | 0.984              | 0.924           | 0.924              |  |
|                 |                 |                    |                 |                    |  |

Note: 25,000 replications for each condition, n = 3 observations per group, within degrees of freedom  $\nu = 20$  for 10 means and  $\nu = 50$  for 25 means, within variance  $\sigma^2 = 3.0$ ,  $\alpha/2 = .025$ .

### 6. Conclusions

The decision rule  $\delta_{DFDR}$  has been shown to be optimal (from a Bayesian perspective) relative to the loss function  $L_{DFDR}$  for a wide class of multiple comparison problems involving sign declarations. It has also been shown to control a Bayesian version of the directional false discovery rate (*DFDR*), as has a Bayesian version of the procedure proposed by Benjamini and Hochberg ( $\delta_{DB\& H}$ ). There is no guarantee that  $\delta_{DFDR}$  or  $\delta_{DB\& H}$  will control a sampling theory *DFDR* for the case

З

| 1  | of pair-wise comparisons, although that appears to occur in the ANOVA examples   | 1  |
|----|--|----|
| 2  | given, where the two rules behave very similarly.  | 2  |
| 3  |  | 3  |
| 4  | References   | 4  |
| 5  |  | 5  |
| 6  | [1] BENJAMINI Y and HOCHBERG Y (1995) Controlling the false discovery rate: A practical and pow-   | 6  |
| 7  | erful approach to multiple testing. Journal of the Royal Statistical Society, Series B, 57, 289–300.   | 7  |
| 8  | [2] BOX, G. E. P. and TIAO, G. C. (1973). Bayesian inference in statistical analysis. Addison-Wesley,  | 8  |
| 9  | Reading, MA.<br>[3] DUNCAN D B (1965) A Bayesian approach to multiple comparisons. <i>Technometrics</i> <b>7</b> 171–222   | 9  |
| 10 | [4] JONES, L. V. and TUKEY, J. W. (2000). A sensible formulation of the significance test. <i>Psychological</i>  | 10 |
| 11 | Methods, 5, 411–414.   | 11 |
| 12 | <ul> <li>[5] LEHMANN, E. L. (1950). Some principles of the theory of testing hypotheses. The Annals of Math-<br/>ematical Statistics, 21, 1–26.</li> </ul>                     | 12 |
| 13 | [6] LEHMANN, E. L. (1957a). A theory of some multiple decision problems. I. The Annals of Mathemat-  | 13 |
| 14 | ical Statistics, <b>28</b> , 1–25.   | 14 |
| 15 | <ul> <li>[1] LEHMANN, E. L. (1957b). A theory of some multiple decision problems. II. The Annals of Mathe-<br/>matical Statistics, 28, 547–572.</li> </ul>                     | 15 |
| 16 | [8] LEWIS, C. and THAYER, D. T. (2004). A loss function related to the FDR for random effects multiple   | 16 |
| 17 | comparisons. Journal of Statistical Planning and Inference, <b>125</b> , 49–58.  | 17 |
| 18 | [9] SAKAR, S. K. and ZHOU, I. (2008). Controlling Bayes directional false discovery rate in random effects model. Journal of Statistical Planning and Inference, 138, 682–693. | 18 |
| 19 | [10] SHAFFER, J. P. (1999). A semi-Bayesian study of Duncan's Bayesian multiple comparison procedure.  | 19 |
| 20 | Journal of Statistical Planning and Inference, <b>82</b> , 197–213.  | 20 |
| 21 | logical Methods, 7, 356–369.   | 21 |
| 22 | [12] WALLER, R. A. and DUNCAN, D. B. (1969). A Bayes rule for symmetric multiple comparisons prob-   | 22 |
| 23 | lems. Journal of the American Statistical Association, <b>64</b> , 1484–1503.  | 23 |
| 24 | with examples from state-to-state differences in educational achievement. Journal of Educational   | 24 |
| 25 | and Behavioral Statistics, 24, 42–69.  | 25 |
| 26 |  | 26 |
| 27 |  | 27 |
| 28 |  | 28 |
| 29 |  | 29 |
| 30 |  | 30 |
| 31 |  | 31 |
| 32 |  | 32 |
| 33 |  | 33 |
| 25 |  | 25 |
| 36 |  | 36 |
| 37 |  | 37 |
| 38 |  | 38 |
| 39 |  | 39 |
| 40 |  | 40 |
| 41 |  | 41 |
| 42 |  | 42 |
| 43 |  | 43 |
| 44 |  | 44 |
| 45 |  | 45 |
| 46 |  | 46 |
| 47 |  | 47 |
| 48 |  | 48 |
| 49 |  | 49 |
| 50 |  | 50 |
| 51 |  | 51 |
|    |  |    |

# The Challenges of Model Objective Selection and Estimation for Ad-hoc Network Data Sets

З

# Farinaz Koushanfar<sup>1,2</sup> and Davood Shamsi<sup>2</sup>

#### $Rice \ University$

Abstract: We introduce a new methodology for determining the difficulty of selecting the modeling objective function and estimating the parameters for an ad-hoc network data set. The method utilizes formulation of the underlying optimization problem instance that consists of an objective function and a set of constraints. The method is illustrated on real distance measurement data used for estimating the locations of wireless nodes that is the most studied and a representative problem for ad-hoc networks estimation. The properties of the data set that could affect the quality of optimization are categorized. In large optimization problems with multiple properties (characteristics) that contribute to the solution quality, it is practically impossible to analytically study the effect of each property. A number of metrics for evaluating the effectiveness of the optimization on each data set are proposed. Using the well known Plackett and Burmann fast simulation methodology, for each metric, the impact of the categorized properties of the data are determined for the specified optimization. A new approach for combining the impacts resulting from different properties on various metrics is described. We emphasize that the method is generic and has the potential to be more broadly applicable to other parameter estimation problems.

#### Contents

| 1 | Introduction 334                                    |
|---|---|
| 2 | Preliminaries 335                                   |
| 3 | Metrics 337   |
| 0 | 3.1 Error Matrice 337                               |
|   | 2.2 Objective Function (OF) Matrice                 |
|   | 2.2.1 Drifting of Objective Function (OF) $(OF)$    |
|   | 3.2.1 Dritting of Objective Function (OF)           |
|   | 3.2.2 Nearest Local Minimum                         |
|   | 3.2.3 Measuring the Slope of OF Around the Solution |
|   | 3.2.4 Depth of the Non-Global Local Minima          |
| 4 | Simulation Methodology                              |
| 5 | Combining Different Ranks                           |
| 6 | Evaluation Results                                  |
| 7 | Conclusion  |
| A | cknowledgement                                      |
| R | eferences   |

<sup>2</sup>Electrical and Computer Engineering Department, Rice University, Houston, TX 77005 <sup>2</sup>Computer Science Department, Rice University, Houston, TX 77005

imsart-coll ver. 2008/08/29 file: Koushanfar.tex date: April 10, 2009

#### 1. Introduction

Wireless adhoc networks consist of multiple wireless nodes distributed in an implementation area. To be power efficient, the wireless nodes only directly communicate with the nodes in their short local range (neighbor nodes). Communication between the non-neighbor nodes is enabled by successive usage of (one or more) local neighbors as forwarding relays. Several problems in this domain include modeling and estimation of data sets that only contain pairwise exchanged data between the neighboring nodes.

Years of continuous research in building statistical models and parameter esti-mation has produced a multitude of readily available methods and tools that can be employed for the problems in ad-hoc networks [10]. One limitation of the available methods is that majority of the ad-hoc modeling and estimation problems concern a large body of data and do not conform with typical assumptions needed to an-alytically declare the known theoretical optimality criteria. In such scenarios, the quality of the modeling and estimation methods are typically evaluated by how they perform on sets of real or simulated data. For example, some statistics of the result-ing prediction error and/or a defined criterion (e.g., Bayesian information criterion (BIC)) is used for experimental evaluation of the method on the adhoc network measurements. A relevant question to answer is if indeed modeling and estimation of the pertinent data set requires introduction of a new model or an estimator, or the data could have been just as well addressed by the other known methods. 

Our objective is to quantify the difficulty of model selection and estimation for a given adhoc network data set. This would provide impetus for inventing newer modeling and estimation objectives and tools that can address the difficult-tocharacterize data. Simultaneously, formation of new tools would depend upon find-ing truly challenging network data sets that need to be addressed, as opposed to building new models that have a limited practical usage. Devising sets of challenging data would also build a foundation for comparing the various modeling objective functions and estimators for the ad-hoc network data sets. The problem of finding challenging data is complicated by variations in properties of the underlying data sets collected by different sources. This includes difference in size, format, wireless ranges, hidden covariates, and the form of noise present in the collected data. Thus, it is not easy to find unique metrics that could be used for comparison of different modeling objective functions and estimation methods. 

In statistics literature, sensitivity of estimation error or other discrepancy metrics to the underlying noise in data has been widely studied for a number of modeling methods [24][3]. Also, the consistency of estimators based on a number of strong assumptions on the distribution of the data has been pursued [14]. However, no generic method or tool for determining the difficulty in modeling a data set free of imposing strong assumptions – such as normality or other closed-form distributions of noise – is available for use in adhoc networks. Note that the runtime complexity of a problem is an orthogonal concept. The complexity measures the worst-case computational time for the algorithm used for addressing the problem. Analyzing the worst-case runtime complexity does not help in understanding the complexity of characterizing a specific data set. 

In adhoc network scenario, after the variable selection is done and the noise
models are assumed, modeling is typically done by selecting a model form (e.g.,
nonlinear regression) and then estimating the model parameters on the data set.
For analyzing the modeling objective function and estimation performance on the
data, we study the pertinent optimization problem that consists of an objective

function (OF) and a number of constraints. The data set is considered as the input to the optimization problem. We introduce a number of metrics that measure the complexity of the optimization problem based on the problem OF properties and constraints. The challenge in most optimization problems is the existence of nonlinearities that make the solution space coarse, causing bumpiness and multiple local minimums. We propose a number of measures for the smoothness of the OF and constraints space that estimate the feasibility of reaching the global minimum.

To enable studying the effectiveness of the optimization on an adhoc network data set, one should characterize the properties of the pertinent data set. The properties are specific to each data set and the problem. In this article, we focus on the problem of finding the location of nodes (localization) in an adhoc wireless network by using erroneous mutual distance measurements between a number of node pairs. However, we emphasize that our method is generic and can be used for determining the challenge in addressing many adhoc data set model objective selection and estimation that includes forming an optimization problem. The lo-calization problem is selected for four reasons. First, it is a very well addressed problem in the literature and there are several methods that are developed for this problem [6][9][19][2][20][26]. Second, there are a number of publicly available data sets for the measured distance data in the networks [5][21][8]. Third, the nonlinear relationship between noise in measurements data and the location of nodes makes the modeling problem extremely challenging. Fourth, localization problem is an NP-complete problem, i.e., in the worst case, there is no algorithm that can solve it in polynomial time [25][6]. Lastly, location discovery is a precursor for a number of other problems in ad hoc networks including sleeping coordination [12, 13], sensor coverage [15], and sensing exposure [16].

We characterize a number of properties of the measurement data set that could affect the quality of location estimation. Studying the interaction between the iden-tified data properties and optimization metrics requires long simulations and anal-ysis. We use the well-known Plackett and Burmann [23] simulation methodology to rapidly study the pairwise linear interactions of properties. A new approach for combining the impacts resulting from different properties of data on various opti-mization metrics is described. The sensitivity of optimization with respect to the various parameter ranks are presented. 

To the best of our knowledge, this is the first work that systematically studies the impact of the adhoc network data set on the optimization employed for finding the modeling objectives and estimations. Most of the previous work are devoted to modeling and analysis of the worst case complexity. The results of our analysis could be directly used for constructing benchmarks for the problem. The proposed work aims at creating a unified framework based on real data that can help evaluation and comparison of desperate efforts that address the same problem.

The remainder of the paper is organized a follows. In the next section, location estimation problem and our notations are formally defined. In Section 3, we devise a number of metrics that are used for OF evaluation. The simulation methodology is described in Section 4. In Section 5, we illustrate how the results of different metrics can be combined. We have applied the derived method on the measurements from a real network in Section 6. We conclude in Section 7.

#### 2. Preliminaries

In this section, we present the formal definition of the problem. We also describe the notations that are used throughout the paper. З

**Location estimation problem:** Given a set of N nodes denoted by  $V = \{v_1, v_2, \ldots, v_N\}$  in  $\mathbb{R}^d$  (d = 2, 3). For a given subset of node pairs denoted by  $E \subset V \times V$ , mutual distance of nodes are measured, *i.e.*, for all  $(v_i, v_j) \in E$ ,  $l(v_i, v_j) = d(v_i, v_j) + \epsilon_{i,j}$  is known;  $d(v_i, v_j)$  is the Euclidean distance between the nodes  $v_i$  and  $v_j$ ;  $\epsilon_{i,j}$  is the distance measurement error. This error is only known if the real and measured location are both available. Moreover, there is a subset with M(> 2) nodes denoted by  $V_B = \{v_1, \ldots, v_M\}$ ,  $V_B \subset V$  such that the nodes in  $V_B$  have their exact location information (coordinates). The nodes in the set  $V_B$  are called the *beacon* nodes.

**Question**: find the location of all possible nodes.

In this paper, we focus on two-dimensional networks. Extension to threedimensional networks is straight forward. Coordinates of the node  $v_i$  are denoted by  $(x_i, y_i)$ .

The location estimation problem can be formulated as an optimization problem. The goal is to find the coordinates of K = N - M non-beacon nodes such that the discrepancy (error) between the measured distance data and the nodes' distances estimated from the final coordinates is minimized. In other words,

(1) 
$$F_L(x_{M+1}, y_{M+1}, x_{M+2}, y_{M+2}, \dots, x_N, y_N) = \sum_{(v_i, v_j) \in E} L(e_{v_i, v_j})$$

$$e_{v_i,v_j} = l(v_i, v_j) - \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$

Where  $L : \mathbb{R} \to \mathbb{R}^+$  is a function that is typically a metric (measure) of error.  $F_L : \mathbb{R}^{2K} \to \mathbb{R}^+$  is known as objective function (OF) of the optimization problem. Note that the OF of the location estimation problem is not necessarily a linear or convex function. There are a number of fast and efficient tools that are developed for linear and convex programming. However, there is no oracle algorithm that can solve all optimization problems. To find the minimum of a nonlinear problem like location estimation, there are a number of heuristic methods that may be employed. The nonlinear system solvers have a tendency to get trapped in a local minimum and do not necessarily lead to the global minimum. Although there are a variety of minimization algorithms, most of them are common in one subcomponent that starts from an initial point and follow the steepest decent to reach the minimum.

The algorithms differ in how they choose the starting point, how they select the direction in the search space, and how they avoid local (non-global) minima. Thus, the shape of the OF around the global minimum is an important factor in finding the solution.

**Data set:** The measurement data used in this problem consists of measured distances between a number of static nodes in the plane. Measurements are noisy; there are multiple measurements for each distance. The true location of the nodes is known and will be known as the ground truth. As explained in Section 1, we sample the data set to obtain instances with specific properties.

Parameters: We will define a number of parameters that can be extracted from
 the data set. The sensitivity of the location estimation to the variations in each
 parameter will be studied. The analysis results will be used for identifying the hard
 instances of measurement data. Ten parameters are studied:

•  $P_1$  – Number of nodes (N): the total number of nodes in the network.

•  $P_2$  – Number of beacons (B): the number of beacon nodes with known locations.

| 1  | • $P_3$ – Mean squared error $(\overline{\epsilon^2})$ : mean squared error of distance measurements.          | 1  |
|----|--|----|
| 2  | • $P_4$ – Maximum allowed squared error (MAX <sub><math>\epsilon_m^2</math></sub> ): the maximum squared error | 2  |
| 3  | that can possibly exist in distance measurements.  | 3  |
| 4  | • $P_5$ – Percent of large errors $(PER_{\epsilon_2})$ : percentage of squared distance mea-                   | 4  |
| 5  | surement noises that are higher than a specific value $\epsilon_0^2$ .   | 5  |
| 6  | • $P_6$ – Mean degree $(\overline{D})$ : mean degree of the nodes in the network. Degree of a                  | 6  |
| 7  | node $v_i$ is define as number of nodes that have their mutual distance to $v_i$ .                             | 7  |
| 8  | • $P_7$ – Minimum length (MINL): possible minimum length of the measured                                       | 8  |
| 9  | distances between nodes in the network.  | 9  |
| 10 | • $P_8$ – Maximum length (MAXL): possible maximum length of the measured                                       | 10 |
| 11 | distances between nodes in the network.  | 11 |
| 12 | • $P_9$ – Mean length ( $\overline{l}$ ): mean length of the measured distances between nodes                  | 12 |
| 13 | in the network.  | 13 |
| 14 | • $P_{10}$ – <i>Minimum degree</i> (MIND): possible minimum degree of the nodes in the                         | 14 |
| 15 | network.   | 15 |
| 16 | To study the effect of the parameters, we construct a variety of network instances                             | 16 |
| 17 | with different properties. The networks are constructed by selecting subsets of an                             | 17 |
| 18 | implemented network. Having specific values for parameters, we use Integer Linear                              | 18 |
| 19 | Programming (ILP) to extract each subset such that it meets specified conditions                               | 19 |
| 20 | To do so, we model parameter constraints as linear equalities and inequalities                                 | 20 |
| 21 | Some parameters such as the mean squared error $\overline{\epsilon^2}$ can be easily stated by linear          | 21 |
| 22 | equalities and inequalities. But some parameters such as the mean degree of the                                | 22 |
| 23 | nodes $\overline{D}$ need a mapping to be stated in linear terms. The description of the exact                 | 23 |
| 24 | procedure of modeling by linear constraints is beyond the scop of this paper [8]                               | 24 |
| 0E | procedure of modeling by model constraints is segond the scop of this paper [0].                               | 25 |

#### 3. Metrics

In this section, we introduce metrics for error and OF that are used for evaluating the importance of different parameters for location estimation. Three error metrics and four OF metrics are presented. Thus, a total of twelve combined metrics are used to evaluate the importance of parameters.

#### 3.1. Error Metrics

The three error metrics studied in this paper are:  $L_1$ ,  $L_2$ , and the maximum likelihood (ML).  $L_1$  and  $L_2$  are the common error norms in the  $L_p$  family defined as:

$$L_p(e_{v_n,v_m} \in E) = (\sum_{(v_n,v_m)\in E} |e_{v_n,v_m}|^p)^{1/p} \quad if \quad 1 \le p < \infty.$$

To find the error metric corresponding to ML, we need to model the noise in distance measurements. To model the noise, the probability density function (PDF) of errors,  $f_m$ , for the distance measurements should be approximated. Different methods are developed to approximate PDF of noise,  $f_m$  [8]. We have used kernel fitting that is a simple and known PDF approximation method [10]. To have the maximum likelihood estimation for the nodes' locations, we find the nodes' coordinates such that they maximize

(2) 
$$\prod_{(v_n,v_m)\in E} f_m(e_{v_n,v_m}) = \exp\{\sum_{(v_n,v_m)\in E} \ln(f_m(e_{v_n,v_m}))\}$$

imsart-coll ver. 2008/08/29 file: Koushanfar.tex date: April 10, 2009


or equivalently minimize

(3) 
$$\sum_{(v_n,v_m)\in E} -\ln(f_m(e_{v_n,v_m})).$$

Note that we assume noise in distance measurements are independently identically distributed. Using the same notations as the Equation 1 and Equation 3, for the ML estimation we consider the following error metric:

(4) 
$$L_{ML}(e_{v_n,v_m}) = -\ln(f_m(e_{v_n,v_m})).$$

З

# 3.2. Objective Function (OF) Metrics

We describe metrics that are used for evaluating OFs. The metrics are introduced based on the properties of OF that are effective in optimization. These metrics are such that they assign larger values to the more difficult-to-optimize OFs. For example, if one selects a convex OF, it may be possible to utilize convex programming depending on the form of the constraints. In defining the OF metrics, we assume that there is a fixed instance of location estimation data. Thus, for a fixed error metric, the OF would be fixed. Metrics of OF are denoted by  $M : \mathcal{C} \to \mathbb{R}^+$  where  $\mathcal{C}$  is the functional space that contains all OFs.

# 3.2.1. Drifting of Objective Function (OF)

47 Since there is noise in distance measurements, true location of the nodes is 48 often not the global minimum of the OF. Location of the OF's global minimum is 49 a measure of the goodness of the OF. Figure 1 illustrates the effect of noise on the 50 OF. For the sake of presentation simplicity, an one-dimensional OF is shown. In 51 this figure,  $p_c$  is the correct nodes' location. However, the global minimum of the OF is displaced at  $p_{gm}$  because of the noise. We consider the distance between  $p_c$ and its displaced location  $p_{gm}$  as an OF metric and denote it by *drifting*. To find the drifting distance, we start from the true locations as the initial point. Next, the steepest descent direction of the OF is followed until a local minimum is reached. The Euclidean distance between the true locations and this local minimum quantifies the drifting metric (denoted by  $M_1$ ) for the pertinent OF.

3.2.2. Nearest Local Minimum

Having a number of local minimums around the global minimum in an OF may cause the optimization algorithm to get trapped in one of the non-global local minimums. It is challenging to minimize such an OF since the global minimum is hard to reach. Figure 1 illustrates the phenomena. The OF has multiple local minima at points  $p_{m1}$ ,  $p_{m1}$  and so on. The steepest decent method leads to the global minimum if and only if we start from a point between  $p_{m1}$  and  $p_{m2}$ . Hence, having a small distance between  $p_{m1}$  and  $p_{m2}$  would complicate the selection of the initial starting point.

We introduce a method to measure the distance of the true locations from the local minimums around the global minimum. Because of curse of dimensionality, it is not possible to find all the local minimums around the global minimum. We randomly sample the OF in multiple directions. The nearest local minimum is computed for each randomly selected direction. We statistically find the distance to the nearest local minimum by using multiple samples.

Assume  $F : \mathbb{R}^{2K} \to \mathbb{R}^+$  is the OF. A random direction in  $\mathbb{R}^{2K}$  is a vector in this space. Let us denote it by  $v \in \mathbb{R}^{2K}$ . First, we define a new function  $h : \mathbb{R}^+ \to \mathbb{R}^+$ such that  $h(t) = F(p_c + tv)$  where  $p_c$  is a vector containing the true locations of nodes. Second, we find the local minimum of h with the smallest positive t and denote it by  $t_1$ . We repeat this procedure for T times and find all  $t_i$ 's. T is the number of samples. Finally, since it is expected that the defined metric has a larger value for more difficult-to-optimize OF, we define the nearest local minimum metric to be

(5) 
$$M_2(F) = \left(\frac{1}{T}\sum_{i=1}^T t_i\right)^{-1}$$

# 3.2.3. Measuring the Slope of OF Around the Solution

The Slope of OF (*i.e.*, the norm of OF's gradient) around the global minimum is a very important parameter in the convergence rate of the optimization algorithm. OFs with a small slope around the true location converge to the global minimum very slowly.

Thus, measuring the slope of the OF around the global minimum can be used to quantify the goodness of OF. Again, we measure slope of the OF in multiple random directions around the true locations, and statistically compute this metric. OFs with sharp slopes around the global minimum are easier to optimize. This can be seen in Figure 2 where the right side of the global minimum,  $p_{gm}$ , has a sharp slope. If the initial point of steepest descent algorithm is between  $p_{gm}$  and  $p_{m2}$ , it converges to the global minimum very fast. However, on the left side of global З

minimum,  $p_{gm}$ , there is a gradual slope. Thus, the steepest descent algorithm would converge very slowly on the left side. We define the true locations' slope metric as

(6) 
$$M_3(F) = \left(\frac{1}{T}\sum_{i=1}^T \text{slope in i-th random direction}\right)^{-1}$$

Note that the slope of the i-th random direction,  $v_i$ , is measured at  $p_{gm} + \sigma v_i$ where  $\sigma$  is a small number and is a user's defined criterion.

#### 3.2.4. Depth of the Non-Global Local Minima

Optimization problems that have an OF with deep local minimums around the global minimum are difficult to solve. A number of heuristic optimization methods take advantage of the shallow local minimums to avoid non-global local minimums, e.g., simulated annealing [11]. In figure 2, avoiding the local minimum at  $p_{l1}$  is much easier than local minimum at  $p_{l2}$ .

We define the forth metric for quantifying the goodness of an OF on the data, as the depth of the non-global local minimums. We randomly select T local minimums around the true locations. Assuming that  $m_i$  is the OF value at the randomly selected local minimums, define

(7) 
$$M_4(F) = \left(\frac{1}{T}\sum_{i=1}^T m_i\right)^{-1}.$$

4. Simulation Methodology

We find the linear effect of each parameter by studying all combinations of parameters. Assume each parameter has just two values. If we have k parameters then we have to study  $2^k$  combinations that is computationally intractable. Instead, we use Plackett and Burman (PB) [23] fast simulation methodology that is a very well known method for reducing the number of simulations. Number of simulation in PB is proportional to the number of parameters. Although the PB method has not been used for the adhoc modeling and estimation problems, it was used for the simulations speedup in a number of other adhoc network problems [1][22][27][28].

In PB design, two values are assigned to each parameter: a normal value and an extreme value. The normal value is the typical value of the parameter while the extreme value is the value that is outside the typical range of the parameter. The extreme value often makes the problem either harder or easier to solve. A number of experiments with normal and extreme values of parameters are conducted.

Experiments are arranged based on a given matrix denoted by the *design matrix*. Design matrix has k columns (k is the number of parameters) and s rows where s is the number of experiments the should be set up as follows. The elements of the design matrix are either 0 or 1. We set up an experiment for each row. Values of the parameters depend on the elements on the row: 0 indicates that the normal value of the parameter is used and 1 indicates that the extreme value of the parameter is used in the experiment corresponding to the row.

Assume that we have selected an error metric,  $L_i$ , and an objective function metric,  $M_j$ . The OF itself denoted by  $F_{L_i}$  would be fixed. For each row of the

design matrix, h, we set up an experiment based on the elements of that row and measure the goodness of the objective function  $M_i(F_{L_i})$  and save it in another array element denoted by  $r_{i,j,h}$ . The corresponding values are summed up for computing the importance factor (IF) of each parameter. For each parameter  $P_t$ , we define

(8) 
$$\operatorname{IF}_{t,i,j} = |\sum_{k=1}^{s} \alpha_{h,t}|$$

 $tr_{i,i,h}$ h=1

where s is the number of experiments (number of rows in the design matrix), and  $\alpha_{h,t}$  is 1 if the extreme value of the parameter  $P_t$  is used in the h-th experiment; otherwise,  $\alpha_{h,t}$  is -1. The absolute value of IF is used to evaluate the effect of each parameter. The largest value indicates the most important parameter. For *i*-th error metric and j-th OF metric,  $IF_{t,i,j} > IF_{u,i,j}$  means that the parameter  $P_t$  is more important than  $P_u$ . Thus, for each error metric,  $L_i$ , and for each objective function metric,  $M_i$ , we can rank parameters based on their effect on the estimated location. This ranking is denoted by  $R_{i,j}$ .

More precise results can be obtained by using the foldover design matrix [18]. In the foldover design matrix, all rows of the single design are repeated after its last row but 0s and 1s are exchanged in the repeated rows.

#### 

# 5. Combining Different Ranks

In this section, we explain how to combine the rankings of the parameters under study to obtain a global order for them. Using the ranking method in the previous section, we would have different rankings for various error metrics and OF metrics. Since there are three error metrics and four objective function metrics, there would be twelve different ranking lists for the importance of parameters; each parameter may have a different rank in each ranking list.

Each rank is obtained based on a specific property of the optimization problem. As it is explained in Section 3, for each error and objective function metric, the parameters are ranked based on the importance factor obtained from PB-design. IFs with large discrepancies lead to a stronger ranking compared to IFs with small discrepancies. Simply summing up the rankings would not necessarily determine which of the importance factors were better differentiating among the parameters. For each ranking,  $R_{i,j}$ , and for each pair of parameters,  $P_s$ ,  $P_t$ , we find the

probability that  $P_s$  is more important than  $P_t$ . Based on the probabilities, we construct the global ranking.

Consider a specific error metric,  $L_i$ , and a specific objective function metric,  $M_j$ . Assume that the importance factor of the parameter  $P_t$ , IF<sub>t,i,j</sub>, is normally distributed  $\mathcal{N}(\lambda_{t,i,j},\sigma^2)$ . The observed value of IF<sub>t,i,j</sub> in a specific experiment is denoted by  $if_{t,i,j}$ . We normalize the importance factors to have a maximum value W. The mean of IFs are assumed to be uniformly distributed in [0, W].

For each two parameters,  $P_s$  and  $P_t$ , given the BP-design experiment importance values  $if_{s,i,j}$ , and  $if_{t,i,j}$ , we find the probability:  $Pr(\lambda_{s,i,j} \geq \lambda_{t,i,j} | IF_{s,i,j} =$  $if_{s,i,j}, IF_{t,i,j} = if_{t,i,j}$ ). The conditional probability can be written in the Bayesian format as

$$\beta_{s,t,i,j} = Pr(\lambda_{s,i,j} \ge \lambda_{t,i,j} | IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j}) =$$

(9) 
$$\frac{Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j} | \lambda_{s,i,j} \ge \lambda_{t,i,j}) Pr(\lambda_{s,i,j} \ge \lambda_{t,i,j})}{Pr(IF_{s,i,j} = if_{s,i,j}, IF_{t,i,j} = if_{t,i,j})}.$$

З

 $Koushanfar \ and \ Shamsi$ 

|  |                                  | 1                                  |  |   | 1  |                                     |  |  |  |  |
|--|----------------------------------|------------------------------------|--|---|--|-------------------------------------|--|--|--|--|
| Parameter  | $N_S$                            | $B_S$                              | $\epsilon_S^2$                                   | $^{MAX}\epsilon_m^2$  | $PER_{\epsilon_0^2}$   | $\overline{D_S}$                    | $MINL_S$                                   | MAXLS  | $\overline{l_S}$                         | $MIND_S$   |
| Normal<br>Value  | 55                               | 12                                 | $10 \ (m^2)$                                     | $200 \ (m^2)$   | 50   | 10                                  | 5(m)                                       | 40 (m)   | 20 (m)                                   | 4  |
| Extreme<br>Value   | 80                               | 3                                  | $50 \ (m^2)$                                     | $500 \ (m^2)$   | 20   | 6                                   | 10 (m)                                     | 60 (m)   | 30 (m)                                   | 3  |
| Since the sume the second seco | here i<br>at Pa                  | s no j $r(\lambda_{s,i})$          | Normal of prior info $_{,j} \ge \lambda_{t,i,j}$ | and extrem prmation $j = \frac{1}{2}$ . For                   | TABLE 1<br>e values for<br>about the<br>urthermo   | or the<br>e dist<br>pre,            | paramete<br>ributior                       | rs.<br>as of $\lambda_{s,i,j}$                           | $_{j} \text{ and } \lambda_{j}$          | t,i,j, we  |
|  | ·                                | ( -,-                              | ,,,,   |   |  | ,                                   |  |  |  |  |
| Pr(  | $(IF_{s,i})$                     | $j_{j} = i$                        | $f_{s,i,j}, IF$                                  | $v_{t,i,j} = v_{t,i}$   | $\lambda_{s,i,j}$  | $\geq \lambda_{t,i}$                | $_{,j}) =$                                 |  |  |  |
| ſ  | ſ                                | Pr(                                | $(IF_{s,i,j} =$                                  | $= i f_{s,i,j}, I$  | $F_{t,i,j} = i$  | $f_{t,i,i}$                         | $\lambda_{s,i,j} =$                        | $y, \lambda_{t,i,j} =$                                   | $= x)\frac{dy}{dx}$                      | $\frac{dx}{dx} =$                                  |
| $J_{x=}$   | $J_{y=}$                         | x                                  | ,•,J   | • -,•, <b>J</b> /   |  |                                     | -,-,J                                      | ,• ,J  | ´ W                                      | N  |
| (10) <u> </u>  | $\int_{-}^{W}$                   | $\int_{0}^{W}$                     |  | $e^{\frac{(y-if_{s,i,j})}{2\sigma^2}}$                        | <sup>2</sup> <u>1</u>  | $e^{\frac{(x-i)}{2}}$               | $\frac{f_{t,i,j})^2}{2\sigma^2} d\eta$     | dx.  |  |  |
| $W^{2}$  | $J_{x=0}$                        | $\int_{y=x}$                       | $v \sqrt{2\pi\sigma^2}$                          | 0   | $\sqrt{2\pi\sigma^2}$  |                                     | a.   | au.  |  |  |
| Similar  | lv. on                           | ie car                             | n find   |   |  |                                     |  |  |  |  |
|  | , 01                             |                                    |  |   |  |                                     |  |  |  |  |
| Pr(II  | $F_{s,i,j}$                      | $= i f_s$                          | $_{,i,j}, IF_{t,i}$                              | $j_{j,j} = i f_{t,i,j}$                                       | ;)   |                                     |  |  |  |  |
|  | = Pr                             | $(IF_{s})$                         | $i,j = if_s$                                     | $_{i,j}, IF_{t,i,j}$  | $= i f_{t,i,j}$  | $ \lambda_{s,i,j} $                 | $\geq \lambda_{t,i,j}$                     | $)Pr(\lambda_{s,i})$                                     | $j \ge \lambda_{t,i}$                    | ,j)  |
|  | + Pr                             | $(IF_{s,s})$                       | $i,j = if_{s,j}$                                 | $_{i,j}, IF_{t,i,j}$  | $= i f_{t,i,j}$  | $ \lambda_{s,i,j} $                 | $<\lambda_{t,i,j}$                         | $)Pr(\lambda_{s,i})$                                     | $_{j} < \lambda_{t,i}$                   | $_{j}).$   |
| NT C   |                                  | 1                                  | ,  | י ת   | C 11   | 1 1                                 | 1.   |  |  |  |
| Now, fo  | or eac                           | en pai                             | rameter,   | $P_t$ , we de   | enne the   | globa                               | 1 impor                                    | tance fac  | tor, $if_t$                              | ,  |
|  |                                  |                                    |  |   |  |                                     |  |  |  |  |
| (11)   |                                  |                                    |  | $f_{i} = \sum_{i=1}^{N_{em}} \sum_{i=1}^{I}$                  | $N_{om}$ $N_p$   | R                                   |  |  |  |  |
| (11)   |                                  |                                    | ı,   | $J_t = \sum_{i=1}^{t} J_i$                                    | $\sum_{i=1}^{\infty} \sum_{s=1}^{\infty} s^{s-1} s^{s$ | , $\rho_{s,t}$<br>$\neq_t$          | $, i, j \cdot$                             |  |  |  |
|  |                                  |                                    |  | <i>v</i> —1 .   | /=1 0=1,0  | 7-0                                 |  |  |  |  |
| Parame   | eters                            | with                               | a larger   | $if_t$ have a   | a higher   | proba                               | bility o                                   | f being in   | nportar                                  | nt com-  |
| bared to the formula $f$   | ne otl                           | ner pa                             | arameter   | s. We sort  | the para   | amete                               | rs basec                                   | l on their   | corresp                                  | onding   |
| Jt values.   |                                  |                                    |  |   |  |                                     |  |  |  |  |
| а <b>г</b> а ч   |                                  | ъ                                  | 14   |   |  |                                     |  |  |  |  |
| b. Evalu   | atior                            | ı Res                              | sults  |   |  |                                     |  |  |  |  |
| We have a  | applie                           | ed the                             | e develor  | ed metho  | od to rea  | al dist                             | ance m                                     | easureme   | ent data                                 | for lo-  |
| cation est   | imati                            | on pr                              | oblem. F   | arameter  | s that w   | ere de                              | scribed                                    | in Sectio  | $n \frac{2}{2}$ are                      | ranked   |
| using our  | meth                             | odolc                              | ogy. We i  | llustrate l   | how the  | variou                              | ıs ranki                                   | ng lists d   | iffer. Tl                                | nen, we  |
| ombine t   | he ra                            | nking                              | s to obt   | ain a glob  | al ranki   | ng.                                 |  |  |  |  |
| The dis  | stance                           | e mea                              | asuremen   | its data fr   | om the   | CENS                                | 5 lab [4]                                  | is used t  | to evalu                                 | ate the  |
| ffect of ea  | ach p                            | aram                               | eter. Thi  | s databas   | e is base  | d on t                              | the real                                   | distance   | measur                                   | ements   |
| or SH4 n   | -                                | F - 7                              |  | -   |  |                                     |  |  |  |  |
| 1  | odes                             | [ <b>8</b> ]. 9                    | 1 nodes  | are locat   | ed in fix  | ed loo                              | cations.                                   | Distance   | e measu                                  | rement   |
| s done m   | odes<br>ultipl                   | [8]. 9<br>e tim                    | 1 nodes<br>es and ir                             | are located different $\begin{bmatrix} 17 \\ 2 \end{bmatrix}$ | ed in fix<br>days. T   | ed loo<br>he dis                    | cations.<br>stance n                       | Distance<br>neasurem                                     | e measu<br>ients are                     | e based  |
| s done my  | odes<br>ultipl<br>ne of          | [8]. 9<br>e tim<br>fligh           | 1 nodes<br>es and ir<br>t (ToF)[                 | are locat<br>different<br>[17] of the                         | ed in fix<br>days. T<br>e signals  | ed loo<br>he dis<br>. In t          | cations.<br>stance n<br>his met            | Distance<br>neasurem<br>hod, the                         | e measu<br>ients are<br>time o           | rement<br>e based<br>f flight                      |
| s done my<br>on the tim<br>of an acou  | odes<br>ultipl<br>ne of<br>ustic | [8]. 9<br>e tim<br>fligh<br>signal | 1 nodes<br>es and ir<br>t (ToF)[<br>l is used    | are locat<br>a different<br>17] of the<br>to determ           | ed in fix<br>days. T<br>e signals<br>nine the  | ed loo<br>he dis<br>. In t<br>dista | cations.<br>stance n<br>his met<br>nce bet | Distance<br>neasurem<br>hod, the<br>ween two<br>strongly | e measu<br>ients arc<br>time o<br>nodes. | rement<br>e based<br>f flight<br>It was<br>tic [7] |

Therefore, parametric methods based on optimizing the results according to a fixed
 noise distribution do not yield good location estimations.

Challenges of Model Objective Selection and Estimation

ML

SMAS

 $\frac{L_2}{2}$ 

ML

NLM

DOI

ML

Parameter

1



imsart-coll ver. 2008/08/29 file: Koushanfar.tex date: April 10, 2009

343

1

DNGLM

 $L_2$ 

ML

| Parameter                  | $N_S$ | $B_S$ | $\epsilon_S^2$ | $^{MAX}\epsilon_m^2$ | $PER_{\epsilon_0^2}$ | $\overline{D_S}$ | $\operatorname{MINL}_S$ | $MAXL_S$ | $\overline{l_S}$ | $\operatorname{MIND}_S$ |
|----------------------------|-------|-------|----------------|----------------------|----------------------|------------------|-------------------------|----------|------------------|-------------------------|
| $N_S$                      | 0     | 0.071 | 0.417          | 0.725                | 0.716                | 0.403            | 0.708                   | 0.691    | 0.598            | 0.748                   |
| BS                         | 0.929 | 0     | 0.899          | 0.993                | 0.984                | 0.884            | 0.977                   | 0.981    | 0.939            | 0.984                   |
| $\overline{\epsilon_c^2}$  | 0.583 | 0.101 | 0              | 0.787                | 0.786                | 0.476            | 0.756                   | 0.754    | 0.660            | 0.798                   |
| MAX <sub>2</sub>           | 0.275 | 0.007 | 0.213          | 0                    | 0.490                | 0.193            | 0.464                   | 0.477    | 0.354            | 0.515                   |
| $PER_{\epsilon_0^2}^{c_m}$ | 0.284 | 0.016 | 0.214          | 0.510                | 0                    | 0.202            | 0.469                   | 0.499    | 0.357            | 0.528                   |
| DS                         | 0.597 | 0.116 | 0.524          | 0.807                | 0.798                | 0                | 0.785                   | 0.795    | 0.678            | 0.821                   |
| MINLS                      | 0.292 | 0.023 | 0.244          | 0.536                | 0.531                | 0.215            | 0                       | 0.519    | 0.392            | 0.545                   |
| MAXLS                      | 0.309 | 0.019 | 0.246          | 0.523                | 0.501                | 0.205            | 0.481                   | 0        | 0.371            | 0.537                   |
| Is ~                       | 0.402 | 0.061 | 0.340          | 0.646                | 0.643                | 0.322            | 0.608                   | 0.629    | 0                | 0.671                   |
| MINDS                      | 0.252 | 0.016 | 0.202          | 0.485                | 0.472                | 0.179            | 0.455                   | 0.463    | 0.329            | 0                       |

#### TABLE 3

Drifting of objective function and  $L_2$  metric:  $Pr(\lambda_{i,j,s} \ge \lambda_{i,j,t} | V_{i,j,s} = v_{i,j,s}, V_{i,j,t} = v_{i,j,t})$ where the first column is  $P_s$  and the first row is  $P_t$ .

The comparative ranks of parameter pairs tend to vary as well. Figure 2 shows the normalized importance factor (IF) for two cases: DOF and SMAS with  $L_2$  error metric. For DOF, the number of beacons  $B_S$  is strongly more important than the mean squared error  $\epsilon_S^2$ . The mean degree of nodes,  $\overline{D}_S$ , is weakly more important than the mean squared error  $\overline{\epsilon}_S^2$ . The same behavior can be seen in SMAS. From our visual inspections, the number of nodes  $N_S$  and the mean degree of nodes  $\overline{D}_S$ are the most important while others almost have the same importance factor (IF). The ranks of the mean squared error  $\overline{\epsilon}_S^2$  and maximum edge length MAXL<sub>S</sub> are 3 and 10 respectively. However, their importance factors are very close.

The discrepancy in the rank and comparative ranks confirms our postulation that averaging the parameter ranks is not the best way for combining them. Thus, we use the combining method that was introduced in Section 5. The probability comparisons for the values in Figure 2 are shown in Tables 3 and 4. The tables compare the importance of parameters. For example, for the DOF- $L_2$ , Figure 2 states that  $B_S$  is strongly more important than  $PER_{\epsilon_0^2}$ . Table 3 shows that the probability that the mean of  $B_S$  is larger than the mean of  $PER_{\epsilon_0^2}$  is 0.984. Similarly,  $MAX_{\epsilon_m^2}$  and  $PER_{\epsilon_n^2}$  have approximately the same importance. The probability that the mean of MAX $_{\epsilon_{m}^{2}}$  is larger than the mean of  $PER_{\epsilon_{m}^{2}}$  is 0.49. This probability value is close to 0.5, meaning that there is not enough information to compare the values.

Table 4 compares the importance factors of SMAS for the  $L_2$  error metric. Table 4 confirms the result. The rows corresponding to  $N_S$ , and  $\overline{D}_S$  have values close to 1 confirming the high importance of the two parameters. When comparing other parameters, the probability that one parameter is greater than the other is about 0.5. It confirms our previous postulation that simple rankings are not sufficient for concluding the global parameter ordering and the importance factors are significant as well.

The global ranking based on the introduced combining method in Section 5 is shown in Table 5. The table indicates that the mean degree of nodes  $\overline{D_S}$  is the most important parameter. This result is consistent with Table 2 where the mean degree of nodes  $\overline{D_S}$  is the most important parameter in the seven scenarios.

The global ranking results could be used to improve the goodness of location estimations in ad-hoc networks. To deploy a network or on an already deployed network, one could exploit the results by considering the analyzed effect of each parameter on the estimated location's accuracy. Based on the constraints of the problem, the best parameters for improving the estimated locations could be determined. For example, when there are limitations for the mean degree of the graph,

Challenges of Model Objective Selection and Estimation

| Nc                                     | 0           | 0.947            | 0.947            | 0.936          | 0.944          | 0.285          | 0.939          | 0.931          | 0.937          | 0.958          |
|--|-------------|------------------|------------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
| $B_S^{o}$                              | 0.053       | 0                | 0.493            | 0.506          | 0.504          | 0.017          | 0.501          | 0.496          | 0.500          | 0.504          |
| 2 C                                    | 0.053       | 0.507            | 0                | 0.509          | 0.505          | 0.018          | 0.504          | 0.516          | 0.507          | 0.511          |
| MAX <sub>e</sub> 2                     | 0.064       | 0.494            | 0.491            | 0              | 0.505          | 0.017          | 0.504          | 0.506          | 0.499          | 0.499          |
| $PER_{\epsilon_0^2}^{-m}$              | 0.056       | 0.496            | 0.495            | 0.495          | 0              | 0.017          | 0.492          | 0.496          | 0.502          | 0.500          |
| $\overline{D_S}$                       | 0.715       | 0.983            | 0.982            | 0.983          | 0.983          | 0              | 0.975          | 0.984          | 0.974          | 0.980          |
| MÎNL <sub>S</sub><br>MAXL <sub>S</sub> | 0.061 0.069 | 0.499<br>0.504   | 0.496<br>0.484   | 0.496<br>0.494 | 0.508<br>0.504 | 0.025<br>0.016 | 0 0.494        | 0.506<br>0     | 0.501<br>0.502 | 0.488<br>0.494 |
| $\overline{l_S}$<br>MIND <sub>S</sub>  | 0.063 0.042 | $0.500 \\ 0.496$ | $0.493 \\ 0.489$ | 0.501<br>0.501 | 0.498<br>0.500 | 0.026<br>0.020 | 0.499<br>0.512 | 0.498<br>0.506 | 0<br>0.495     | 0.505<br>0     |
| ~~ ~ ~ ~                               |             |                  |                  | > )  IT        | Table 4        | ¥.7            | ,              |                |                | , .            |

| Parameter | $N_S$ | $B_S$ | $\epsilon_S^2$ | $^{\mathrm{MAX}}\epsilon_m^2$ | $PER_{\epsilon_0^2}$ | $\overline{D_S}$ | $\operatorname{MINL}_S$ | $\mathrm{MAXL}_S$ | $\overline{\iota_S}$ | $\operatorname{MIND}_S$ |
|-----------|-------|-------|----------------|-------------------------------|----------------------|------------------|-------------------------|-------------------|----------------------|-------------------------|
| Rank      | 2     | 3     | 4              | 8                             | 10                   | 1                | 6                       | 9                 | 7                    | 5                       |
|           |       |       |                |                               |                      |                  |                         |                   |                      |                         |

| TAB    | le 5   |
|--------|--------|
| Global | ranks. |

one can increase the number of nodes in the network to increase the accuracy of the estimated location. Note that, changing one parameter typically only improves the accuracy up to a certain point; further changing the parameter would not yield an improvement in the estimation accuracy.

7. Conclusion

We introduce a systematic methodology for determining the challenge of modeling a pertinent adhoc network data set. The complex modeling problem is studied as an instance of a nonlinear optimization problem that consists of an objective function (OF) and a set of constraints. The data set is the optimization input and the estimated model is the output. We characterize the input by a set of its characteristic parameters. We define four new metrics that can be used to evaluate the goodness of an input for being optimized by a specific OF. The introduced metrics are: (1) drifting of the OF, (2) distance to the nearest local minimum, (3) the slope of the OF around the solution, and (4) the depth of the non-global local minima. We employ Plackett and Burmann simulation methodology to systematically evaluate the linear impact of various input parameters under each metric. Finally, we present a method for combining the effect of parameters under different metrics to determine the global impact of each parameter. We utilize the new methodology for estimating the locations of the nodes in an ad-hoc network where the distance measurement data is available. Three common forms of OF are considered:  $L_1$ ,  $L_2$ and  $L_{\infty}$ . Our evaluations show that the mean degree on the nodes and the number of nodes in the network are the two most important parameters for estimating the locations.

Acknowledgement

This work is partly supported by the National Science Foundation (NSF) CAREER Award under grant number 0644289. 

#### Koushanfar and Shamsi

| 1  | Re         | eferences   | 1  |
|----|------------|---|----|
| 2  |            |   | 2  |
| 3  | [1]        | BARRETT, C., MARATHE, A., MARATHE, M. V. and DROZDA, M. (2002). Characterizing the interaction  | 3  |
| 4  | [0]        | between routing and MAC protocols in ad-hoc networks. <i>MobiHoc</i> , 92-103.  | 4  |
| 5  | [2]        | ization. Information Processing in Sensor Networks (IPSN), 2673-2684.   | 5  |
| 6  | [3]        | BULLARD, C. and SEBALD, A. (1988). Monte carlo sensitivity analysis of input-output models. The   | 6  |
| 7  |            | Review of Economics and Statistics, <b>22.4</b> , 708-712.  | 7  |
| 8  | [4]<br>[5] | CENS: Center for Embedded Networked Sensing at UCLA. http://research.cens.ucla.edu/.  | 8  |
| 9  | [6]        | EREN, T., GOLDENBERG, D., WHITELEY, W., YANG, Y. R., MORSE, A., ANDERSON, B. and BELHUMEUR, P.  | 9  |
| 10 |            | (2004). Rigidity, computation, and randomization in network localization. <i>INFOCOM</i> , 2673-2684.   | 10 |
| 11 | [7]        | FENG, J., GIROD, L. and POTKONJAK, M. (2006). Consistency-based on-line localization in sensor networks. Distributed Computing in Sensor Systems (DCOSS) 520-545                                    | 11 |
| 12 | [8]        | FENG, J., GIROD, L. and POTKONJAK, M. (2006). Location discovery using data-driven statistical error  | 12 |
| 13 |            | modeling. INFOCOM, 1-14.  | 13 |
| 14 | [9]        | GOLDENBERG, D., KRISHNAMURTHY, A., MANESS, W., YANG, Y., YOUNG, A., MORSE, A. and SAVVIDES,   | 14 |
| 15 | [10]       | HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. (2001). The Elements of Statistical Learning: Data  | 15 |
| 16 | . ,        | Mining, Inference, and Prediction. Springer-Verlag, Germany.  | 16 |
| 17 | [11]       | KIRKPATRICK, S., GELATT, C. D. JR. and VECCHI, M.P. (1983). Optimization by simulated annealing.  | 17 |
| 18 | [12]       | Science, 220.4598, 671-680.<br>Koushanfar, F., Davare, A., Nguyen, D., Sangiovanni-Vincentelli, A. and Potkonjak, M. (2007).  | 18 |
| 10 | [+-]       | Techniques for maintaining connectivity in wireless ad-hoc networks under energy constraints. ACM   | 10 |
| 20 | [10]       | Trans. on Embedded Computing Systems, <b>6.3</b> , 16–37.   | 20 |
| 20 | [13]       | KOUSHANFAR, F., TAFT, N. and POTKONJAK, M. (2006). Sleeping coordination for comprehensive sensing using isotonic regression and domatic partitions. <i>INFOCOM</i> 1-13                            | 20 |
| 21 | [14]       | LEHMANN, E. (2006). Nonparametrics: Statistical Methods Based on Ranks. Springer-Verlag, Ger-   | 21 |
| 22 | r          | many.   | 22 |
| 23 | [15]       | MEGERIAN, S., KOUSHANFAR, F., QU, G., VELTRI, G. and POTKONJAK, M. (2002). Exposure in wireless sensor networks: Theory and practical solutions. <i>ACM Journal of Wireless Networks</i> 85, 443454 | 23 |
| 24 | [16]       | MEGERIAN, S., KOUSHANFAR, F., POTKONJAK, M. and SRIVASTAVA, M. (2005). Worst- and best-case   | 24 |
| 20 | r 1        | coverage in sensor networks. <i>IEEE Transactions on Mobile Computing</i> , <b>4.1</b> , 84-92.   | 20 |
| 20 | [17]       | LANZISERA, S., LIN, D. and PISTER, K. (2006). Rt time of flight ranging for wireless sensor network   | 20 |
| 21 | [18]       | MONTGOMERY, D. (2001). Design and Analysis of Experiments, 5th ed. Wiley, New York.   | 21 |
| 28 | [19]       | PATWARI, N., ASH, J., KYPEROUNTAS, S., HERO, A. III, MOSES, R. and CORREAL, R. (2005). Locating   | 28 |
| 29 |            | the nodes: cooperative localization in wireless sensor networks. <i>IEEE Signal Processing Magazine</i> , <b>22.4</b> , 54.60   | 29 |
| 30 | [20]       | PRIYANTHA, N., CHAKRABORTY, A. and BALAKRISHNAN, H. (2000). The cricket location-support sys-   | 30 |
| 31 |            | tem. <i>MOBICOM</i> , 32–43.  | 31 |
| 32 | [21]       | OSL: Open Systems Laboratory at UIUC. http://www-osl.cs.uiuc.edu/research?action=topic  | 32 |
| 33 | [22]       | PERKINS, D., HUGHES, H. and OWEN, C. (2002). Factors affecting the performance of ad hoc networks.  | 33 |
| 34 |            | <i>ICC</i> , 2048-2052.   | 34 |
| 35 | [23]       | PLACKETT, R.L. and BURMAN, J.P. (1946). The design of optimum multifactorial experiments.   | 35 |
| 36 | [24]       | SALTELLI, A., TARANTOLA, S., CAMPOLONGO, F. and RATTO, M. (2004). Sensitivity Analysis in Prac-   | 36 |
| 37 |            | tice: A Guide to Assessing Scientific Models. Wiley, England.   | 37 |
| 38 | [25]       | SAXE, J. (1979) Embeddability of weighted graphs in k-space is strongly np-hard. Allerton Conf.   | 38 |
| 39 | [26]       | SLIJEPCEVIC, S., MEGERIAN, S. and POTKONJAK, M. (2002). Location errors in wireless embedded sen-   | 39 |
| 40 |            | sor networks: sources, models, and effects on applications. Mobile Computing and Communications   | 40 |
| 41 | [07]       | Review, <b>6.3</b> , 67-78.   | 41 |
| 42 | [27]       | ad hoc networks. MSWiM, 159-168.  | 42 |
| 43 | [28]       | VADDE, K. and SYROTIUK, V. (2004). Factor interaction on service delivery in mobile ad hoc networks.  | 43 |
| 44 |            | IEEE Journal on Selected Areas in Communications (JSAC), 22.7, 13351346.  | 44 |
| 45 |            |   | 45 |
| 46 |            |   | 46 |
| 47 |            |   | 47 |
| 48 |            |   | 48 |
| 49 |            |   | 49 |
| 50 |            |   | 50 |
| 51 |            |   | 51 |
|    |            |   |    |

З

| a I   | Minimum-Variance Equity Portfolio<br>Wilhelmine von Türk <sup>1,*</sup>  |
|---|--|
|   | American Century Investments   |
| A<br>id<br>p<br>h<br>la<br>c  | <b>Abstract:</b> It is shown that the capitalization-weighted portfolio is mathemat-<br>cally required to coincide with the minimum-variance portfolio, provided both<br>ortfolios are defined with respect to the same (arbitrary) collection of equities<br>aving linearly independent returns. This result is a logical consequence of the<br>aw of iterated expectations and has important implications for equity return<br>ovariance structure.  |
| Content   | s  |
| 1 Introd<br>2 Notat<br>3 The F<br>4 The L<br>5 Proof<br>6 Concl<br>Technical<br>Reference   | luction  |
| 1. Intro  | duction  |
| Samuelso<br>Expectat<br>can be ap<br>LIE revea<br>this law o<br>This p<br>in the co<br>shown th<br>matically<br>portfolios<br>ing linear<br>return co | n (1965) was the first to recognize the relevance of the Law of Iterated<br>ions (LIE) in the area of finance. His celebrated paper shows how this law<br>oplied in the context of the price of a single risky asset. Essentially, the<br>als that forecasting error is not predictable. An illuminating discussion of<br>can be found, e.g., in Campbell, Lo and MacKinlay (1997).<br>aper extends Samuelson's insight by developing an application of the LIE<br>ntext of an arbitrary collection of equity returns. Specifically, it will be<br>at the LIE implies that the capitalization-weighted portfolio is mathe-<br>required to coincide with the minimum-variance portfolio, provided both<br>are defined with respect to the same (arbitrary) collection of equities hav-<br>ly independent returns. This result has important implications for equity<br>variance structure, as summarized in technical appendix A. |
| <sup>1</sup> 43 Del<br>*The au<br>partment c<br>with the g<br>comings in<br>AMS 2.<br>Keywor  | Mar Avenue, Berkeley, CA 94708, email: wvonturk@sbcglobal.net<br>athor gratefully acknowledges her intellectual indebtedness to the tradition of the De-<br>f Statistics, University of California at Berkeley, where she completed the dissertation<br>uidance and advice of Professor Erich L. Lehmann. The responsibility for any short-<br>the present paper rests squarely on the shoulders of the author.<br>000 subject classifications: Primary 91B28; secondary 62H25<br>ds and phrases: minimum-variance portfolio, capitalization-weighted portfolio, equity  |

З

return covariance structure, Law of Iterated Expectations, equity portfolio optimization.

The paper is organized as follows. Section 2 introduces the notation and presents a sufficient condition for proving the main result. The population view of the surprises (forecasting errors) is developed in Section 3. This population view yields an important building block for the proof of the main result. Section 4 applies the LIE in the context of a cross-section of equities, and exhibits the implication of this law from which the main result logically flows. The proof of the main result is summarized in Section 5. Section 6 concludes.

Technical appendix A summarizes the implications of this paper's main result for equity return covariance structure.

#### 2. Notation and Preliminary Results

Consider an arbitrary collection of n equities with linearly independent returns. Denote their returns by  $\mathbf{R} = (R_1 \dots R_n)'$  and denote the covariance matrix of  $\mathbf{R}$ by  $\Sigma$ . Since the returns are assumed to be linearly independent,  $\Sigma$  is nonsingular, so that its inverse  $\Sigma^{-1}$  is well defined. Let **1** denote the n x 1 vector all of whose elements are equal to 1.

Let  $\mathbf{M} = (M_1 \dots M_n)'$  denote the market capitalizations of the equities. Let  $w_i =$  $M_i/\mathbf{M'1}$  and let  $\mathbf{w} = (w_1 \dots w_n)'$ . The vector  $\mathbf{w}$  corresponds to the *capitalization* weights. Note that all of the elements of **w** are positive. Let  $\mathbf{p} = (p_1 \dots p_n)'$  denote an  $n \ge 1$  vector of constants such that  $\mathbf{p}'\mathbf{1} = 1$ . Then the elements of  $\mathbf{p}$  correspond to the investment proportions of a portfolio fully invested in the n equities, and  $\mathbf{p'R}$  denotes its return. If  $\mathbf{p} = \mathbf{w}$ , the investment proportions are those of the capitalization-weighted portfolio. The expression  $\mathbf{w'R}$  denotes the return of the capitalization-weighted portfolio.

Among all possible fully invested portfolios that can be formed from this collection of equities, there is one portfolio whose return has minimum variance. This portfolio is the minimum-variance portfolio. It is well known that, if  $\Sigma$  is nonsingular, the investment proportions of the minimum-variance portfolio are given by  $\Sigma^{-1}\mathbf{1}/\mathbf{1}'\Sigma^{-1}\mathbf{1}$  (see, e.g., Roll (1977), Campbell, Lo and MacKinlay (1997) or Grinold and Kahn (2000)).

The desired result is that

$$\mathbf{w} = \Sigma^{-1} \mathbf{1} / \mathbf{1}' \Sigma^{-1} \mathbf{1}.$$

For the purpose of finding a proof of this result, it is convenient to define

 $\beta_i = \operatorname{Cov}(R_i, \mathbf{w'R}) / \operatorname{Var}(\mathbf{w'R}).$ 

Let  $\boldsymbol{\beta} = (\beta_1 \dots \beta_n)'$ . Then

Note that  $\mathbf{w}'\mathbf{1} = 1$ , so that

$$\mathbf{w}' \Sigma \mathbf{w} = 1/eta' \Sigma^{-1} \mathbf{1}$$

and

$$\mathbf{w} = \Sigma^{-1} \boldsymbol{\beta} / \boldsymbol{\beta}' \Sigma^{-1} \mathbf{1}.$$

47  
48With this notation it is possible to establish the following preliminary result.47  
4849Lemma 2.1. The following assertions are equivalent:49  
50  
(i) 
$$\mathbf{w} = \Sigma^{-1} \mathbf{1} / \mathbf{1}' \Sigma^{-1} \mathbf{1}$$
50  
50  
5151(ii)  $\beta = \mathbf{1}$ 51

 $\boldsymbol{\beta} = \boldsymbol{\Sigma} \mathbf{w} / \mathbf{w}' \boldsymbol{\Sigma} \mathbf{w}.$ 

*Proof.*  $\mathbf{w} = \Sigma^{-1} \boldsymbol{\beta} / \boldsymbol{\beta}' \Sigma^{-1} \mathbf{1}$  implies that  $\Sigma \mathbf{w} = \boldsymbol{\beta} / \boldsymbol{\beta}' \Sigma^{-1} \mathbf{1}$ , and  $\mathbf{w} = \Sigma^{-1} \mathbf{1} / \mathbf{1}' \Sigma^{-1} \mathbf{1}$ implies that  $\Sigma \mathbf{w} = \mathbf{1}/\mathbf{1}'\Sigma^{-1}\mathbf{1}$ . Therefore,  $\mathbf{w} = \Sigma^{-1}\mathbf{1}/\mathbf{1}'\Sigma^{-1}\mathbf{1}$  if and only if  $\boldsymbol{\beta} = c\mathbf{1}$ for  $c = \beta' \Sigma^{-1} \mathbf{1} / \mathbf{1}' \Sigma^{-1} \mathbf{1}$ . Since  $\mathbf{w}' \boldsymbol{\beta} = \mathbf{w}' \mathbf{1} = 1$ , this constant c is equal to 1. З Next, consider an investment period which begins at time t = 0 and ends at time t = T. Note that at time t = 0, w is a vector of known constants while R is a vector of random variables. Let  $E[\mathbf{R}]$  denote the  $n \ge 1$  vector whose elements correspond to  $E[R_i]$ , where  $E[R_i]$  is the expected value of  $R_i$  at time t = 0 (i = 1, ..., n). The expression  $E[R_i]$  corresponds to the best forecast at time t = 0 of the realized value of  $R_i$  observed at time t = T. Now shift attention from the returns  $R_i$  to the surprises  $S_i$ , where  $S_i = R_i - R_i$  $E[R_i]$ . The surprises correspond to forecasting errors. Let  $\mathbf{S} = (S_1, \ldots, S_n)'$ , let  $\mathcal{A}$ denote the  $\sigma$ -field induced by **S**, and let  $\mathcal{F}$  denote the  $\sigma$ -field induced by **w'S**. Since  $\mathbf{w'S}$  is a function of  $\mathbf{S}, \mathcal{F} \subset \mathcal{A}$ . The conditional expectation of  $S_i$  given  $\mathbf{w'S}$  can then be written as  $E^{\mathcal{F}}[S_i]$ , which is defined almost surely (a.s.) in the sense that any two versions agree, except possibly on a null set in  $\mathcal{A}$ . With this notation, it is possible to simplify the task of proving the desired result as follows. **Lemma 2.2.** For the purpose of showing that  $\mathbf{w} = \Sigma^{-1} \mathbf{1} / \mathbf{1}' \Sigma^{-1} \mathbf{1}$ , it is sufficient to show that  $\mathbb{E}^{\mathcal{F}}[S_i] = \mathbf{w}'\mathbf{S}$  (a.s.)  $(i = 1, \dots, n)$ . *Proof.* In light of Lemma 2.1, it is sufficient to show that  $E^{\mathcal{F}}[S_i] = \mathbf{w}'\mathbf{S}$  (a.s.)  $(i = 1, \ldots, n)$  implies that  $\beta = 1$ . By definition,  $\operatorname{Cov}(R_i, \mathbf{w'R}) = \operatorname{E}[S_i \mathbf{w'S}]$ and  $\operatorname{Var}(\mathbf{w}'\mathbf{R}) = \operatorname{E}\left[(\mathbf{w}'\mathbf{S})^2\right].$ Therefore it suffices to show that  $E^{\mathcal{F}}[S_i] = \mathbf{w}'\mathbf{S}$  (a.s.) (i = 1, ..., n) implies that  $\mathbf{E}[S_i \mathbf{w}' \mathbf{S}] = \mathbf{E}\left[ (\mathbf{w}' \mathbf{S})^2 \right] \quad (i = 1, \dots, n).$ This is easily accomplished by recalling the usual properties of conditional expec-tation operators, which yields  $\mathbf{E}[S_i \mathbf{w}' \mathbf{S}] = \mathbf{E}\left[\mathbf{w}' \mathbf{S} \mathbf{E}^{\mathcal{F}}[S_i]\right],$ so that  $E^{\mathcal{F}}[S_i] = \mathbf{w}'\mathbf{S}$  (a.s.) implies  $E[S_i\mathbf{w}'\mathbf{S}] = E\left[(\mathbf{w}'\mathbf{S})^2\right]$  (i = 1, ..., n). Before it can be shown that the assertion  $E^{\mathcal{F}}[S_i] = \mathbf{w}'\mathbf{S}$  (a.s.) (i = 1, ..., n)flows logically from the LIE, it is first necessary to develop the population view of the surprises (forecasting errors). This is done in the next section. 3. The Population View Consider an arbitrary collection of n publicly traded companies and an investment period which begins at time t = 0 and ends at time t = T. The development of the population view of the surprises depends on the operation of repricing the 

shares of each company at the beginning of the investment period, while adjusting
the shares outstanding so as to leave the market capitalization of each company
unchanged. This operation, very familiar to equity investors, guarantees that there
is no loss of generality in repricing the shares of each company at the beginning of
the investment period to have a value of one dollar.

W. von Türk

As in Section 2, let  $M_i$  denote the market capitalizations at the beginning of the investment period (i = 1, ..., n). If the price of one share is one dollar (and if the market capitalizations are rounded to the nearest dollar), then  $M_i$  corresponds to the number  $N_i$  of shares outstanding for the *i*-th company. Each share of the *i*-th company has return  $R_i$  and associated surprise (forecasting error)  $S_i$ . Therefore, at the end of the investment period, the observed values  $s_i$  of the random variables  $S_i$ form a population in which  $s_i$  occurs with frequency  $N_i$ . The capitalization weights  $w_i$  then correspond to the relative frequencies with which the  $s_i$  are observed. This population has mean w's, where  $\mathbf{s} = (s_1, \ldots, s_n)'$  is the observed value of the random vector  $\mathbf{S} = (S_1, \ldots, S_n)'$ . 

It is well known that the population mean is the expected value of a randomly selected element from that population. From the standpoint of time t = T, the population mean is observed with certainty, while from the standpoint of time t = 0, the population mean corresponds to the random variable  $\mathbf{w}'\mathbf{S}$ . This means that

#### $E_T[a randomly selected surprise] = \mathbf{w}'\mathbf{s},$

where " $\mathbf{E}_T$ " is to be read as "the expected value at time t = T, conditional on the observed value  $\mathbf{w's}$  of  $\mathbf{w'S}$ ." Translating this into  $\sigma$ -field notation yields the following lemma.

**Lemma 3.1.**  $E^{\mathcal{F}}[a \text{ randomly selected surprise}] = \mathbf{w'S} (a.s.), where \mathcal{F} denotes the <math>\sigma$ -field induced by  $\mathbf{w'S}$ .

As will be seen in Section 5, Lemma 3.1 is an important building block for the proof of the main result.

# 4. The Law of Iterated Expectations

A key insight provided by the LIE is that forecasting error is not predictable. The present paper contemplates an arbitrary collection of n equities with returns  $R_i$  and associated forecasting errors  $S_i = R_i - \mathbb{E}[R_i]$  (i = 1, ..., n). In this context, the LIE implies that the expected value of any surprise  $S_i$  corresponds to the expected value of a randomly selected surprise. From the standpoint of time t = 0, this yields the set of equations

$$E[S_i] = E[a \text{ randomly selected surprise}] (i = 1, ..., n).$$

Conditional on  $\mathbf{w}'\mathbf{S}$ , the LIE similarly implies the following lemma.

**Lemma 4.1.**  $E^{\mathcal{F}}[S_i] = E^{\mathcal{F}}[a \text{ randomly selected surprise}] (a.s.) (i = 1, ..., n), where$  $\mathcal{F}$  denotes the  $\sigma$ -field induced by  $\mathbf{w}'\mathbf{S}$ .

Lemma 4.1 exhibits the implication of the LIE from which the main result of the present paper logically flows, as summarized in the next section.

#### 5. Proof of the Main Result

The main result is stated as the proposition below. PROPOSITION: For an arbitrary collection of n equities having linearly independent returns,  $\mathbf{w} = \Sigma^{-1} \mathbf{1} / \mathbf{1}' \Sigma^{-1} \mathbf{1}$ .

imsart-coll ver. 2008/08/29 file: Von\_Turk.tex date: April 10, 2009

*Proof.* From Lemma 4.1,

 $\mathbf{E}^{\mathcal{F}}[S_i] = \mathbf{E}^{\mathcal{F}}[a \text{ randomly selected surprise}] (a.s.) (i = 1, \dots, n).$ 

From Lemma 3.1,

$$E^{\mathcal{F}}[a \text{ randomly selected surprise}] = \mathbf{w'S} (a.s.).$$

Combining Lemmas 4.1 and 3.1 yields

$$\mathbf{E}^{\mathcal{F}}[S_i] = \mathbf{w}'\mathbf{S} \text{ (a.s.) } (i = 1, \dots, n).$$

In light of Lemma 2.2, this completes the proof.

#### 6. Conclusion

It has been shown that the capitalization-weighted portfolio is mathematically required to coincide with the minimum-variance portfolio, provided both portfolios are defined with respect to the same (arbitrary) collection of equities having linearly independent returns. This result is a logical consequence of the LIE, and has important implications for equity return covariance structure, as summarized in technical appendix A.

#### Technical Appendix A

The main result of this paper has important implications for equity return covariance structure, as summarized in the following proposition.

**PROPOSITION:** For any collection of n equities with linearly independent returns, the covariance matrix  $\Sigma$  of  $\mathbf{R} = (R_1, \ldots, R_n)'$  is of the form

$$\Sigma = \mathbf{11}'k + \mathbf{U}^2$$

where  $\mathbf{U}^2$  is a diagonal matrix such that the *i*-th diagonal element is positive and inversely proportional to  $M_i$  (i = 1, ..., n), and where k corresponds to a constant which can be positive, negative or zero.

Note that for positive values of k, general equity return covariance structure exhibited in the proposition above corresponds to the covariance matrix of a 1factor model in which

- (i) the factor loadings of the unique common factor are all equal to one; and
- (ii) the variances of the specific factors are inversely proportional to the market capitalizations of the equities.

#### References

- [1] CAMPBELL, J. Y., LO, A. W. and MACKINLAY, A. C. (1997). The Econometrics of Financial Markets. Princeton University Press, Princeton, New Jersey. GRINOLD, R. C. and KAHN, R. N. (2000). Active Portfolio Management. McGraw-Hill, New York.
- ROLL, R. (1977). A critique of the asset pricing theory's tests. Journal of Financial Economics 4, 129-176.
- SAMUELSON, P. (1965). Proof that properly anticipated prices fluctuate randomly. Industrial Man-[4] agement Review 6, 41-49.

imsart-coll ver. 2008/08/29 file: Von\_Turk.tex date: April 10, 2009

З