## The Many Flavors of Penalized Linear Discriminant Analysis

Daniela M. Witten Assistant Professor of Biostatistics University of Washington

May 9, 2011 Fourth Erich L. Lehmann Symposium Rice University

## Overview

There has been a great deal of interest in the past 15+ years in penalized regression,

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \{ ||\mathbf{y} - \mathbf{X}\boldsymbol{\beta}||^2 + P(\boldsymbol{\beta}) \},$$

especially in the setting where the number of features p exceeds the number of observations n.

- ► *P* is a penalty function. Could be chosen to promote
  - ▶ sparsity: e.g. the lasso,  $P(\beta) = ||\beta||_1$
  - smoothness
  - piecewise constancy...
- ► How can we extend the concepts developed for regression when p > n to other problems?
- ► A Case Study: Penalized linear discriminant analysis.

The Normal Model Fisher's Discriminant Problem Optimal Scoring

## The classification problem

- ► The Set-up:
  - ► We are given *n* training observations x<sub>1</sub>,..., x<sub>n</sub> ∈ ℝ<sup>p</sup>, each of which falls into one of K classes.
  - Let y ∈ {1,...,K}<sup>n</sup> contain class memberships for the training observations.
  - Let  $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1' \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}$ .
  - ► Each column of X (feature) is centered to have mean zero.

The Normal Model Fisher's Discriminant Problem Optimal Scoring

# The classification problem

- ► The Set-up:
  - ► We are given *n* training observations x<sub>1</sub>,..., x<sub>n</sub> ∈ ℝ<sup>p</sup>, each of which falls into one of K classes.
  - ► Let  $\mathbf{y} \in \{1, ..., K\}^n$  contain class memberships for the training observations.

• Let 
$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n^T \end{pmatrix}$$

- ► Each column of X (feature) is centered to have mean zero.
- ► The Goal:
  - We wish to develop a classifier based on the training observations x<sub>1</sub>,..., x<sub>n</sub> ∈ ℝ<sup>p</sup>, that we can use to classify a test observation x<sup>\*</sup> ∈ ℝ<sup>p</sup>.
  - A classical approach: linear discriminant analysis.

The Normal Model Fisher's Discriminant Problem Optimal Scoring

#### Linear discriminant analysis



The Normal Model Fisher's Discriminant Problem Optimal Scoring

#### LDA via the normal model

Fit a simple normal model to the data:

$$\mathbf{x}_i | y_i = k \sim N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_w)$$

Apply Bayes' Theorem to obtain a classifier: assign x\* to the class for which δ<sub>k</sub>(x\*) is largest:

$$\delta_k(\mathbf{x}^*) = \mathbf{x}^{*T} \mathbf{\Sigma}_w^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^T \mathbf{\Sigma}_w^{-1} \boldsymbol{\mu}_k + \log \pi_k$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

The Normal Model Fisher's Discriminant Problem Optimal Scoring

## Fisher's discriminant

A geometric perspective: project the data to achieve good classification.



√ へ (~
6 / 29

The Normal Model Fisher's Discriminant Problem Optimal Scoring

### Fisher's discriminant

A geometric perspective: project the data to achieve good classification.



≣ •ી લ (∿ 6 / 29

The Normal Model Fisher's Discriminant Problem Optimal Scoring

### Fisher's discriminant

A geometric perspective: project the data to achieve good classification.



√ へ (~
6 / 29

The Normal Model Fisher's Discriminant Problem Optimal Scoring

### Fisher's discriminant

A geometric perspective: project the data to achieve good classification.



6 / 29

The Normal Model Fisher's Discriminant Problem Optimal Scoring

#### Fisher's discriminant and the associated criterion

Look for the discriminant vector  $oldsymbol{eta} \in \mathbb{R}^p$  that maximizes

$$\boldsymbol{\beta}^{T} \hat{\boldsymbol{\Sigma}}_{b} \boldsymbol{\beta}$$
 subject to  $\boldsymbol{\beta}^{T} \hat{\boldsymbol{\Sigma}}_{w} \boldsymbol{\beta} \leq 1$ .

- $\hat{\Sigma}_b$  is an estimate for the between-class covariance matrix.
- $\hat{\Sigma}_{w}$  is an estimate for the within-class covariance matrix.
- This is a generalized eigen problem; can obtain multiple discriminant vectors.
- To classify, multiply data by discriminant vectors and perform nearest centroid classification in this reduced space.
- ► If we use K 1 discriminant vectors then we get the LDA classification rule. If we use fewer than K 1, we get reduced-rank LDA.

SOR

The Normal Model Fisher's Discriminant Problem Optimal Scoring

## LDA via optimal scoring

- Classification is such a bother. Isn't regression so much nicer?
- It wouldn't make sense to solve

$$\underset{\boldsymbol{\beta}}{\operatorname{minimize}} \{ || \mathbf{y} - \mathbf{X} \boldsymbol{\beta} ||^2 \}.$$

But can we formulate classification as a regression problem in some other way?

The Normal Model Fisher's Discriminant Problem Optimal Scoring

## LDA via optimal scoring

• Let **Y** be a  $n \times K$  matrix of dummy variables;  $Y_{ik} = 1_{y_i = k}$ .

minimize {
$$||\mathbf{Y}\boldsymbol{\theta} - \mathbf{X}\boldsymbol{\beta}||^2$$
} subject to  $\boldsymbol{\theta}^T \mathbf{Y}^T \mathbf{Y} \boldsymbol{\theta} = 1$ .

- We are choosing the optimal scoring of the class labels in order to recast the classification problem as a regression problem.
- The resulting β is proportional to the discriminant vector in Fisher's discriminant problem.
- ► Can obtain the LDA classification rule, or reduced-rank LDA.

The Normal Model Fisher's Discriminant Problem Optimal Scoring

#### Linear discriminant analysis



The Normal Model Optimal Scoring Fisher's Discriminant Problem

### LDA when $p \gg n$

When  $p \gg n$ , we cannot apply LDA directly, because the within-class covariance matrix is singular.

There is also an interpretability issue:

- ► All *p* features are involved in the classification rule.
- ► We want an interpretable classifier. For instance, a classification rule that is a
  - sparse,
  - smooth, or
  - piecewise constant

linear combination of the features.

The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Penalized LDA

- We could extend LDA to the high-dimensional setting by applying (convex) penalties, in order to obtain an interpretable classifier.
- ► For concreteness, in this talk: we will use l<sub>1</sub> penalties in order to obtain a sparse classifier.
- ▶ Which version of LDA should we penalize, and does it matter?



토▶ ∢토▶ 토 ∽੧. 12/29 Linear Discriminant Analysis **The Normal Model Penalized LDA** Optimal Scoring Connections Fisher's Discriminant Problem

#### Penalized LDA via the normal model

The classification rule for LDA is

$$\mathbf{x}^{*T} \hat{\mathbf{\Sigma}}_{w}^{-1} \hat{\boldsymbol{\mu}}_{k} - \frac{1}{2} \hat{\boldsymbol{\mu}}_{k}^{T} \hat{\mathbf{\Sigma}}_{w}^{-1} \hat{\boldsymbol{\mu}}_{k},$$

where  $\hat{\boldsymbol{\Sigma}}_w$  and  $\hat{\boldsymbol{\mu}}_k$  denote MLEs for  $\boldsymbol{\Sigma}_w$  and  $\boldsymbol{\mu}_k$ .

- When  $p \gg n$ , we cannot invert  $\hat{\Sigma}_w$ .
- Can use a regularized estimate of  $\Sigma_w$ , such as

$$\boldsymbol{\Sigma}_{w}^{D} = \begin{pmatrix} \hat{\sigma}_{1}^{2} & 0 & \dots & 0 \\ 0 & \hat{\sigma}_{2}^{2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \hat{\sigma}_{p}^{2} \end{pmatrix}$$

・ロト ・ ロ ト ・ ヨ ト ・ 日 ト ・ の へ ()

٠

13/29

Linear Discriminant Analysis The Normal Model Penalized LDA Optimal Scoring Connections Fisher's Discriminant Problem

Interpretable class centroids in the normal model

- For a sparse classifier, we need zeros in estimate of  $\Sigma_w^{-1}\mu_k$ .
- An interpretable classifier:
  - Use  $\Sigma_w^D$ , and estimate  $\mu_k$  according to

$$\underset{\boldsymbol{\mu}_k}{\text{minimize}} \left\{ \sum_{j=1}^p \sum_{i:y_i=k} \frac{(X_{ij} - \mu_{kj})^2}{\sigma_j^2} + \lambda ||\boldsymbol{\mu}_k||_1 \right\}.$$

- ► Apply Bayes' Theorem to obtain a classification rule.
- This is the nearest shrunken centroids proposal, which yields a sparse classifier because we are using a diagonal estimate of the within-class covariance matrix and a sparse estimate of the class mean vectors.

Citation: Tibshirani et al. 2003, Stat Sinica

Linear Discriminant Analysis Penalized LDA Connections The Normal Model Optimal Scoring Fisher's Discriminant Problem

### Penalized LDA via optimal scoring

► We can easily extend the optimal scoring criterion:

$$\underset{\boldsymbol{\beta},\boldsymbol{\theta}}{\text{minimize}} \{ \frac{1}{n} || \mathbf{Y}\boldsymbol{\theta} - \mathbf{X}\boldsymbol{\beta} ||^2 + \lambda ||\boldsymbol{\beta}||_1 \} \text{ subject to } \boldsymbol{\theta}^T \mathbf{Y}^T \mathbf{Y} \boldsymbol{\theta} = 1.$$

- ► An efficient iterative algorithm will find a local optimum.
- ► We get sparse discriminant vectors, and hence classification using a subset of the features.

Citation: Clemmensen Hastie Witten and Ersboll 2011, Submitted

Linear Discriminant Analysis The Normal Model Penalized LDA Optimal Scoring Connections Fisher's Discriminant Problem

Penalized LDA via Fisher's discriminant problem

• A simple formulation:

$$\underset{\boldsymbol{\beta}}{\operatorname{maximize}} \{ \boldsymbol{\beta}^{\mathsf{T}} \hat{\boldsymbol{\Sigma}}_{b} \boldsymbol{\beta} - \lambda || \boldsymbol{\beta} ||_{1} \} \text{ subject to } \boldsymbol{\beta}^{\mathsf{T}} \tilde{\boldsymbol{\Sigma}}_{w} \boldsymbol{\beta} \leq 1$$

where  $\tilde{\boldsymbol{\Sigma}}_{w}$  is some full rank estimate of  $\boldsymbol{\Sigma}_{w}$ .

- A non-convex problem, because  $\beta^T \hat{\Sigma}_b \beta$  isn't concave in  $\beta$ .
- Can we find a local optimum?

Citation: Witten and Tibshirani 2011, JRSSB

The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Maximizing a function via minorization



うへで 17/29

The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Maximizing a function via minorization



The Normal Model Optimal Scoring Fisher's Discriminant Problem

### Maximizing a function via minorization



The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Maximizing a function via minorization



The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Maximizing a function via minorization



The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Maximizing a function via minorization



The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Maximizing a function via minorization



The Normal Model Optimal Scoring Fisher's Discriminant Problem

### Maximizing a function via minorization



The Normal Model Optimal Scoring Fisher's Discriminant Problem

## Minorization

 Key point: Choose a minorizing function that is easy to maximize.



 Minorization allows us to efficiently find a local optimum for Fisher's discriminant problem with any convex penalty.

Normal +  $\ell_1$  and Fisher's +  $\ell_1$ Fisher's +  $\ell_1$  and Optimal Scoring +  $\ell_1$ 

### Connections between flavors of penalized LDA



つへで 19/29

Normal +  $\ell_1$  and Fisher's +  $\ell_1$ Fisher's +  $\ell_1$  and Optimal Scoring +  $\ell_1$ 

#### Connections between flavors of penalized LDA

- 1. Normal Model  $+ \ell_1$ : use a diagonal estimate for  $\Sigma_w$  and then apply  $\ell_1$  penalty to the class mean vectors.
- 2. Optimal scoring  $+ \ell_1$ : apply  $\ell_1$  penalty to discriminant vectors.
- 3. Fisher's discriminant problem  $+ \ell_1$ : apply  $\ell_1$  penalty to discriminant vectors.

So are (1) and (3) different? And are (2) and (3) the same?

Normal +  $\ell_1$  and Fisher's +  $\ell_1$ Fisher's +  $\ell_1$  and Optimal Scoring +  $\ell_1$ 

```
Normal Model + \ell_1 and Fisher's + \ell_1
```





つへで 21/29

Normal +  $\ell_1$  and Fisher's +  $\ell_1$ Fisher's +  $\ell_1$  and Optimal Scoring +  $\ell_1$ 

#### Normal Model $+ \ell_1$ and Fisher's $+ \ell_1$

$\mu_{11}$	$\mu_{21}$	$\mu_{K1}$	
$\mu_{12}$	$\mu_{22}$	$\mu_{K2}$	
$\boldsymbol{\mu}_{13}$	$\mu_{23}$	$\mu_{K3}$	
•	•	•	
•	•	•	
:	:	:	
•	•		
•	•	•	
$\mu_{1p}$	$\mu_{2p}$	$\mu_{Kp}$	

- ► Normal model + l<sub>1</sub> penalizes the elements of this matrix.
- ► Fisher's + l<sub>1</sub> penalizes the left singular vectors.
- ► Clearly these are different...
- ► ...but if K = 2, then they are (essentially) the same.

Normal +  $\ell_1$  and Fisher's +  $\ell_1$ Fisher's +  $\ell_1$  and Optimal Scoring +  $\ell_1$ 

```
Normal Model + \ell_1 and Fisher's + \ell_1
```





୬ < ୯ 23 / 29

Normal  $+ \ell_1$  and Fisher's  $+ \ell_1$ Fisher's  $+ \ell_1$  and Optimal Scoring  $+ \ell_1$ 

```
Fisher's+\ell_1 and Optimal Scoring+\ell_1
```



Both problems involve "penalizing the discriminant vectors" so they must be the same, right?

Normal  $+ \ell_1$  and Fisher's  $+ \ell_1$ Fisher's  $+ \ell_1$  and Optimal Scoring  $+ \ell_1$ 

## Fisher's+ $\ell_1$ and Optimal Scoring+ $\ell_1$

**Theorem:** For any value of the tuning parameter for  $FD+\ell_1$ , there exists some tuning parameter for  $OS+\ell_1$  such that the solution to one problem is a critical point of the other.

- In other words there is a correspondence between the critical points, though not necessarily the solutions.
- So the resulting "sparse discriminant vectors" may be different!

Normal +  $\ell_1$  and Fisher's +  $\ell_1$ Fisher's +  $\ell_1$  and Optimal Scoring +  $\ell_1$ 

# Connections



## Pros and Cons

#### Penalized LDA via normal model:

- ► (+) In the case of a diagonal estimate for Σ<sub>w</sub> and ℓ<sub>1</sub> penalties on mean vectors, it is well-motivated and simple.
- (-) No obvious extension to non-diagonal estimates of  $\Sigma_w$ .
- ► (-) Cannot obtain a "low-rank" classifier.

#### Penalized LDA via Fisher's discriminant problem:

- ► (+) Any convex penalties can be applied to discriminant vectors.
- (+) Can use any full-rank estimate of  $\Sigma_w$ .
- ► (+) Can obtain a "low-rank" classifier.

#### Penalized LDA via optimal scoring:

- (+) An extension of regression.
- $\blacktriangleright$  (+) Any convex penalties can be applied to discriminant vectors.
- ► (+) Can obtain a "low-rank" classifier.
- (-) Cannot use any estimate of  $\Sigma_w$ .
- ► (-) Usual motivation for OS is that it yields the same discriminant vectors as Fisher's problem. Not true when penalized!

# Conclusions

• A sensible way to regularize regression when  $p \gg n$ :

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \{ || \mathbf{y} - \mathbf{X}\boldsymbol{\beta} ||^2 + P(\boldsymbol{\beta}) \}.$$

- One could argue that this is *the* way to penalize regression.
- But as soon as we step away from regression, even to a closely related problem like LDA, the situation becomes much more complex – there is no longer a "single way" to approach the problem.
- And the situation becomes only more complex for more complex statistical methods!
- Need a principled framework to determine which penalized extension of established statistical methods is "best".

## References

- Witten and Tibshirani (2011) Penalized classification using Fisher's linear discriminant. To appear in *Journal of the Royal Statistical Society, Series B.*
- Clemmensen, Hastie, Witten, and Ersboll (2011) Sparse discriminant analysis. Submitted.