

INCORPORATING BIOLOGICAL INFORMATION INTO LINEAR MODELS: A BAYESIAN APPROACH TO THE SELECTION OF PATHWAYS AND GENES

BY FRANCESCO C. STINGO, YIAN A. CHEN, MAHLET G. TADESSE
AND MARINA VANNUCCI¹

*Rice University, Moffitt Cancer Center, Georgetown University
and Rice University*

The vast amount of biological knowledge accumulated over the years has allowed researchers to identify various biochemical interactions and define different families of pathways. There is an increased interest in identifying pathways and pathway elements involved in particular biological processes. Drug discovery efforts, for example, are focused on identifying biomarkers as well as pathways related to a disease. We propose a Bayesian model that addresses this question by incorporating information on pathways and gene networks in the analysis of DNA microarray data. Such information is used to define pathway summaries, specify prior distributions, and structure the MCMC moves to fit the model. We illustrate the method with an application to gene expression data with censored survival outcomes. In addition to identifying markers that would have been missed otherwise and improving prediction accuracy, the integration of existing biological knowledge into the analysis provides a better understanding of underlying molecular processes.

1. Introduction. DNA microarrays have been used successfully to identify gene expression signatures characteristic of disease subtypes [Golub et al. (1999)] or distinct outcomes to therapy [Shipp et al. (2002)]. Many statistical methods have been developed to select genes for disease diagnosis, prognosis and therapeutic targets. However, gene selection alone may not be sufficient. In cancer pharmacogenomics, for instance, cancer drugs are increasingly designed to target specific pathways to account for the complexity of the oncogenic process and the complex relationships between genes [Downward (2006)]. Metabolic pathways, for example, are defined as a series of chemical reactions in a living cell that can be activated or inhibited at multiple points. If a gene at the top of a signaling cascade is selected as a target, it is not guaranteed that the reaction will be successfully inactivated, because multiple genes downstream can still be activated or inhibited. Signals are generally relayed via multiple signaling routes or networks. Even if a branch of the pathway is completely blocked by inhibition or activation of multiple

Received October 2009; revised February 2011.

¹Supported in part by NIH Grant R01-HG0033190-05 and NSF Grant DMS-10-07871.

Key words and phrases. Bayesian variable selection, gene expression, Markov chain Monte Carlo, Markov random field prior, pathway selection.

genes, the signal may still be relayed through an alternative branch or even through a different pathway [Bild et al. (2006)]. Downward (2006) pointed out that targeting a single pathway or a few signaling pathways might not be sufficient. Thus, the focus is increasingly on identifying both relevant genes and pathways. Genes and/or gene products generally interact with one another and they often function together concertedly. Here we propose a Bayesian model that addresses this question by incorporating information of pathway memberships and gene networks in the analysis of DNA microarray data. Such information is used to define pathway summaries, specify prior distributions, and structure the MCMC moves.

Several public and commercial databases have been developed to structure and store the vast amount of biological knowledge accumulated over the years into functionally or biochemically related groups. These databases focus on describing signaling, metabolic or regulatory pathways. Some examples include Gene Ontology (GO) [Ashburner et al. (2000)], Kyoto Encyclopedia of Genes and Genomes (KEGG) [Kanehisa and Goto (2000)], MetaCyc [Krieger et al. (2004)], PathDB, Reactome KnowledgeBase [Joshi-Tope et al. (2005)], Invitrogen iPath (www.invitrogen.com) and Cell Signaling Technology (CST) Pathway (www.cellsignal.com). The need to integrate gene expression data with the biological knowledge accumulated in these databases is well recognized. Several software packages that query pathway information and overlay DNA microarray data on pathways have been developed. Nakao et al. (1999) implemented a visualization tool that color-codes KEGG pathway diagrams to reflect changes in their gene expression levels. GenMAPP [Dahlquist et al. (2002)] is another graphical tool that allows visualization of microarray data in the context of biological pathways or any other functional grouping of genes. Doniger et al. (2003) use GenMAPP to view genes involved in specific GO terms. Another widely used method that relates pathways to a set of differentially expressed genes is the gene set enrichment analysis (GSEA) [Subramanian et al. (2005)]. Given a list of genes, GSEA computes an enrichment score to reflect the degree to which a predefined pathway is over-represented at the top or bottom of the ranked list. These procedures are useful starting points to observe gene expression changes for known biological processes.

Recent studies have gone a step further and focused on incorporating pathway information or gene–gene network information into the analysis of gene expression data. For example, Park, Hastie and Tibshirani (2007) have attempted to incorporate GO annotation to predict survival time, first grouping genes based on their GO membership, calculating the first principal component to form a super-gene within each cluster and then applying a Cox model with L_1 penalty to identify super-genes, that is, GO terms related to the outcome. Wei and Li (2007) have considered a small set of 33 preselected signaling pathways and used the implied relationships among genes to infer differentially expressed genes, and Wei and Li (2008) have extended this work by including a temporal dimension. Li and Li (2008) and Pan, Xie and Shen (2010) have proposed different procedures that use the gene–gene

network to build penalties in a regression model for gene selection. Bayesian approaches have also been developed. [Li and Zhang \(2010\)](#) have incorporated the dependence structure of transcription factors in a regression model with gene expression outcomes. There, a network is defined based on the Hamming distance between candidate motifs and used to specify a Markov random field prior for the motif selection indicator. [Telesca et al. \(2008\)](#) have proposed a model for the identification of differentially expressed genes that takes into account the dependence structure among genes from available pathways while allowing for correction in the gene network topology. [Stingo and Vannucci \(2011\)](#) use a Markov random field prior that captures the gene–gene interaction network in a discriminant analysis setting.

These methods use the gene–pathway relationships or gene network information to identify either the important pathways or the genes. Our goal is to develop a more comprehensive method that selects both pathways and genes using a model that incorporates pathway–gene relationships and gene dependence structures. In order to identify relevant genes and pathways, latent binary vectors are introduced and updated using a two-stage Metropolis–Hastings sampling scheme. The gene networks are used to define a Markov random field prior on the gene selection indicators and to structure the Markov chain Monte Carlo (MCMC) moves. In addition, the pathway information is used to derive pathway expression measures that summarize the group behavior of genes within pathways. In this paper we make use of the first latent components obtained by applying partial least squares (PLS) regressions on the selected genes from each pathway. PLS is an efficient statistical regression technique that was initially proposed in the chemometrics literature [[Wold \(1966\)](#)] and more recently used for the analysis of genomic and proteomic data; see [Boulesteix and Strimmer \(2007\)](#). We apply the model to simulated and real data using the pathway structure from the KEGG database.

Our simulation studies show that the MRF prior leads to a better separation between relevant and nonrelevant pathways, and to less false positives in a model with fairly small regression coefficients. Other authors have reported similar results. [Li and Zhang \(2010\)](#), in particular, comment on the effect of the MRF prior on the selection power in their linear regression setting. They also notice that adding the MRF prior implies a relatively small increase in computational cost. [Wei and Li \(2007, 2008\)](#) report that their method is quite effective in identifying genes and modified subnetworks and that it has higher sensitivity than commonly used procedures that do not use the pathway structure, with similar and, in some cases, lower false discovery rates. Furthermore, in our model formulation we use the network information not only for prior specification but also to structure the MCMC moves. This is helpful for arriving at promising models faster by proposing relevant configurations. In real data applications the integration of pathway information may allow the identification of relevant predictors that could be missed otherwise, aiding the interpretation of the results, in particular, for the selected

genes that are connected in the MRF, and also improving the prediction accuracy of selected models.

The paper is organized as follows. Section 2 contains the model formulation and prior specification. Section 3 describes the MCMC procedure and strategies for posterior inference. In Section 4 performances are evaluated on simulated data and an application of the method to gene expression data with survival outcomes is presented. Section 5 concludes the paper with a brief discussion.

2. Model specification. We describe how we incorporate pathway and gene network information into a Bayesian modeling framework for gene and pathway selection. Figure 1 represents a schematic representation of our approach and model.

2.1. Regression on latent measures of pathway activity. Our goal is to build a model for identifying pathways related to a particular phenotype while simultaneously locating genes from these selected pathways that are involved in the biological process of interest. The data we have available for analysis consist of the following:

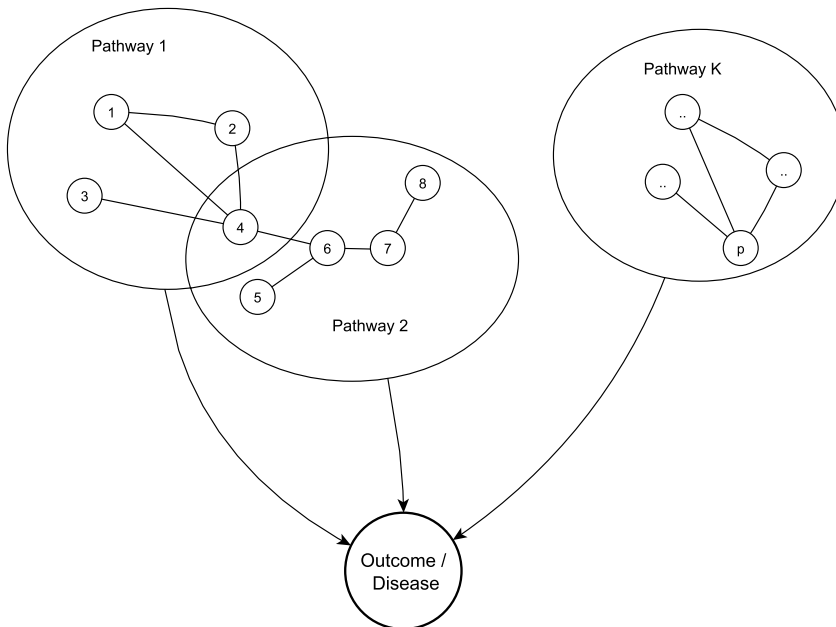


FIG. 1. Schematic representation of our proposed approach. Information on known pathways and gene–gene networks is obtained from available databases. PLS components restricted to known pathways serve as possible regressors to predict a disease outcome, according to model (1). The goal of the inference is to identify the pathways to be included in the model and the genes to be included within those pathways.

- (1) Y , an $n \times 1$ vector of outcomes.
- (2) \mathbf{X} , an $n \times p$ matrix of gene expression levels. Without loss of generality, \mathbf{X} is centered so that its columns sum to 0.
- (3) \mathbf{S} , a $K \times p$ matrix indicating membership of genes in pathways, with elements $s_{kj} = 1$ if gene j belongs to pathway k , and $s_{kj} = 0$ otherwise.
- (4) \mathbf{R} , a $p \times p$ matrix describing relationships between genes, with $r_{ij} = 1$ if genes i and j have a direct link in the gene network, and $r_{ij} = 0$ otherwise.

The matrices \mathbf{S} and \mathbf{R} are constructed using information retrieved from pathway databases; see the application in Section 4.2 for details.

Since the goal of the analysis is to study the association between the response variable and the pathways, we need to derive a score as a measure of “pathway expression” that summarizes the group behavior of included genes within pathways. We do this by using the latent components from a PLS regression of Y on selected subsets of genes from each pathway. A number of recent studies have, in fact, applied dimension reduction techniques to capture the group behavior of multiple genes. Pittman et al. (2004), for instance, first apply k -means clustering to identify subsets of potentially related genes, then use as regressors the first principal components obtained from applying principal component analysis (PCA) to each cluster. Bair et al. (2006) start by removing genes that have low univariate correlation with the outcome variable, then apply PCA on the remaining genes to form clusters or conceptual pathways, which are used as regressors. In our method, instead of attempting to infer conceptual pathways, we use the existing pathway information. We compute a pathway activity measure by applying PLS regression of Y on a subset of selected genes from the pathway. PLS has the advantage of taking into account the covariance between regressors and the response variable Y , whereas PCA focuses solely on the variability in the covariate data. The selection of a subset of gene expressions to form the PLS components is similar in spirit to the sparse PCA method proposed by Zou, Hastie and Tibshirani (2006), which selects variables to form the principal components.

To identify both relevant groups and important genes, we introduce two binary vector indicators, a $K \times 1$ vector $\boldsymbol{\theta}$ for the inclusion of the groups and a $p \times 1$ vector $\boldsymbol{\gamma}$ for the inclusion of genes, that is, $\gamma_j = 1$ if gene j is selected for at least one pathway score, and $\gamma_j = 0$ otherwise. Assuming that the response Y is continuous, the linear regression model that relates Y to the selected pathways and genes is

$$(1) \quad Y = \mathbf{1}\alpha + \sum_{k=1}^{K_\theta} T_{k(\boldsymbol{\gamma})} \beta_{k(\boldsymbol{\gamma})} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}),$$

where $K_\theta = \sum_{k=1}^K \theta_k$ is the number of selected pathways and where $T_{k(\boldsymbol{\gamma})}$ corresponds to the first latent PLS component generated based on the expression levels of selected genes belonging to pathway k , that is, using the X_j 's corresponding

to $s_{kj} = 1$ and $\gamma_j = 1$. To be more precise, let pathway k contain $p_k = \sum_{j=1}^p s_{kj}$ genes and let $p_{k\gamma} = \sum_{j=1}^p s_{kj}\gamma_j$ denote the number of selected genes (i.e., genes included in the model) that belong to pathway k . Then $T_{k(\gamma)}$ corresponds to the first latent PLS component generated by applying PLS to the expression data of the $p_{k\gamma}$ genes, denoted as $\mathbf{X}_{k(\gamma)}$,

$$T_{k(\gamma)} = \mathbf{X}_{k(\gamma)}U_1,$$

where U_1 is the $p_{k\gamma} \times 1$ eigenvector corresponding to the largest eigenvalue of $C_{xy}C_{xy}^T$, with $C_{xy} = \text{cov}(\mathbf{X}_{k(\gamma)}, Y)$ [see, e.g., Lindgren, Geladi and Wold (1993)]. Thus, $T_{k(\gamma)}$ is an $n \times 1$ vector and $\beta_k(\gamma)$ is a scalar. Model (1) can therefore be seen as a PLS regression model with PLS components restricted to available pathways, and where the goal of the inference is to identify the pathways to be included in the model, and the genes to be included within those pathways.

2.2. *Models for categorical or censored outcomes.* In the construction above, we have assumed a continuous response. However, our model formulation can easily be extended to handle categorical or censored outcome variables.

When Y is a categorical variable taking one of G possible values, $0, \dots, G - 1$, a probit model can be used, as done by Albert and Chib (1993), Sha et al. (2004) and Kwon et al. (2007). Briefly, each outcome Y_i is associated with a vector $(p_{i,0}, \dots, p_{i,G-1})$, where $p_{ig} = P(Y_i = g)$ is the probability that subject i falls in the g th category. The probabilities p_{ig} can be related to the linear predictors using a data augmentation approach. Let \mathbf{Z}_i be latent data corresponding to the unobserved propensities of subject i to belong to one of the classes. When the observed outcomes Y_i correspond to nominal values, the relationship between Y_i and $\mathbf{Z}_i = (z_{i,1}, \dots, z_{i,G-1})$ can be defined as

$$(2) \quad Y_i = \begin{cases} 0, & \text{if } \max_{1 \leq l \leq G-1} \{z_{i,l}\} \leq 0, \\ g, & \text{if } \max_{1 \leq l \leq G-1} \{z_{i,l}\} > 0 \text{ and } z_{i,g} = \max_{1 \leq l \leq G-1} \{z_{i,l}\}. \end{cases}$$

A multivariate normal model can then be used to associate \mathbf{Z}_i to the predictors

$$(3) \quad \mathbf{Z}_i = \mathbf{1}\alpha + \sum_{k=1}^{K_\theta} T_{i,k(\gamma)}\beta_{k(\gamma)} + \boldsymbol{\varepsilon}_i, \quad \boldsymbol{\varepsilon}_i \sim \mathcal{N}(0, \boldsymbol{\Sigma}), i = 1, \dots, n.$$

If the observed outcomes Y_i correspond, instead, to ordinal categories, the latent variable Z_i is defined such that $Y_i = g$ if $\delta_g < Z_i \leq \delta_{g+1}$, $g = 0, \dots, G - 1$, where the boundaries δ_g are unknown and $-\infty = \delta_0 < \delta_1 < \dots < \delta_{G-1} < \delta_G = \infty$. The latent variable Z_i is associated with the predictors through the linear model

$$(4) \quad Z_i = \alpha + \sum_{k=1}^{K_\theta} T_{i,k(\gamma)}\beta_{k(\gamma)} + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2), i = 1, \dots, n.$$

For censored survival outcomes, an accelerated failure time (AFT) model can be used [Wei (1992); Sha, Tadesse and Vannucci (2006)]. In this case, the observed data are $Y_i = \min(T_i, C_i)$ and $\delta_i = I\{Y_i \leq C_i\}$, where T_i is the survival time for subject i , C_i is the censoring time, and δ_i is a censoring indicator. A data augmentation approach can be used and latent variables Z_i can be introduced such that

$$(5) \quad \begin{cases} Z_i = \log(Y_i), & \text{if } \delta_i = 1, \\ Z_i > \log(Y_i), & \text{if } \delta_i = 0. \end{cases}$$

The AFT model can then be written in terms of the latent Z_i similarly to (4) where the ε_i 's are independent and identically distributed random variables that may take one of several parametric forms. Sha, Tadesse and Vannucci (2006) consider cases where ε_i follows a normal or a t -distribution.

2.3. Prior for regression parameters. The regression coefficient β_k in (1) measures the effect of the PLS latent component summarizing the effect of pathway k on the response variable. However, not all pathways are related to the phenotype and the goal is to identify the predictive ones. Bayesian methods that use mixture priors for variable selection have been thoroughly investigated in the literature, in particular, for linear models; see George and McCulloch (1997) for multiple regression, Brown, Vannucci and Fearn (1998) for extensions to multivariate responses and Sha et al. (2004) for probit models. A comprehensive review on features of the selection priors and on computational aspects of the method can be found in Chipman, George and McCulloch (2001). Similarly, we use the latent vector $\boldsymbol{\theta}$ to specify a scale mixture of a normal density and a point mass at zero for the prior on each β_k in (1):

$$(6) \quad \beta_k | \theta_k, \sigma^2 \sim \theta_k \cdot \mathcal{N}(\beta_0, h\sigma^2) + (1 - \theta_k) \cdot \delta_0(\beta_k), \quad k = 1, \dots, K,$$

where $\delta_0(\beta_k)$ is a Dirac delta function. The hyperparameter h in (6) regulates, together with the hyperparameters of $p(\boldsymbol{\theta}, \boldsymbol{\gamma} | \eta)$ defined in Section 2.4 below, the amount of shrinkage in the model. We follow the guidelines provided by Sha et al. (2004) and specify h in the range of variability of the data so as to control the ratio of prior to posterior precision. For the intercept term, α , and the variance, σ^2 , we take conjugate priors $\alpha | \sigma^2 \sim \mathcal{N}(\alpha_0, h_0\sigma^2)$ and $\sigma^2 \sim \text{Inv-Gamma}(\nu_0/2, \nu_0\sigma_0^2/2)$, where α_0 , β_0 , h_0 , h , ν_0 and σ_0^2 are to be elicited.

2.4. Priors for pathway and gene selection indicators. In this section we define the prior distributions for the pathway selection indicator, $\boldsymbol{\theta}$, and gene selection indicator, $\boldsymbol{\gamma}$. These priors are first defined marginally then jointly, taking into account some necessary constraints.

Each element of the latent K -vector $\boldsymbol{\theta}$ is defined as

$$(7) \quad \theta_k = \begin{cases} 1, & \text{if pathway } k \text{ is represented in the model,} \\ 0, & \text{otherwise} \end{cases}$$

for $k = 1, \dots, K$. We assume independent Bernoulli priors for the θ_k 's,

$$(8) \quad p(\boldsymbol{\theta}|\varphi_k) = \prod_{k=1}^K \varphi_k^{\theta_k} (1 - \varphi_k)^{1-\theta_k},$$

where φ_k determines the proportion of pathways expected a priori in the model. A mixture prior can be further specified for φ_k to achieve a better discrimination in terms of posterior probabilities between significant and nonsignificant pathways by inflating $p(\theta_k = 0)$ toward 1 for the nonrelevant pathways, as first suggested by Lucas, Carvalho, Wang, Bild, Nevins and West (2006),

$$(9) \quad p(\varphi_k) = \rho \delta_0(\varphi_k) + (1 - \rho) \mathcal{B}(\varphi_k|a_0, b_0),$$

where $\mathcal{B}(\varphi_k|a_0, b_0)$ is a Beta density function with parameters a_0 and b_0 . Since inference on φ_k is not of interest, it can be integrated out to simplify the MCMC implementation. This leads to the following marginal prior for $\boldsymbol{\theta}$:

$$(10) \quad p(\boldsymbol{\theta}) = \prod_k \left[\rho \cdot (1 - \theta_k) + (1 - \rho) \cdot \frac{B(a_0 + \theta_k, b_0 + 1 - \theta_k)}{B(a_0, b_0)} \right],$$

where $B(\cdot, \cdot)$ is the Beta function. Prior (10) corresponds to a product of Bernoulli distributions with parameter $\varphi_k^* = \frac{a_0(1-\rho)}{a_0+b_0}$.

For the latent p -vector \boldsymbol{y} we specify a prior distribution that is able to take into account not only the pathway membership of each gene but also the biological relationships between genes within and across pathways, which are captured by the matrix R . Following Li and Zhang (2010), we model these relations using a Markov random field (MRF), where genes are represented by nodes and relations between genes by edges. A MRF is a graphical model in which the distribution of a set of random variables follow Markov properties that can be described by an undirected graph. In particular, two unconnected genes are considered conditionally independent given all other genes [Besag (1974)]. Relations on the MRF are represented by the following probabilities:

$$(11) \quad p(\gamma_j|\eta, \gamma_i, i \in N_j) = \frac{\exp(\gamma_j F(\gamma_j))}{1 + \exp(F(\gamma_j))},$$

where $F(\gamma_j) = (\mu + \eta \sum_{i \in N_j} \gamma_i)$ and N_j is the set of direct neighbors of gene j in the MRF using only pathways represented in the model, that is, pathways with $\theta_k = 1$. The corresponding global distribution on the MRF is given by

$$(12) \quad p(\boldsymbol{y}|\boldsymbol{\theta}, \mu, \eta) \propto \exp(\mu \mathbf{1}'_p \boldsymbol{y} + \eta \boldsymbol{y}' \mathbf{R} \boldsymbol{y}),$$

with $\mathbf{1}_p$ the unit vector of dimension p and \mathbf{R} the matrix introduced in Section 2.1. The parameter μ controls the sparsity of the model, while η regulates the smoothness of the distribution of \boldsymbol{y} over the graph by controlling the prior probability of selecting a gene based on how many of its neighbors are selected. In particular, higher values of η encourage the selection of genes with neighbors already selected

into the model. If a gene does not have any neighbor, then its prior distribution reduces to an independent Bernoulli with parameter $p = \exp(\mu)/[1 + \exp(\mu)]$, which is a logistic transformation of μ .

Here, unlike Li and Zhang (2010), who fix both parameters of the MRF prior, we specify a hyperprior for η . We give positive probability to values of η bigger than 0, which is biologically more intuitive than negative values of this parameter (which would favor neighboring genes to have different inclusion status). Such restriction on the domain of η also minimizes the “phase transition” problem that typically occurs with MRF parameterizations of type (11), where the dimension of the selected model increases massively for small increments of η . When the phase transition occurs the number of selected genes increases substantially. Here, after having detected the phase transition value η_{PT} , by simulating from (12) over a grid of η values, we specify a Beta distribution $\text{Beta}(c_0, d_0)$ on η/η_{PT} .

Constraints need to be imposed to ensure both interpretability and identifiability of the model. We essentially want to avoid the following:

- (1) empty pathways, that is, selecting a pathway but none of its member genes;
- (2) orphan genes, that is, selecting a gene but none of the pathways that contain it;
- (3) selection of identical subsets of genes by different pathways, a situation that generates identical values $T_{k(\gamma)}$ and $T_{k'(\gamma)}$ to be included in the model.

These constraints imply that some combinations of θ and γ values are not allowed. The joint prior probability for (θ, γ) taking into account these constraints is given by

$$p(\theta, \gamma | \eta) \propto \begin{cases} \prod_{k=1}^K \varphi_k^{*\theta_k} (1 - \varphi_k^*)^{1-\theta_k} \exp(\mu \mathbf{1}'_p \gamma + \eta \gamma' \mathbf{R} \gamma), & \text{for valid configurations,} \\ 0, & \text{for invalid configurations.} \end{cases}$$

3. Model fitting. We now describe our MCMC procedure to fit the model and discuss strategies for posterior inference with huge posterior spaces, as in this model. In the Bayesian literature on variable selection for standard linear regression models stochastic search algorithms have been designed to explore the posterior space, and have been successfully employed in genomic applications with prohibitive settings, handling models with thousands of genes. A key to these applications is the assumption of sparsity of the model, that is, the belief that the response is associated with a small number of regressors. A stochastic search then allows one to explore the posterior space in an effective way, quickly finding the most probable configurations, that is, those corresponding to coefficients with high marginal probabilities, while spending less time in regions with low posterior probability.

We describe below the MCMC algorithm we have designed for our problem. In particular, borrowing from the literature on stochastic searches for variable selection, we work with a marginalized model and design a Metropolis–Hastings algorithm that updates the indicator parameters for the inclusion of pathways and genes with a set of moves that add and/or delete a single gene and a single pathway. Also, we update the parameter η of the MRF from its posterior distribution by employing the general method proposed by Møller et al. (2006). In Stingo et al. (2011) we discuss how our Bayesian stochastic search variable selection kernel generates an ergodic Markov chain over the restricted space. In applications, we have found that a good way to assess if the stochastic exploration can be considered satisfactory is to check the concordance of the posterior probabilities obtained from different chains started from different initial points.

3.1. *Marginal posterior probabilities.* The model parameters consist of $(\alpha, \beta, \sigma^2, \boldsymbol{\gamma}, \boldsymbol{\theta}, \eta)$. The MCMC procedure can be made more efficient by integrating out some of the parameters. Here, we integrate out the regression parameters, α , β and σ^2 . This leads to a multivariate t -distribution

$$(13) \quad f(Y|\mathbf{T}, \boldsymbol{\theta}, \boldsymbol{\gamma}) \sim \mathcal{T}_{\nu_0}(\alpha_0 \mathbf{1}_n + \mathbf{T}_{(\boldsymbol{\theta}, \boldsymbol{\gamma})} \beta_0, \sigma_0^2 (\mathbf{I}_n + h_0 \mathbf{1}_n \mathbf{1}_n' + \mathbf{T}_{(\boldsymbol{\theta}, \boldsymbol{\gamma})} \boldsymbol{\Sigma}_0 \mathbf{T}_{(\boldsymbol{\theta}, \boldsymbol{\gamma})}')),$$

with ν_0 degrees of freedom and $\mathbf{1}_n$ an n -vector of ones, and where $\boldsymbol{\Sigma}_0 = h \mathbf{I}_{K_\theta}$, with \mathbf{I}_n the $n \times n$ identity matrix, and $\mathbf{T}_{(\boldsymbol{\theta}, \boldsymbol{\gamma})}$ the $n \times K_\theta$ matrix derived from the first PLS latent components for the selected pathways using the selected genes. In the notation (13) the two arguments of the t -distribution represent the mean and the scale parameter of the distribution, respectively. The posterior probability distribution of the pathway and gene selection indicators is then given by

$$(14) \quad f(\boldsymbol{\theta}, \boldsymbol{\gamma}, \eta | \mathbf{T}, Y) \propto f(Y | \mathbf{T}, \boldsymbol{\theta}, \boldsymbol{\gamma}) \cdot p(\boldsymbol{\theta}, \boldsymbol{\gamma} | \eta) \cdot p(\eta).$$

3.2. *MCMC sampling.* The MCMC steps consist of the following: (I) sampling pathway and gene selection indicators from $p(\boldsymbol{\theta}, \boldsymbol{\gamma} | \text{rest})$; (II) sampling the MRF parameter from $p(\eta | \text{rest})$; (III) sampling additional parameters introduced when fitting probit models for categorical outcomes or AFT models for survival data.

- (I) The parameters $(\boldsymbol{\theta}, \boldsymbol{\gamma})$ are updated using a Metropolis–Hastings algorithm in a two-stage sampling scheme. The pathway–gene relationships are used to structure the moves and account for the constraints specified in Section 2.4. Details of the MCMC moves to update $(\boldsymbol{\theta}, \boldsymbol{\gamma})$ are given in Stingo et al. (2011). They consist of randomly choosing one of the following move types:
- (1) change the inclusion status of gene and pathway by randomly choosing between adding a pathway and a gene or removing them both;
 - (2) change the inclusion status of gene but not pathway by randomly choosing between adding a gene or removing a gene;

- (3) change the inclusion status of pathway but not gene by randomly choosing between adding a pathway or removing a pathway.
- (II) At this step we want to draw the MRF parameter η from the posterior density $p(\eta|\boldsymbol{y}) \propto p(\eta)p(\boldsymbol{y}|\eta)$. The prior distribution on \boldsymbol{y} is of the form

$$(15) \quad p(\boldsymbol{y}|\eta) = q_\eta(\boldsymbol{y})/Z_\eta$$

with unnormalized density $q_\eta(\boldsymbol{y})$ and a normalizing constant Z_η which is not available analytically. When calculating the Metropolis–Hastings ratio to determine the acceptance probability of a new value η^p ,

$$(16) \quad H(\eta^p|\eta^o) = \frac{p(\eta^p)q_{\eta^p}(\boldsymbol{y})q(\eta^o|\eta^p)}{p(\eta^o)q_{\eta^o}(\boldsymbol{y})q(\eta^p|\eta^o)} \bigg/ \frac{Z_{\eta^p}}{Z_{\eta^o}},$$

with η^o the current value for η , one needs to take into account that $Z_{\eta^p}/Z_{\eta^o} \neq 1$. Following Møller et al. (2006), we introduce an auxiliary variable w , defined on the same state space as that of \boldsymbol{y} , which has conditional density $f(w|\eta, \boldsymbol{y})$, and consider the posterior $p(\eta, w|\boldsymbol{y}) \propto f(w|\eta, \boldsymbol{y})p(\eta)q_\eta(\boldsymbol{y})/Z_\eta$, which of course still involves the unknown Z_η . Obviously, marginalization over w of $p(\eta, w|\boldsymbol{y})$ gives the desired distribution $p(\eta|\boldsymbol{y})$. Now, if (η^o, w^o) is the current state of the algorithm, we first propose η^p with density $q(\eta^p|\eta^o)$, then w^p with density $q(w^p|w^o, \eta^p, \eta^o)$. As usual, the choice of these proposal densities is arbitrary from the point of view of the equilibrium distribution of the chain of η values. The choice of $f(w|\eta, \boldsymbol{y})$ is also arbitrary. The key idea of the method proposed by Møller et al. (2006) is to take the proposal density for the auxiliary variable w to be of the same form as (15), but dependent on η^p rather than η^o , that is,

$$(17) \quad q(w^p|w^o, \eta^p, \eta^o) = p(w^p|\eta^p) = q_{\eta^p}(w^p)/Z_{\eta^p}.$$

Then the Metropolis–Hastings ratio becomes

$$(18) \quad H(\eta^p, w^p|\eta^o, w^o) = \frac{f(w^p|\eta^p, \boldsymbol{y})p(\eta^p)q_{\eta^p}(\boldsymbol{y})q_{\eta^o}(w^o)q(\eta^o|\eta^p)}{f(w^o|\eta^o, \boldsymbol{y})p(\eta^o)q_{\eta^o}(\boldsymbol{y})q_{\eta^p}(w^p)q(\eta^p|\eta^o)},$$

and no longer depends on Z_{η^p}/Z_{η^o} . The new value w^p for the auxiliary variable w is drawn from (17) by perfect simulation using the algorithm proposed by Propp and Wilson (1996).

- (III) In the case of classification or survival outcomes, the augmented data Z need to be updated from their full conditionals using Gibbs sampling; see Sha et al. (2004), Sha, Tadesse and Vannucci (2006) and Kwon et al. (2007) for details.

3.3. *Posterior inference.* The MCMC procedure results in a list of visited models with included pathways indexed by θ and selected genes indexed by γ , and their corresponding relative posterior probabilities. Pathway selection can be based on the marginal posterior probabilities $p(\theta_k | \mathbf{T}, Y)$. A simple strategy is to compute Monte-Carlo estimates by counting the number of appearances of each pathway across the visited models. Relevant pathways are identified as those with largest marginal posterior probabilities. Then relevant genes from these pathways are identified based on their marginal posterior probabilities conditional on the inclusion of a pathway of interest, $p(\gamma_j | \mathbf{T}, Y, I\{\theta_k s_{kj} = 1\})$. An alternative inference for gene selection is to focus on a subset of pathways, \mathcal{P} , and consider the marginal posterior probability conditional on at least one pathway the gene belongs to being represented in the model, $p(\gamma_j | \mathbf{T}, Y, I\{\sum_{k \in \mathcal{P}} \theta_k s_{kj} > 0\})$. We note that Rao–Blackwellized estimates have been employed in standard linear regression models, in place of frequency estimates, by averaging the full conditional posterior probabilities of the inclusion indicators. These estimates are computationally quite expensive, though they may have better precision, as noted by Guan and Stephens (2011). Because of our strategy for inference, that selects first pathways and then genes conditional on selected pathways, Rao–Blackwellized estimates of marginal probabilities may not be straightforward to derive. In all simulations and examples reported in this paper we have obtained satisfactory results by simply estimating the marginal posterior probabilities with the corresponding relative frequencies of inclusion in the visited models.

Inference for a new set of observations, (\mathbf{X}_f, Y_f) , can be done via least squares prediction, $\hat{Y}_f = \mathbf{1}_n \tilde{\alpha} + \mathbf{T}_{f(\theta, \gamma)} \tilde{\beta}_{(\theta, \gamma)}$, where $\mathbf{T}_{f(\theta, \gamma)}$ is the first principal component based on selected genes from relevant pathways and where $\tilde{\alpha} = \bar{Y}$ and $\tilde{\beta}_{(\theta, \gamma)} = (\mathbf{T}'_{(\theta, \gamma)} \mathbf{T}_{(\theta, \gamma)} + h^{-1} \mathbf{I}_{K_\theta})^{-1} \mathbf{T}'_{(\theta, \gamma)} Y$, with Y the response variable in the training and $\mathbf{T}_{(\theta, \gamma)}$ the scores obtained from the training data using selected pathways and genes included in the model. Note that for prediction purposes, since we do not know the future Y_f , a PLS regression cannot be fit. Therefore, we generate $T_{f(\theta, \gamma)}$ by considering the first latent component obtained by applying PCA to each selected pathway using the included genes.

In the case of categorical or censored survival outcomes, the sampled latent variables Z would be used to estimate \hat{Z}_f , then the correspondence between Z and the observed outcome outlined in Section 2.2 can be invoked to predict Y_f [Sha et al. (2004, 2006); Kwon et al. (2007)].

4. Application. We assess performances on simulated data, then illustrate an application to microarrays using the KEGG pathway database to define the MRF.

4.1. *Simulation studies.* We investigated the performance of our model using simulated data based on the gene-pathway relations, \mathbf{S} , and gene network, \mathbf{R} , of 70 pathways and 1,098 genes from the KEGG database. The relevant pathways were

defined by selecting 4 pathways at random. For each of the 4 selected pathways, one gene was picked at random and its direct neighbors that belong to the selected pathways were chosen. This resulted in the selection of 4 pathways and 15 genes: 7 out of 30 from the first pathway, 3 out of 35 from the second, 3 out of 105 from the third, and 2 out of 47 from the fourth pathway. Gene expressions for $n = 100$ samples were simulated for these 15 genes using an approach similar to Li and Li (2008). This was accomplished by first creating an ordering among the 15 selected genes by changing the undirected edges in the gene networks into directed edges. The first node on the ordering, which we denote by X_{F_1} , was selected from each pathway and drawn from a standard normal distribution; note that this node has no parents. Then all child nodes directly connected only to X_{F_1} and denoted by X_{F_2} were drawn from $X_{F_2} \sim \mathcal{N}(X_{F_1}\rho, 1)$. Subsequent child nodes at generation j , X_{F_j} , were drawn using all parents from $X_{F_j} \sim \mathcal{N}(\rho \mathbf{X}_{pa(F_j)} \mathbf{1}_{|pa(F_j)|}, 1)$, where $pa(F_j)$ indicates the set of parents of node j and $\mathbf{X}_{pa(F_j)}$ is a matrix containing the expressions of all the $|pa(F_j)|$ parents for node j . The expression levels of the remaining 1,073 genes deemed irrelevant were simulated from a standard normal density. The response variables for the $n = 100$ samples were generated from

$$Y_i = \sum_{j=1}^{15} X_{ij}\beta + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, 1), i = 1, \dots, 100.$$

For the first data set we set $\beta = \pm 0.5$, with the same sign for genes belonging to the same pathways. For the second and third data sets we used $\beta = \pm 1$ and $\beta = \pm 1.5$, respectively. Note how the generating process is different from model (1) being fit.

We report results obtained by choosing, when possible, hyperparameters that lead to weakly informative prior distributions. A vague prior is assigned to the intercept α by setting h_0 to a large value tending to ∞ . For σ^2 , the shape parameter can be set to $\nu_0/2 = 3$, the smallest integer such that the variance of the inverse-gamma distribution exists, and the scale parameter $\nu_0\sigma_0^2/2$ can be chosen to yield a weakly informative prior. For the vector of regression coefficients, β_k , we set the prior mean to $\beta_0 = 0$ and choose h in the range of variability of the covariates, as suggested in Section 2.3. Specifically, we set $h_0 = 10^6$, $\alpha_0 = \beta_0 = 0$, $\nu_0\sigma_0/2 = 0.5$, and $h = 0.02$. For the pathway selection indicators, θ_k , we set $\varphi_k^* = 0.01$. As for the prior at the gene level, we set $\mu = -3.5$, corresponding to setting the proportion of genes expected a priori in the model to, at least, 3% of the total number of genes. Parameters φ_k^* and μ influence the sparsity of the model and consequently the magnitude of the marginal posterior probabilities. Some sensitivity is, of course, to be expected. However, in our simulations we have noticed that the ordering of pathways and genes based on posterior probability remains roughly the same and, therefore, the final selections are unchanged as long as one adjusts the threshold on the posterior probabilities. Also, for the hyperprior on η , we set $\eta_{PT} = 0.092$, to avoid the phase transition problem, and $c_0 = 5$ and $d_0 = 2$, to obtain a prior distribution that favors bigger values of η in the interval $0 \leq \eta \leq \eta_{PT}$. In our simulations we did not notice sensitivity to the specification of c_0 and d_0 .

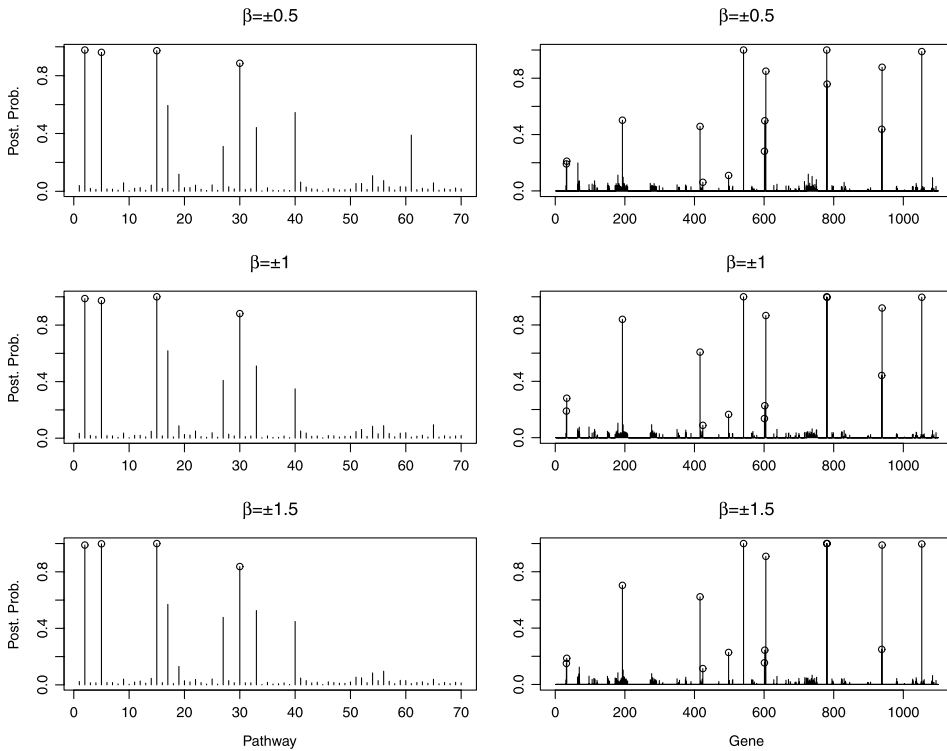


FIG. 2. Simulated data: Marginal posterior probabilities for pathway selection, $p(\theta_k|\mathbf{T}, Y)$, and conditional posterior probabilities for gene selection, $p(\gamma_j|\mathbf{T}, Y, I\{\sum_{k \in \mathcal{P}} \theta_k s_{kj} > 0\})$, for the three simulated data sets. Open circles indicate pathways and genes used to generate the outcome variable.

The MCMC sampler was run for 300,000 iterations with the first 50,000 used as burn-in. We computed the marginal posterior probabilities for pathway selection, $p(\theta_k = 1|Y, \mathbf{T})$, and the conditional posterior probabilities for gene selection given a subset of selected pathways, $p(\gamma_j|\mathbf{T}, Y, I\{\sum_{k \in \mathcal{P}} \theta_k s_{kj} > 0\})$. Figure 2 displays the marginal posterior probabilities of inclusion for all 70 pathways and the conditional posterior probabilities of inclusion for all 1,098 genes.

Important pathways and genes can be selected as those with highest posterior probabilities. For example, in all 3 scenarios all four relevant pathways were selected with a marginal posterior probability cutoff of 0.8. Reducing the selection threshold to a marginal posterior probability of 0.5 pulls in two false positive pathways, for all the three simulated scenarios considered. One of the false positives is the pathway with index 17 in Figure 2, which contains more than 100 genes. A closer investigation of the MCMC output reveals that different subsets of its member genes are selected whenever it is included in the model, resulting in a high marginal posterior of inclusion for the pathway but low marginal posterior probabilities for all its member genes. The second false positive pathway appears

to be selected often because it contains two or three of the relevant genes that were used to simulate the response variable and were also included in the model with high marginal posterior probabilities; all its other member genes have very low probabilities of selection. As expected, the identification of the relevant genes is easier when the signal-to-noise ratio is higher. Conditional upon the best 4 selected pathways, a marginal posterior probability cutoff of 0.5 on the marginal probability of gene inclusion leads to the selection of 7, 8 and 8 relevant genes, for the three scenarios, respectively, and no false positives. With a marginal probability threshold of 0.1, 14 of the relevant genes are selected with 4 false positives for the scenario with $\beta = \pm 0.5$, while 13 relevant genes are selected with only two false positives for the simulated data with $\beta = \pm 1$. In the simulated setting with $\beta = \pm 1.5$ all the 15 relevant genes are selected without any false positive at a threshold of 0.12.

Generally speaking, the effect of the MRF prior depends on the concordance of the prior network with the data. For the simulated data, we found that the model with the MRF prior, compared to the same model without the MRF, performs better in terms of pathway selection, as it provides a clearer separation between relevant and nonrelevant pathways. In particular, the average difference, over the three scenarios, between the relevant pathway with the lowest posterior probability and the nonrelevant pathway with the highest posterior probability is 0.28, while without the MRF prior it is only 0.18. In addition, we have observed increased sensitivity of the MRF prior in selecting the true variables. For example, for the simulated case with $\beta \pm 1.5$, in order to select all 15 relevant genes, the marginal probability cutoff must be reduced to 0.088 at the expense of including 3 false positives. Other authors have reported similar results [Li and Zhang (2010)]. In the real data application we describe below, employing information on gene–gene networks aids the interpretation of the results, in particular, for those selected genes that are connected in the MRF, and improves the prediction accuracy.

4.2. *Application to microarray data.* We consider the van't Veer et al. (2002) breast cancer microarray data.² Gene expression measures were collected on each patient using DNA microarray with 24,481 probes. Missing expressions were imputed using a k -nearest neighbor algorithm with $k = 10$. The procedure consists of identifying the k closest genes to the one with missing expression in array j using the other $n - 1$ arrays, then imputing the missing value by the average expression of the k neighbors [Troyanskaya et al. (2001)]. We focus on the 76 sporadic lymph-node-negative patients, 33 of whom developed distant metastasis within 5 years; the remaining 43 are viewed as censored cases. We randomly split the patients into a training set of 38 samples and a test set of the same size using a fairly balanced split of metastatic/nonmetastatic cases. The goal is to identify a subset of pathways and genes that can predict time to distant metastasis.

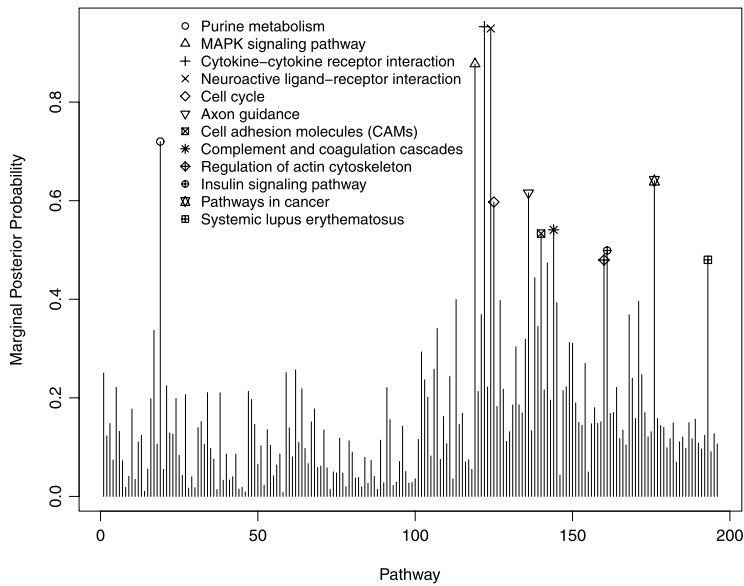
²Available at www.rii.com/publications/2002/vantveer.htm.

The gene network and pathway information were obtained from the KEGG database. This was accomplished by mapping probes to pathways using the links between pathway node identifiers and LocusLink ID. Using the R package *KEGG-graph* [Zhang and Wiemann (2009)], we first downloaded the gene network for each pathway, then merged all networks into a single one with all genes. A total of 196 pathways and 3,592 probes were included in the analysis, with each pathway containing multiple genes and with most genes associated with several pathways.

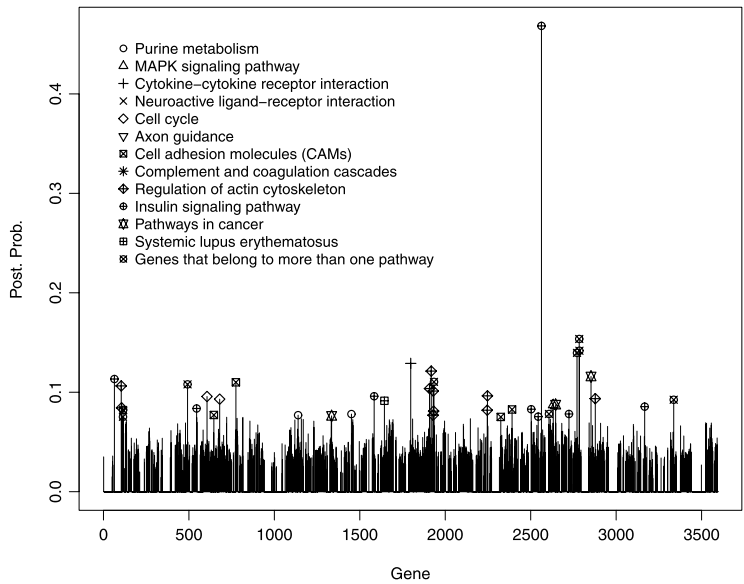
We ran two MCMC chains with different starting numbers of included variables, 50 and 80, respectively. We used 600,000 iterations with a burn-in of 100,000 iterations. We incorporated the first latent vector of the PLS for each pathway into the analysis as described in Section 2.1 and set the number of pathways expected a priori in the model to 10% of the total number. For the gene selection, we set the hyperparameter of the Markov random field to $\mu = -3.5$, indicating that a priori at least 3% of genes are expected to be selected. We set $\eta_{PT} = 0.09$, to avoid the phase transition problem, and $c_0 = 1$ and $d_0 = 1$, to obtain a noninformative prior distribution. A sensitivity analysis showed that the posterior inference is not affected by different values of c_0 and d_0 . We set $\alpha_0 = \beta_0 = 0$, $h_0 = 10^6$ and $h = 0.1$ for the prior on the regression parameters and obtained a vague prior for σ^2 by choosing $\nu_0/2 = 3$ and $\nu_0\sigma_0^2/2 = 0.5$.

The trace plots for the number of included pathways and the number of selected genes showed good mixing (figures not shown). The MCMC samplers mostly visited models with 20–45 pathways and 50–90 genes. To assess the agreement of the results between the two chains, we looked at the correlation between the marginal posterior probabilities for pathway selection, $p(\theta_k|\mathbf{T}, Y)$, and found good concordance between the two MCMC chains with a correlation coefficient of 0.9933. Concordance among the marginal posterior probabilities was confirmed by looking at a scatter plot of $p(\theta_k|\mathbf{T}, Y)$ across the two MCMC chains (figure not shown).

The model also showed good predictive performance. Sha, Tadesse and Vannucci (2006) already analyzed these data using an AFT model with 3,839 probes as predictors and obtained a predictive MSE of 1.9317 using the 11 probe sets with highest marginal probabilities. Our model incorporating pathway information achieved a predictive MSE of 1.4497 on the validation set, using 12 selected pathways and 41 probe sets with highest posterior probabilities. The selected pathways and genes are clearly indicated in the marginal posterior probability plots displayed in Figure 3. If we increase the marginal probability thresholds for selection and consider a model with 7 selected pathways and 14 genes, to make the comparison more fair with the results of Sha, Tadesse and Vannucci (2006), we obtain a MSE of 1.7614. As a reminder, our model selects relevant pathways and relevant genes simultaneously, while the model of Sha, Tadesse and Vannucci (2006) selects genes only. Of course, one can always select pathways post-hoc, as those that contain the selected genes. However, as single genes belong to multiple pathways, we expect our approach to give a more precise selection.



(a)



(b)

FIG. 3. Microarray data: Plot (a): Marginal posterior probabilities for pathway selection, $p(\theta_k|\mathbf{T}, Y)$. The 12 pathways with largest probabilities are marked with symbols. Plot (b): Conditional posterior probabilities for gene selection, $p(\gamma_j|\mathbf{T}, Y, I\{\sum_{k \in \mathcal{P}} \theta_k s_{kj} > 0\})$. The 41 probes with largest probability that belong to the 12 selected pathways in plot (a) are marked with Δ .

TABLE 1

The 41 selected genes divided by islands and with associated pathway indices (in parenthesis)

Singleton genes (no direct neighbor selected)

ACACB (10), C4A (8, 12), CALM1 (10), CCNB2 (5), CD4 (7), CDC2 (5), CLDN11 (7), FZD9 (11), GYS2 (10), HIST1H2BN (12), IFNA7 (3), NFASC (7), NRCAM (7), PCK1 (10), PFKP (10), PPARGC1A (10), PXN (9)

Island 1

ACTB (9), ACTG1 (9), ITGA1 (9), ITGA7 (9), ITGB3 (9), ITGB4 (9), ITGB6 (9), ITGB8 (7, 10), MYL5 (9), MYL9 (9), PDPK1 (10), PIK3CD (9, 10, 11), PLA2G4A (2), PLCG1 (11), PRKCA (2, 11), PRKY (2, 10), PRKY (2, 10), PTGS2 (11), SOCS3 (10)

Island 2

ACVR1B (2, 3, 11), ACVR1B (2, 3, 11), TGFB3 (2, 3, 5, 11)

Island 3

ENTPD3 (1), GMPS (1)

Notes: The pathway indices correspond to the following: 1-Purine metabolism, 2-MAPK signaling pathway, 3-Cytokine–cytokine receptor interaction, 4-Neuroactive ligand-receptor interaction, 5-Cell cycle, 6-Axon guidance, 7-Cell adhesion molecules (CAMs), 8-Complement and coagulation cascades, 9-Regulation of actin cytoskeleton, 10-Insulin signaling pathway, 11-Pathways in cancer, 12-Systemic lupus erythematosus.

From a practical point of view, researchers can use the posterior probabilities produced by our selection algorithm as a way to prioritize the relevant pathways and genes for further experimental work. For example, the genes corresponding to the best 41 selected probe sets, conditional upon the best 12 selected pathways, are listed in Table 1 divided by islands, which correspond to sets of connected genes in the Markov random field. The islands help with the biological interpretation by locating relevant branches of pathways. A subset of the selected pathways along with islands and singletons are displayed in Figure 4. Several of the identified pathways are involved in tumor formation and progression. For instance, the mitogen-activated protein kinase (MAPK) signaling pathway, involved in various cellular functions, including cell proliferation, differentiation and migration, has been implicated in breast cancer metastasis [Lee et al. (2007)]. The KEGG pathway in cancers was also selected with high posterior probability. Other interesting pathways are the insulin signaling pathway, which has been linked to the development, progression and outcome of breast cancer, and purine metabolism, involved in nucleotide biosynthesis and affects cell cycle activity of tumor cells.

In addition, several genes with known association to breast cancer were also selected. Protein kinase C alpha (PKCA), which belongs to the MAPK pathway and the KEGG pathways in cancer, has been reported to play roles in many different cellular processes, including cell functions associated with breast cancer progression. It has been shown to be overexpressed in some antiestrogen resistant breast cancer cell lines and to be involved in the growth of tamoxifen resistant human

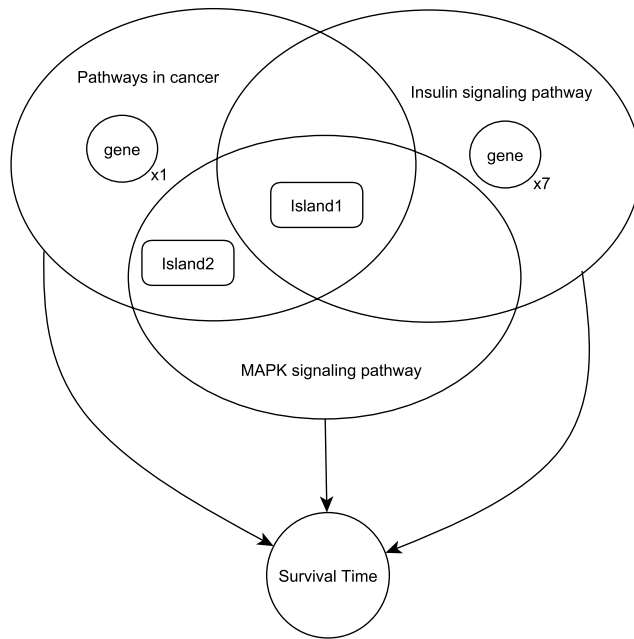


FIG. 4. *Microarray data: Graphical representation of a subset of selected pathways with islands and singletons. The genes in the islands are listed in Table 1.*

breast cancer cells [Frankel et al. (2007)]. Patients with PKCA-positive tumors have been shown to have worse survival than patients with PKCA-negative tumors, independently of other factors [Lønne et al. (2010)]. Prostaglandin-endoperoxide synthase-2 (PTGS2, also known as cyclooxygenase-2 or COX2) has also been related to breast cancer. Denkert, Winzer and Hauptmann (2004) observed COX2 overexpression in breast cancer and strong association with indicators of poor prognosis, such as lymph node metastasis, poor differentiation and large tumor size. This was further confirmed by Gupta et al. (2007), who showed that the expression of COX2 in human breast cancer cells facilitates the assembly of new tumor blood vessels, the release of tumor cells into the circulation, and the breaching of lung capillaries by circulating tumor cells to seed pulmonary metastasis. This is an important finding, as the majority of breast cancer deaths result from metastases rather than direct effects of the primary tumor. Another gene previously shown to be predictive of breast cancer lung metastasis is integrin, beta-8 (ITGB8) [Landemaine et al. (2008)]. We also identified integrin, beta-4 (ITGB4) which regulates key signaling pathways related to carcinoma progression, and is linked to aggressive tumor behavior and poor prognosis in certain breast cancer subtypes [Guo et al. (2006)].

5. Discussion. We have proposed a model that incorporates biological knowledge from pathway databases into the analysis of DNA microarrays to identify

pathways and genes related to a phenotype. Information on pathway membership and gene networks are used to define pathway summaries, specify prior distributions that account for the dependence structure between genes, and define the MCMC moves to fit the model. The gene network prior and the synthesis of the pathway information through PLS bring in additional information that is especially useful in microarray data, due to the low sample size and large measurement error. Performances of the method were evaluated on simulated data and a breast cancer gene expression study with survival outcomes was used to illustrate its application.

Our simulation studies show the effect of the MRF prior on the posterior inference. In general, as expected, the effect of the prior depends on the data and, in particular, on the concordance of the prior network with the data. In our simulations, employing the MRF prior allows us to achieve a better separation of the relevant pathways from those not relevant (in particular, we have found a larger average difference, over three scenarios, between the relevant pathway with the lowest posterior probability and the nonrelevant pathway with the highest posterior probability). In addition, in the simulated setting with fairly small regression coefficients the model with the MRF prior was able to select all the correct genes without any false positive, while the model without MRF includes 3 false positives. Other authors have reported improvements on selection power and sensitivity with respect to commonly used procedures that do not use the pathway structure, with similar, and in some cases, lower false discovery rates. In addition, in our formulation of the model we have used biological information not only for prior specification but also to structure the MCMC moves. This is helpful in arriving at promising models avoiding visiting invalid configurations. Finally, in real data applications, we have found that employing information on gene–gene networks can lead to the selection of significant genes that would have been missed otherwise, aiding the interpretation of the results, and achieving better predictions compared to models that do not treat genes as connected elements that work in groups or pathways.

Several MRF priors for gene selection indicators have been proposed in the literature. It is interesting to compare the parametrization of the MRF used in this paper and in [Li and Zhang \(2010\)](#) to the parametrization used in [Wei and Li \(2007, 2008\)](#), where the prior on $\boldsymbol{\gamma}$ is defined as

$$(19) \quad P(\boldsymbol{\gamma}|\cdot) \propto \exp(dn_1 - gn_{01}),$$

where n_1 is the number of selected genes and n_{01} is the number of edges linking genes with different values of γ_j , that is, edges linking included and nonincluded genes among all pathways,

$$n_1 = \sum_{j=1}^p \gamma_j, \quad n_{01} = \frac{1}{2} \sum_{i=1}^p \left[\sum_{j=1}^p r_{ij} - \left| \sum_{j=1}^p r_{ij}(1 - \gamma_j) - \sum_{j=1}^p r_{ij}\gamma_j \right| \right].$$

While d plays the same role as μ in (12), the parametrization using g has a different effect from η on the probability of selection of a gene. This is evident from the conditional probability $P(\gamma_j|\cdot, \gamma_i, i \in N_j) = \frac{\exp(\gamma_j F(\gamma_j))}{1 + \exp(F(\gamma_j))}$, where

$F(\gamma_j) = d + g \sum_{i \in N_j} (2\gamma_i - 1)$. Higher values of g encourage neighboring genes to take on the same γ_j value, and, consequently, genes with nonselected neighbors have lower prior probability of being selected than genes with no neighbors. We felt that parametrization (12) was a better choice for our purposes. First, in a context of sparsity, where only few nodes are supposed to take value 1, a prior that assigns larger probability of inclusion to genes with selected neighbors than to isolated genes seems more appropriate. Second, the exact simulation algorithm of [Propp and Wilson \(1996\)](#) cannot be used to simulate from (19). While any other method to draw from (19) would be acceptable, as said by [Møller et al. \(2006\)](#), Markov chain methods, to sample from a MRF, require to check at each step that the chain has converged to the equilibrium distribution, to avoid introducing additional undesirable stochasticity. On the other hand, one advantage of parametrization (19) is that no phase transition problem is associated to the distribution.

Pathway databases are incomplete and the gene network information is often unavailable for many genes. Thus, there may be situations where the dependence structure and the MRF prior specification on the gene selection indicator, $\boldsymbol{\gamma}$, cannot be used for all genes. When the only information available is the pathway membership of genes, the prior on $\boldsymbol{\gamma}$ could be elicited to capture other interesting characteristics. For example, a gene can have a priori higher probability of being selected when several pathways that contain it are included in the model. We may also want to avoid favoring the selection of a large pathway just because of its size. In such cases, conditional on $\boldsymbol{\theta}$, independent Bernoulli priors can be specified for γ_j relating the probability of selection to the proportion of included pathways that contain gene j , adjusting for the pathway sizes, p_k , that is, $\gamma_j | \boldsymbol{\theta} \sim \text{Bernoulli}(c \cdot \frac{\sum_{k=1}^K \theta_k s_{kj} / p_k}{\sum_{k=1}^K s_{kj} / p_k})$, with c a hyperparameter to be elicited.

In our approach we have made use of PLS components as summary measures of the expression of genes belonging to known pathways and then applied a fully Bayesian approach for the selection of the pathways to be included in the model, and the genes to be included within those pathways. Penalized techniques, including lasso [[Tibshirani \(1996\)](#)], elastic net [[Zou and Hastie \(2005\)](#)] and group lasso [[Yuan and Lin \(2006\)](#)] have been studied extensively in the literature and have been successfully applied to gene expression data. The group lasso, in particular, defines sets of variables, then selects either all the variables in the group or none of them. Recently, a modification of the method was proposed by [Friedman, Hastie and Tibshirani \(2010\)](#) using a more general penalty that yields sparsity at both the group and individual feature levels to select groups and predictors within each group. Our understanding of group lasso is that the method works best in situations where variables belonging to the same group are highly correlated, while covariates in different groups do not exhibit high correlation. However, genes belonging to the same pathway often do not exhibit high correlation in their expression levels. Also, in our case there are genes belonging to different pathways that have high correlation, as well as genes that belong to more than one pathway. Initial

investigations suggest that, in terms of prediction MSE, Bayesian formulations of lasso methods perform similarly to and, in some cases, better than the frequentist lasso [see, e.g., [Kyung et al. \(2010\)](#)]. Particularly relevant to our approach is the work of [Guan and Stephens \(2011\)](#), who apply Bayesian variable selection (BVS) and stochastic search methods in a regression model for genome-wide data. In simulations they find that, in spite of the apparent computational challenges, BVS produces better power and predictive performance compared with standard lasso techniques.

SUPPLEMENTARY MATERIAL

Supplement (DOI: [10.1214/11-AOAS463SUPP](https://doi.org/10.1214/11-AOAS463SUPP); .pdf). Description of the MCMC steps for (θ, γ) and discussion on ergodicity of the Markov chain on the restricted space.

REFERENCES

- ALBERT, J. H. and CHIB, S. (1993). Bayesian analysis of binary and polychotomous response data. *J. Amer. Statist. Assoc.* **88** 669–679. [MR1224394](#)
- ASHBURNER, M., BALL, C. A., BLAKE, J. A., BOTSTEIN, D., BUTLER, H., CHERRY, J. M., DAVIS, A. P., DOLINSKI, K., DWIGHT, S. S., EPPIG, J. T., HARRIS, M. A., HILL, D. P., ISSEL-TARVER, L., KASARSKIS, A., LEWIS, S., MATESE, J. C., RICHARDSON, J. E., RINGWALD, M., RUBIN, G. M. and SHERLOCK, G. (2000). Gene ontology: Tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.* **25** 25–29.
- BAIR, E., HASTIE, T., PAUL, D. and TIBSHIRANI, R. (2006). Prediction by supervised principal components. *J. Amer. Statist. Assoc.* **101** 119–137. [MR2252436](#)
- BESAG, J. (1974). Spatial interaction and the statistical analysis of lattice systems. *J. Roy. Statist. Soc. Ser. B* **36** 192–236. [MR0373208](#)
- BILD, A. H., YAO, G., CHANG, J. T., WANG, Q., POTTI, A., CHASSE, D., JOSHI, M.-B., HARPOLE, D., LANCASTER, J. M., BERCHUCK, A., OLSON, J. A. JR., MARKS, J. R., DRESSMAN, H. K., WEST, M. and NEVINS, J. R. (2006). Oncogenic pathway signatures in human cancers as a guide to targeted therapies. *Nature* **439** 353–357.
- BOULESTEIX, A.-L. and STRIMMER, K. (2007). Partial least squares: A versatile tool for the analysis of high-dimensional genomic data. *Brief. Bioinformatics* **8** 32–44.
- BROWN, P. J., VANNUCCI, M. and FEARN, T. (1998). Multivariate Bayesian variable selection and prediction. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **60** 627–641. [MR1626005](#)
- CHIPMAN, H., GEORGE, E. I. and MCCULLOCH, R. E. (2001). The practical implementation of Bayesian model selection. In *Model Selection. Institute of Mathematical Statistics Lecture Notes—Monograph Series* **38** 65–134. IMS, Beachwood, OH. [MR2000752](#)
- DAHLQUIST, K. D., SALOMONIS, N., VRANIZAN, K., LAWLOR, S. C. and CONKLIN, B. R. (2002). GenMAPP, a new tool for viewing and analyzing microarray data on biological pathways. *Nat. Genet.* **31** 19–20.
- DENKERT, C., WINZER, K.-J. and HAUPTMANN, S. (2004). Prognostic impact of cyclooxygenase-2 in breast cancer. *Clin. Breast Cancer* **4** 428–433.
- DONIGER, S., SALOMONIS, N., DAHLQUIST, K., VRANIZAN, K., LAWLOR, S. and CONKLIN, B. (2003). MAPPFinder: Using Gene Ontology and GenMAPP to create a global gene-expression profile for microarray data. *Genome Biology* **41** R7.
- DOWNWARD, J. (2006). Cancer biology: Signatures guide drug choice. *Nature* **439** 274–275.

- FRANKEL, L. B., LYKKESFELDT, A. E., HANSEN, J. B. and STENVANG, J. (2007). Protein Kinase C alpha is a marker for antiestrogen resistance and is involved in the growth of tamoxifen resistant human breast cancer cells. *Breast Cancer Res. Treat.* **104** 165–179.
- FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2010). A note on the group lasso and a sparse group lasso. Technical report, Dept. Stat., Stanford Univ.
- GEORGE, E. I. and MCCULLOCH, R. E. (1997). Approaches for Bayesian variable selection. *Statist. Sinica* **7** 339–373.
- GOLUB, T. R., SLONIM, D. K., TAMAYO, P., HUARD, C., GAASENBEEK, M., MESIROV, J. P., COLLER, H., LOH, M. L., DOWNING, J. R., CALIGIURI, M. A., BLOOMFIELD, C. D. and LANDER, E. (1999). Molecular classification of cancer: Class discovery and class prediction by gene expression monitoring. *Science* **286** 531–537.
- GUAN, Y. and STEPHENS, M. (2011). Bayesian variable selection regression for genome-wide association studies, and other large-scale problems. *Ann. Appl. Stat.* To appear.
- GUO, W., PYLAYEVA, Y., PEPE, A., YOSHIOKA, T., MULLER, W. J., INGHIRAMI, G. and GIANCOTTI, F. G. (2006). Beta 4 integrin amplifies ErbB2 signaling to promote mammary tumorigenesis. *Cell* **126** 489–502.
- GUPTA, G. P., NGUYEN, D. X., CHIANG, A. C., BOS, P. D., KIM, J. Y., NADAL, C., GOMIS, R. R., MANOVA-TODOROVA, K. and MASSAGUÉ, J. (2007). Mediators of vascular remodelling co-opted for sequential steps in lung metastasis. *Nature* **446** 765–770.
- JOSHI-TOPE, G., GILLESPIE, M., VASTRIK, I., D'EUSTACHIO, P., SCHMIDT, E., DE BONO, B., JASSAL, B., GOPINATH, G. R., WU, G. R., MATTHEWS, L., LEWIS, S., BIRNEY, E. and STEIN, L. (2005). Reactome: A knowledgebase of biological pathways. *Nucleic Acids Res.* **33** D428–D432.
- KANEHISA, M. and GOTO, S. (2000). Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **28** 27–30.
- KRIEGER, C., ZHANG, P., MUELLER, L., WANG, A., PALEY, S., ARNAUD, M., PICK, J., RHEE, S. and KARP, P. (2004). MetaCyc: A multiorganism database of metabolic pathways and enzymes. *Nucleic Acids Res.* **32** D438–442.
- KWON, D., TADESSE, M. G., SHA, N., PFEIFFER, R. M. and VANNUCCI, M. (2007). Identifying biomarkers from mass spectrometry data with ordinal outcome. *Cancer Inform.* **3** 19–28.
- KYUNG, M., GILL, J., GHOSH, M. and CASELLA, G. (2010). Penalized regression, standard errors, and Bayesian lassos. *Bayesian Anal.* **5** 369–412.
- LANDEMAINE, T., JACKSON, A., BELLAHCÈNE, A., RUCCI, N., SIN, S., ABAD, B. M., SIERRA, A., BOUDINET, A., GUINEBRETIERE, J.-M., RICEVUTO, E., NOGUÈS, C., BRIFOD, M., BIÈCHE, I., CHEREL, P., GARCIA, T., CASTRONOVO, V., TETI, A., LIDEREAU, R. and DRIOUCH, K. (2008). A six-gene signature predicting breast cancer lung metastasis. *Cancer Res.* **68** 6092–6099.
- LEE, S., JEONG, Y., IM, H. G., KIM, C., CHANG, Y. and LEE, I. (2007). Silibinin suppresses PMA-induced MMP-9 expression by blocking the AP-1 activation via MAPK signaling pathways in MCF-7 human breast carcinoma cells. *Biochemical and Biophysical Research Communications* **354** 65–171.
- LI, C. and LI, H. (2008). Network-constrained regularization and variable selection for analysis of genomics data. *Bioinformatics* **24** 1175–1182.
- LI, F. and ZHANG, N. (2010). Bayesian Variable selection in structured high-dimensional covariate space with application in genomics. *J. Amer. Statist. Assoc.* **105** 1202–1214.
- LINDGREN, F., GELADI, P. and WOLD, S. (1993). The kernel algorithm of PLS. *Journal of Chemometrics* **7** 45–59.
- LØNNE, G. K., CORNMARK, L., ZAHIROVIC, I. O., LANDBERG, G., JIRSTRÖM, K. and LARSSON, C. (2010). PKCalpha expression is a marker for breast cancer aggressiveness. *Mol. Cancer* **9** 76.

- LUCAS, J., CARVALHO, C., WANG, Q., BILD, A. Nevins, J. and West, M. (2006). Sparse statistical modelling in gene expression genomics. In *Bayesian Inference for Gene Expression and Proteomics* (K. Do, P. Mueller and M. Vannucci, eds.) 155–176. Cambridge Univ. Press, Cambridge.
- MØLLER, J., PETTITT, A. N., REEVES, R. and BERTHELSEN, K. K. (2006). An efficient Markov chain Monte Carlo method for distributions with intractable normalising constants. *Biometrika* **93** 451–458. [MR2278096](#)
- NAKAO, M., BONO, H., KAWASHIMA, S., KAMIYA, T., SATO, K., GOTO, S. and KANEHISA, M. (1999). Genome-scale gene expression analysis and pathway reconstruction in KEGG. *Genome Informatics Series: Workshop on Genome Informatics* **10** 94–103.
- PAN, W., XIE, B. and SHEN, X. (2010). Incorporating predictor network in penalized regression with application to microarray data. *Biometrics* **66** 474–484.
- PARK, M. Y., HASTIE, T. and TIBSHIRANI, R. (2007). Averaged gene expressions for regression. *Biostatistics* **8** 212–227.
- PITTMAN, J., HUANG, E., DRESSMAN, H., HORNG, C., CHENG, S., TSOU, M., CHEN, C., BILD, A., IVERSEN, E., HUANG, A., NEVINS, J. and WEST, M. (2004). Integrated modeling of clinical and gene expression information for personalized prediction of disease outcomes. *Proc. Natl. Acad. Sci. USA* **101** 8431–8436.
- PROPP, J. G. and WILSON, D. B. (1996). Exact sampling with coupled Markov chains and applications to statistical mechanics. In *Proceedings of the Seventh International Conference on Random Structures and Algorithms (Atlanta, GA, 1995)* **9** 223–252. [MR1611693](#)
- SHA, N., TADESSE, M. G. and VANNUCCI, M. (2006). Bayesian variable selection for the analysis of microarray data with censored outcomes. *Bioinformatics* **22** 2262–2268.
- SHA, N., VANNUCCI, M., TADESSE, M. G., BROWN, P. J., DRAGONI, I., DAVIES, N., ROBERTS, T. C., CONTESTABILE, A., SALMON, M., BUCKLEY, C. and FALCIANI, F. (2004). Bayesian variable selection in multinomial probit models to identify molecular signatures of disease stage. *Biometrics* **60** 812–828. [MR2089459](#)
- SHIPP, M. A., ROSS, K. N., TAMAYO, P., WENG, A. P., KUTOK, J. L., AGUIAR, R. C. T., GAASENBEEK, M., ANGELO, M., REICH, M., PINKUS, G. S., RAY, T. S., KOVAL, M. A., LAST, K. W., NORTON, A., LISTER, T. A., MESIROV, J., NEUBERG, D. S., LANDER, E. S., ASTER, J. C. and GOLUB, T. R. (2002). Diffuse large B-cell lymphoma outcome prediction by gene-expression profiling and supervised machine learning. *Nat. Med.* **8** 68–74.
- STINGO, F. and VANNUCCI, M. (2011). Variable selection for discriminant analysis with Markov random field priors for the analysis of microarray data. *Bioinformatics* **27** 495–501.
- STINGO, F., CHEN, Y., TADESSE, M. and VANNUCCI, M. (2011). Supplement to: “Incorporating biological information into linear models: A Bayesian approach to the selection of pathways and genes.” [DOI:10.1214/11-AOAS463SUPP](#).
- SUBRAMANIAN, A., TAMAYO, P., MOOTHA, V. K., MUKHERJEE, S., EBERT, B. L., GILLETTE, M. A., PAULOVICH, A., POMEROY, S. L., GOLUB, T. R., LANDER, E. S. and MESIROV, J. P. (2005). Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* **102** 15545–15550.
- TELESKA, D., MULLER, P., PARMIGIANI, G. and FREEDMAN, R. (2008). Modeling dependent gene expression. Technical report, Dept. of Biostatistics, Univ. Texas M.D. Anderson Cancer Center.
- TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58** 267–288. [MR1379242](#)
- TROYANSKAYA, O., CANTOR, M., SHERLOCK, G., BROWN, P., HASTIE, T., TIBSHIRANI, R., BOTSTEIN, D. and ALTMAN, R. B. (2001). Missing value estimation methods for DNA microarrays. *Bioinformatics* **17** 520–525.
- VAN’T VEER, L., DAI, H., VAN DE VIJVER, M., HE, Y., HART, A., MAO, M., PETERSE, H., VAN DER KOOPY, K., MARTON, M., WITTEVEEN, A., SCHREIBER, G., KERKHOVEN, R.,

- ROBERTS, C., LINSLEY, P., BERNARDS, R. and FRIEND, S. (2002). Gene expression profiling predicts clinical outcome of breast cancer. *Nature* **415** 530–536.
- WEI, L. J. (1992). The accelerated failure time model: A useful alternative to the Cox regression model in survival analysis. *Stat. Med.* **11** 1871–1879.
- WEI, Z. and LI, H. (2007). A Markov random field model for network-based analysis of genomic data. *Bioinformatics* **23** 1537–1544.
- WEI, Z. and LI, H. (2008). A hidden spatial-temporal Markov random field model for network-based analysis of time course gene expression data. *Ann. Appl. Stat.* **2** 408–429. [MR2415609](#)
- WOLD, H. (1966). Estimation of principal components and related models by iterative least squares. In *Multivariate Analysis (Proc. Internat. Sympos., Dayton, Ohio, 1965)* (P. Krishnaiah, ed.) 391–420. Academic Press, New York. [MR0220397](#)
- YUAN, M. and LIN, Y. (2006). Model selection and estimation in regression with grouped variables. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **68** 49–67. [MR2212574](#)
- ZHANG, J. D. and WIEMANN, S. (2009). KEGGgraph: A graph approach to KEGG PATHWAY in R and bioconductor. *Bioinformatics* **25** 1470–1471.
- ZOU, H. and HASTIE, T. (2005). Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **67** 301–320. [MR2137327](#)
- ZOU, H., HASTIE, T. and TIBSHIRANI, R. (2006). Sparse principal component analysis. *J. Comput. Graph. Statist.* **15** 265–286. [MR2252527](#)

F. C. STINGO
M. VANNUCCI
DEPARTMENT OF STATISTICS
RICE UNIVERSITY
HOUSTON, TEXAS 77005
USA
E-MAIL: fcs1@rice.edu
E-MAIL: marina@rice.edu

Y. A. CHEN
MOFFITT CANCER CENTER
TAMPA, FLORIDA 33612
USA
E-MAIL: Ann.Chen@moffitt.org

M. G. TADESSE
DEPARTMENT OF MATHEMATICS AND STATISTICS
GEORGETOWN UNIVERSITY
WASHINGTON, DC 20057
USA
E-MAIL: mgt26@georgetown.edu