

BAYESIAN SPARSE GRAPHICAL MODELS FOR CLASSIFICATION WITH APPLICATION TO PROTEIN EXPRESSION DATA

BY VEERABHADRAN BALADANDAYUTHAPANI^{1,3,*}, RAJESH TALLURI^{3,*},
YUAN JI^{2,†}, KEVIN R. COOMBES[‡], YILING LU^{*}, BRYAN T. HENNESSY^{§,4},
MICHAEL A. DAVIES^{*} AND BANI K. MALLICK[¶]

*The University of Texas M.D. Anderson Cancer Center**, *NorthShore
University HealthSystem and University of Chicago[†],*
*The Ohio State University[‡], Beaumont Hospital[§] and
Texas A&M University[¶]*

Reverse-phase protein array (RPPA) analysis is a powerful, relatively new platform that allows for high-throughput, quantitative analysis of protein networks. One of the challenges that currently limit the potential of this technology is the lack of methods that allow for accurate data modeling and identification of related networks and samples. Such models may improve the accuracy of biological sample classification based on patterns of protein network activation and provide insight into the distinct biological relationships underlying different types of cancer. Motivated by RPPA data, we propose a Bayesian sparse graphical modeling approach that uses selection priors on the conditional relationships in the presence of class information. The novelty of our Bayesian model lies in the ability to draw information from the network data as well as from the associated categorical outcome in a unified hierarchical model for classification. In addition, our method allows for intuitive integration of a priori network information directly in the model and allows for posterior inference on the network topologies both within and between classes.

Received February 2013; revised October 2013.

¹Supported in part by NIH Grant R01 CA160736 and the Cancer Center Support Grant (CCSG) (P30 CA016672).

²Supported by NIH R01 CA132897.

³Equal contributors.

⁴Supported by TRA (translational research award-TRA-2010-8) from the Health Research Board Ireland (HRB) and Science Foundation Ireland (SFI).

Key words and phrases. Bayesian methods, protein signaling pathways, graphical models, mixture models.

<p>This is an electronic reprint of the original article published by the Institute of Mathematical Statistics in <i>The Annals of Applied Statistics</i>, 2014, Vol. 8, No. 3, 1443–1468. This reprint differs from the original in pagination and typographic detail.</p>

Applying our methodology to an RPPA data set generated from panels of human breast cancer and ovarian cancer cell lines, we demonstrate that the model is able to distinguish the different cancer cell types more accurately than several existing models and to identify differential regulation of components of a critical signaling network (the PI3K-AKT pathway) between these two types of cancer. This approach represents a powerful new tool that can be used to improve our understanding of protein networks in cancer.

1. Introduction.

1.1. *Protein signaling pathways in cancer.* The treatment of cancer is rapidly evolving due to an improved understanding of the signaling pathways that are activated in tumors. Global profiling of DNA mutations, chromosomal copy number changes, DNA methylations and gene expression have greatly improved our appreciation of the heterogeneity of cancer [Nishizuka et al. (2003), Blower et al. (2007), Gaur et al. (2007), Shankavaram et al. (2007), Ehrlich et al. (2008)]. However, the characterization of protein signaling networks has proven to be much more challenging. Several reasons underscore the critical importance of overcoming this challenge: first, changes in cellular DNA and RNA both ultimately result in changes in protein expression and/or function, thus, protein networks represent the summation of changes that happen at the DNA and RNA levels. Second, research has demonstrated that many of the most common oncogenic genetic changes activate proteins in kinase signaling pathways. Numerous studies of protein networks and expression analysis have shown promising results. Due to the hyperactivation of kinase signaling pathways, numerous kinase inhibitors have been used in clinical trials, frequently with dramatic clinical activity. Inhibitors that target protein signaling pathways have been approved by the U.S. Food and Drug Administration for a variety of cancer types, including chronic myelogenous leukemia, breast cancer, colon cancer, renal cell carcinoma and gastrointestinal stromal tumors [as reviewed in Davies, Hennessy and Mills (2006)].

Protein networks need to be assessed directly, as DNA or RNA analyses often do not accurately reflect or predict the activation status of protein networks. Many proteins are regulated by post-translational modifications, such as phosphorylation or cleavage events, that are not detected by the analysis of DNA or RNA. Several studies have also demonstrated marked discordance between mRNA and protein expression levels, particularly for genes in kinase signaling and cell cycle regulation pathways [Varambally et al. (2005), Shankavaram et al. (2007)]. It has been demonstrated recently, in both cancer cell lines and tumors, that different genetic mutations in the

same signaling pathway can result in significant differences in the quantitative activation levels of downstream pathway effectors [Stemke-Hale et al. (2008), Davies et al. (2009), Vasudevan et al. (2009), Park et al. (2010)]. Although these observations support the suggestion that direct measurements are essential to measure protein network activation, a number of studies have demonstrated that signaling pathways are frequently regulated by complex feed-forward and feedback regulatory loops, as well as cross-talk between different pathways [Mirzoeva et al. (2009), Zhang et al. (2009), Halaban et al. (2010)]. Thus, developing an accurate understanding of the regulation of protein signaling networks will be optimized by approaches that: (1) assess multiple pathways simultaneously for different tumor types and/or conditions, and (2) allow for the use of rigorous statistical approaches to identify differential functional networks.

1.2. Reverse-phase protein lysate arrays. As explained, there is a strong rationale for methods that will directly assess the activation status of protein signaling networks in cancer. Traditional protein assays include immunohistochemistry (IHC), Western blotting, enzyme-linked immunosorbent assay (ELISA) and mass spectroscopy. Although IHC is a very powerful technique for the detection of protein expression and location, it is critically limited in network analyses by its non- to semi-quantitative nature. Western blotting can also provide important information, but due to its requirement for relatively large amounts of protein, it is difficult to use when comprehensively assessing protein networks, and also is semi-quantitative in nature. The ELISA method provides quantitative analysis, but is similarly limited by requirements of relatively high amounts of specimen and by the high cost of analyzing large pools of specimens. Mass spectroscopy is a powerful, quantitative approach, but its utility is mainly limited by the cost and time required to analyze individual samples, which limits the ability to run large sample sets that are needed to appropriately assess characteristics of disease heterogeneity and protein networks. Reverse-phase protein array (RPPA) analysis is a relatively new technology that allows for quantitative, high-throughput, time- and cost-efficient analysis of protein networks using small amounts of biological material [Paweletz et al. (2001); Tibes et al. (2006)].

RPPA data collection. We provide a brief overview of the RPPA experiment and data collection. In order to perform RPPA, proteins are isolated from the biological specimens such as cell lines, tumors or serum using standard laboratory-based methods. The protein concentrations are then determined for the samples and, subsequently, serial 2-fold dilutions prepared from each sample are then arrayed on a glass slide. Each slide is then probed with an antibody that recognizes a specific protein epitope that reflects the

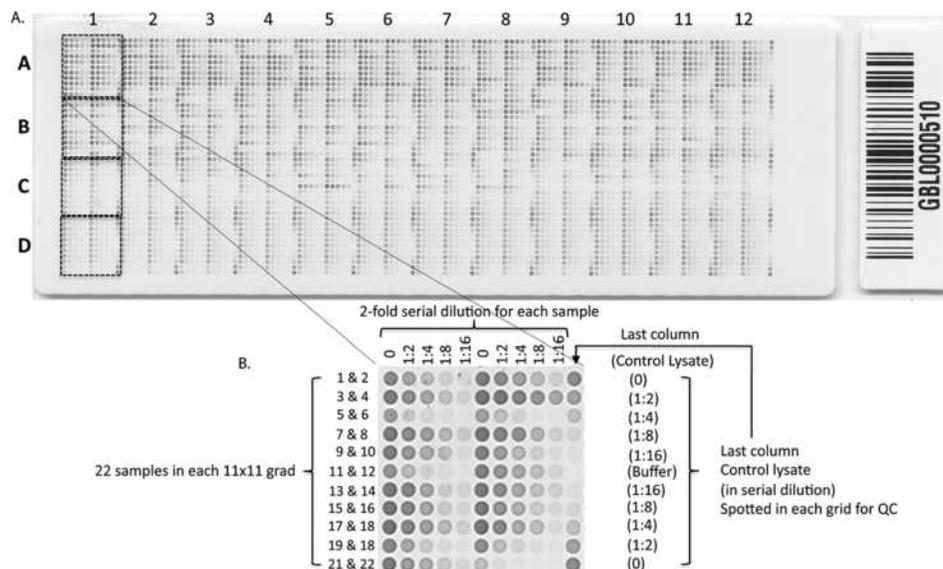


FIG. 1. An example of a reverse-phase protein array (RPPA) slide. (A) Each slide is comprised of 4 rows (A–D) of 12 columns (1–12) grids of 11X11 spots. (B) Each grid has 22 individual samples and 11 controls. Each row of the grid consists of 2 individual samples (each with 5 serial 2-fold dilutions) and one control spot. Reproduced with permission from Tabchy et al. (2011).

activation status of the protein. A visible signal is then generated through the use of a signal amplification system and staining. The signal reflects the relative amount of that epitope in each spot on the slide, as shown in Figure 1. The arrays are then scanned and the resulting images are analyzed with an imaging software specifically designed for the quantification of RPPA analysis (MicroVigene, VigeneTech Inc., Carlisle, MA). The relative signal intensities of the dilution series for each sample on the array are used to calculate the relative protein concentrations [Neeley et al. (2009), Zhang et al. (2009)]. Background correction is used to separate the signal from the noise by subtracting the extracted background intensity from the foreground intensity for each individual spot. Relative protein amount is calculated using a joint estimation method that utilizes the logistic model of Tabus et al. (2006). This method overcomes quenching at high levels and background noise at low levels. An R package, SuperCurve, developed to use with this joint estimation method is available at <http://bioinformatics.mdanderson.org/Software/OOMPA>. As with most high-throughput technologies, the normalization of the resulting intensities is conducted before any downstream analysis in order to adjust for sources of systematic variation not attributable to biological variation. Technical differences in protein loading for each sample are determined by first dividing the results for each protein measured by the

average value among all the specimens, and then by determining the average value for each sample across all of the measured proteins. This relative loading factor is then used to normalize the raw data for each sample, to correct for any differences in protein loading between specimens. We refer the reader to Paweletz et al. (2001) and Hennessy et al. (2010) for more biological and technical details concerning RPPAs.

Biological researchers typically choose specific targeted pathways containing 50–200 proteins, usually assayed using the same number of arrays, with each array hybridized against one protein. Because of the reverse design (as compared to conventional gene expression microarrays), RPPAs allow much larger sample sizes than the traditional microarrays, thus allowing *detailed* and *integrated* analyses of protein signaling networks with higher statistical power. Furthermore, this makes it possible to use RPPAs to measure protein expression for multiple tumor classes and/or cell conditions. The scientific aims we address using RPPA data in this paper are threefold: to infer differential networks between tumor classes/subtypes; to utilize a priori information in inferring protein network topology within tumor classes/subtypes; and, finally, to utilize network information in designing optimal classifiers for tumor classification. We believe this will improve our understanding of the regulation of protein signaling networks in cancer. Understanding the differences in protein networks between various cancer types and subtypes may allow for improved therapeutic strategies for each specific type of tumor. Such information may also be relevant when determining the origin of a tumor, which is clinically important in cases with indeterminate histologic analysis, particularly for patients who have more than one type of cancer.

1.3. *Graphical models for network analysis.* A convenient and coherent statistical representation of protein networks is accorded by graphical models [Lauritzen (1996)]. By “protein network” we mean any graph with proteins as nodes, where the edges between proteins may code for various biological information. For example, an edge between two proteins may represent the fact that their products interact physically (protein–protein interaction network), the presence of an interaction such as a synthetic-lethal or suppressor interaction [Kelley and Ideker (2005)], or the fact that these proteins code for enzymes that catalyze successive chemical reactions in a pathway [Vert and Kanehisa (2003)].

Our focus is on undirected graphical models and on Gaussian graphical models (GGM) in particular [Whittaker (1990)]. These models provide representations of the conditional independence structure of the multivariate distribution—to develop and infer protein networks. In such models, the nodes represent the variables (proteins) and the edges represent pairwise dependencies, with the edge set defining the global conditional independence structure of the distribution. We develop an adaptive modeling approach

for the covariance structure of high-dimensional distributions with a focus on sparse structures, which are particularly relevant in our setting in which the number of variables/proteins (p) can exceed the number of observations (n).

GGMs have been under intense methodological development over the past few years in both frequentist [Meinshausen and Bühlmann (2006), Chaudhuri, Drton and Richardson (2007), Yuan and Lin (2007), Friedman, Hastie and Tibshirani (2008), Bickel and Levina (2008)] and Bayesian settings [Giudici and Green (1999), Roverato (2002), Carvalho and Scott (2009)]. Wong, Carter and Kohn (2003) proposed a reversible jump MCMC-based Bayesian model for covariance selection. In high-dimensional settings, Dobra et al. (2004) used regression analysis to find directed acyclic graphs and converted them to undirected (sparse) graphs to explore the underlying network structure, and Rodríguez, Lenkoski and Dobra (2011) proposed a new approach for sparse covariance estimation in heterogeneous samples. However, most of the approaches we have cited focused on inferring the conditional independence structure of the graph and did not consider classification, which is one of the foci of our article. Rapaport et al. (2007) used spectral decomposition to detect the underlying network structure and classify genetic data using support vector machines (SVM). More recently, Monni and Li (2010) proposed a graph-based regression approach incorporating pathway information as a prior for classification procedures, however, their method does not detect differential networks based on available data. Zhu, Shen and Pan (2009) proposed network-based classification for microarray data using support vector machines. This was extended to network-based sparse Bayesian classifiers by Miguel Hernández-Lobato, Hernández-Lobato and Suárez (2011), but these approaches do not estimate the network and also do not take into account the differences in network structure between the two classes. Another recent method is that of Fan, Jin and Yao (2013), who propose a two-stage approach wherein they first select features and then subsequently use the retained features and Fisher’s LDA for classification using only one covariance matrix for both the classes.

In this article, we propose a constructive method for sparse graphical models using selection priors on the conditional relationships in the presence of class information. Our method has several advantages over classical approaches. First, we incorporate (integrate) the uncertainty of the parameters in deriving the optimal rule via Bayesian model mixing. Second, our network model provides an adaptively regularized estimate of the covariance matrix and hence is capable of handling $n < p$ situations. More importantly, our model uses this information in deriving the optimal classification boundary. The novelty of our Bayesian model lies in the ability to draw information from the network data from all the classes as well as from the associated categorical outcomes in a unified hierarchical model for classification. Through

this process, it offers the advantages of sparse Bayesian modeling of GGM, as well as the simplicity of a Bayesian classification model. In addition, with available online databases containing tens of thousands of reactions and interactions, there is a pressing need for methods integrating a priori pathway knowledge in the proteomic data analysis models. We integrate prior information directly in the model in an intuitive way such that the presence of an edge can be specified by providing the probability of an edge being present in the correlation matrix. Our method is fully Bayesian and allows for posterior inference on the network topologies both within and between classes. After fitting the Bayesian model, we obtain the posterior probabilities of the edge inclusion, which leads to false discovery rate (FDR)-based calls on significant edges.

The structure of our paper is as follows. In Section 2 we outline our Bayesian graph-based model for classification of RPPA data. Section 3 focuses on Bayesian FDR-based determination of significant networks. Section 4 presents the results of our case study using an RPPA experiment. We end with a discussion and conclusion in Section 5. All technical details and additional analysis results are presented in the supplementary material [Baladandayuthapani et al. (2014)].

2. Probability model. Our data construct for modeling is as follows. We observe a tuple: $(Z_i, \mathbf{Y}_i), i = 1, \dots, n$, where Z_i is a categorical outcome denoting the type or subtype of cancer (binary or multi category) and $\mathbf{Y}_i = (Y_i^{(1)}, \dots, Y_i^{(p)})$ is a p -dimensional vector of proteins assayed for the i th sample/patient/array. We detail the model here for binary classification (when Z_i is a binary variable), noting that generalization to multi-class classification can be achieved in an analogous manner. We factorize the joint distribution (likelihood) of the data $p(\mathbf{Y}_i, Z_i), \forall i$ in the following manner

$$p(\mathbf{Y}_i, Z_i) = p(\mathbf{Y}_i|Z_i)p(Z_i),$$

where the first component models the joint distribution of the p proteins given the class variable Z_i and the second component models the marginal distribution of the class variables. We model the first component as a mixture of the multivariate normal distributions as

$$p(\mathbf{Y}_i|Z_i, \boldsymbol{\mu}, \boldsymbol{\Omega}) \sim Z_i N(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Sigma}^{(1)}) + (1 - Z_i) N(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Sigma}^{(2)}),$$

where $\boldsymbol{\mu}^{(\bullet)}$ and $\boldsymbol{\Sigma}^{(\bullet)}$ are the corresponding means and covariances for the two classes. To specify the marginal component, we note that in the classification framework only a fraction of Z 's, say Z^u , will be unobserved (specifically in the case of prediction, as shown in Section 2.2) and they will be further modeled as

$$p(Z^u|h) \sim \text{Bernoulli}(h),$$

where we assign a Beta prior on probability h as $h \sim \text{Beta}(\eta, \zeta)$. Note that this prior can be generalized to be class-specific by allowing h to depend on the class k by changing the corresponding hyperparameters η_k, ζ_k .

Our main constructs of interest in this framework are $(\boldsymbol{\mu}^{(k)}, \boldsymbol{\Sigma}^{(k)})$, $k = 1, 2$ for each of the classes, where the latter provides a dependence structure between the proteins, which we model in a GGM framework. The key idea behind GGMs is rather to model the precision matrix $\boldsymbol{\Omega}^{(k)} = \boldsymbol{\Sigma}^{(k)^{-1}}$, which dictates the network structure between the variables. In this framework of particular interest is the identification of zero entries in the precision matrix—a zero entry at the ij th element of $\boldsymbol{\Omega}$ indicates conditional independence between the two random variables \mathbf{Y}_i and \mathbf{Y}_j , given all other variables. This is often referred to as the covariance selection problem in GGMs [Dempster (1972), Cox and Wermuth (2002)]. In the section below we provide a constructive method for sparse estimation (identification of many zeros) of the precision matrix in high-dimensional settings, but also allow for borrowing strength between classes to estimate the class-specific precision matrices for conducting classification.

2.1. Parameterization of the precision matrix. Given the number of variables p , the size of the precision matrix ($p \times p$) is potentially of high dimension. Instead of specifying a global (joint) distribution on the precision matrix, we explore local dependencies by breaking it down into components. For some applications, it is desirable to work directly with standard deviations and correlations [Barnard, McCulloch and Meng (2000), Liechty, Liechty and Müller (2004)] that do not correspond to any type of parameterization (e.g., Cholesky, etc.). This parameterization has a practical motivation because most biologists think in terms of correlations between the proteins, thus easing prior elicitation, as we show below. To this end, we parameterize the precision matrix (for each class k , suppressing the superscript for ease of notation) as $\boldsymbol{\Omega} = \mathbf{S} \times \mathbf{C} \times \mathbf{S}$, where \mathbf{S} is a diagonal matrix with nonzero diagonal elements that contains the inverse of the partial standard deviations and \mathbf{C} is a matrix that contains partial correlation coefficients. Note that the correlation matrix \mathbf{C} satisfies the properties of a correlation matrix, that is, the partial correlation coefficients (ρ_{ij}) between variables i, j share a one-to-one correspondence to the elements C_{ij} as

$$\rho_{ij} = \frac{-\Omega_{ij}}{(\Omega_{ii}\Omega_{jj})^{1/2}} = -C_{ij}.$$

Due to this correspondence, sparse estimation of $\boldsymbol{\Omega}$ directly implies the identification of zeros in the elements of \mathbf{C} . Thus, we model \mathbf{C} as a convolution,

$$\mathbf{C} = \mathbf{A} \odot \mathbf{R},$$

where \odot is the Hadamaard operator indicating element-wise multiplication between the two (stochastic) matrices: a *selection* matrix \mathbf{A} and the corresponding *correlation* matrix \mathbf{R} with the following properties:

- Both \mathbf{A} and \mathbf{R} are symmetric.
- Both \mathbf{A} and \mathbf{R} have ones as their diagonal elements.
- The off-diagonal elements of \mathbf{A} are either 0 or 1 and the off-diagonal elements of \mathbf{R} lie in the range $[-1, 1]$.
- Both \mathbf{A} and \mathbf{R} *need not* be positive definite, but the convolution \mathbf{C} *has to be* positive definite.

In essence, \mathbf{A} is a binary selection matrix that selects which of the elements in \mathbf{R} are zero or nonzero. In other words, \mathbf{A} performs variable selection on the elements of the matrix \mathbf{R} by shrinking the nonrequired variables (edges) exactly to zero and thus inducing sparsity in the estimation of the resulting precision matrix governing the GGM. We discuss hereafter the estimation and prior specifications for each of these matrices.

Prior construction. \mathbf{R} is a matrix with all of its off-diagonal elements in the range $[-1, 1]$, therefore, we assign an independent uniform prior over $[-1, 1]$ for all R_{ij} , $i < j$. Correspondingly, since the off-diagonal elements of \mathbf{A} are binary (0 or 1), we assign an independent Bernoulli prior with probability q_{ij} for the element A_{ij} , $i < j$. Note that this element-wise prior specification on \mathbf{A} and \mathbf{R} does not ensure that the \mathbf{C} ($=\mathbf{A} \odot \mathbf{R}$) is positive definite—hence a valid graph. Thus, a key ingredient of our modeling scheme is that we need an additional constraint: $\mathbf{C} \in \mathbb{C}_p$ where \mathbb{C}_p is the space of all proper correlation matrices of dimension p , such that the joint convolved prior on \mathbf{A} and \mathbf{R} can be written as

$$\mathbf{A}, \mathbf{R} | \mathbf{q} \sim \prod_{i < j} \{\text{Uniform}_{R_{ij}}[-1, 1] \text{Bernoulli}_{A_{ij}}(q_{ij})\} I(\mathbf{A} \odot \mathbf{R} \in \mathbb{C}_p),$$

where $I(\bullet)$, the indicator function, ensures that the correlation matrix is positive definite and introduces dependence among the elements of the matrices \mathbf{R}, \mathbf{A} , and q_{ij} is the probability of the ij th element being selected as 1.

We ensure the positive-definiteness constraint in our posterior sampling schemes. Specifically, we perform MCMC sampling in such a way that the constraint $\mathbf{C} \in \mathbb{C}_p$ is satisfied—to search over the possible space of valid correlation matrices. To implement the constraint, we draw R_{ij}, A_{ij} , sequentially conditioned on all other elements of \mathbf{R} and \mathbf{A} such that the realized value of C_{ij} ensures \mathbf{C} is positive definite given all other parameter values. Briefly, we follow the method of Barnard, McCulloch and Meng (2000) to find the range $[u_{ij}, v_{ij}]$ on the individual elements of R that will guarantee

the positive definiteness of \mathbf{C} . The resulting form of the conditional prior on the off-diagonal elements R_{ij} can be written as

$$R_{ij}|a_{ij}, A_{-ij}, R_{-ij} \sim \text{Uniform}(u_{ij}, v_{ij})I(-1 < R_{ij} < 1), \quad i \neq j, i < j,$$

where R_{-ij} contains all other off-diagonal elements of \mathbf{R} except the ij th element and A_{-ij} contains all elements of \mathbf{A} except the ij th element. The limits of the Uniform distribution u_{ij} and v_{ij} are chosen such that $\mathbf{C} = \mathbf{A} \odot \mathbf{R}$ is positive definite and (conditionally) u_{ij} and v_{ij} are functions of R_{-ij} and A_{-ij} (see Appendix A in the supplementary material [Baladandayuthapani et al. (2014)] for the detailed proof).

In this construction, the parameter probability q_{ij} controls the degree of sparsity in the GGM in an adaptive manner by element-wise selection of the entries of the correlation matrix. We assign a beta hyperprior for the probabilities q_{ij} as

$$q_{ij} \sim \text{Beta}(a_{ij}, b_{ij}), \quad i \neq j,$$

where the hyperparameters a_{ij}, b_{ij} can be set to induce prior information on the graph structure (see Section 2.3). To complete the hierarchical specification, we choose an (exchangeable) inverse-gamma prior on the inverse of the partial standard deviations S , which is a diagonal matrix containing entries $S_i = \Omega_{ii}^{1/2}$ as $S_i \sim IG(g, h)$, $i = 1, 2, \dots, p$.

Borrowing strength between classes. Note that in the above construction all the parameters are class-specific, that is, are different for each class k , and thus model fitting and estimation can be done for each class separately. But the main advantage of Bayesian methodology lies in borrowing strength between the classes for both estimation of the graphical structure and subsequent prediction/classification. This can be accomplished by having a variable that introduces dependence between the classes linking the selection matrix \mathbf{A} . We introduce a latent variable λ_{ij} defined as

$$\lambda_{ij} = \begin{cases} 1, & \text{if } A_{ij}^1 \neq A_{ij}^2, \\ 0, & \text{if } A_{ij}^1 = A_{ij}^2, \end{cases}$$

where \mathbf{A}^1 and \mathbf{A}^2 are the class-specific selection matrices. The binary variables λ_{ij} 's imply the presence or absence of the same edge in the graphical model of both classes. In other words, $\lambda_{ij} = 1$ signifies a *differential* edge (i.e., the relation between the covariates i, j is significant in only one class but not the other), whereas $\lambda_{ij} = 0$ signifies a *conserved* edge (i.e., the relation between the covariates i, j is significant in both classes). Thus, the λ 's serve a dual purpose in our model setup. They not only introduce dependence between the classes, since they are shared between both classes, but also have a distinct interpretation in terms of differential/conserved patterns

of dependence between the graphs for the classes. This information is vital for understanding the biological processes and inferring conclusions from the analysis, as we show in Section 4.

Since the λ_{ij} 's are binary random variables, we propose a Bernoulli prior on λ_{ij} as

$$\lambda_{ij} \sim \text{Bernoulli}(\pi_{ij}), \quad i < j,$$

where the parameter π_{ij} is the probability that the relation between the i th and j th variables is different. We further assign a beta hyperprior for the probabilities π_{ij} as

$$\pi_{ij} \sim \text{Beta}(e_{ij}, f_{ij}), \quad i \neq j.$$

To complete the prior specification on the graphical model, we propose a normal prior on the means $(\boldsymbol{\mu}^{(1)}, \boldsymbol{\mu}^{(2)})$ as

$$\boldsymbol{\mu}^{(k)} \sim N(\boldsymbol{\mu}_0^{(k)}, \mathbf{B}_0^{-1(k)}), \quad k = 1, 2.$$

2.2. Prediction. In this section we lay out our graph-based prediction (classification) scheme. Suppose the class variables \mathbf{Z} (of size $n \times 1$) are partitioned into a vector of training samples \mathbf{Z}^t (of size $n_t \times 1$) and a vector of (unknown) test/validation cases \mathbf{Z}^u (of size $n_u \times 1$) to be predicted. The corresponding observed variables are also partitioned as $[\mathbf{Y}^t; \mathbf{Y}^u]$. Denote the observed data by $\mathcal{D} = \{\mathbf{Y}^t, \mathbf{Z}^t, \mathbf{Y}^u\}$. In Bayesian prediction, for a new sample with protein expression information \mathbf{Y}^u , we have to obtain the posterior predictive probability that its class variable \mathbf{Z}^u , given all observed data \mathcal{D} , is $p(\mathbf{Z}^u | \mathcal{D})$.

To estimate these probabilities, we treat $\mathbf{Z}^u \equiv \{Z_o^u : o = 1, \dots, n_u\}$ as a parameter in the model and extend the MCMC analysis to sample these values at each iteration. Specifically, we draw \mathbf{Z}^u from the corresponding conditional posterior distribution in each MCMC iteration (see Appendix B in the supplementary material [Baladandayuthapani et al. (2014)] for the full conditional distribution). The way our model is specified, the posterior distribution of \mathbf{Z}^u is analyzed conditional not only on all the data from both classes \mathcal{D} , but also on the parameters that are shared between the classes. Thus, the predictions are obtained in a single MCMC fitting procedure along with all other parameters, thereby accounting for all sources of variation. We note that the limitation of this method is that training and test splits of the data must be contemplated prior to analysis (as is usually done) and/or analysis fully repeated if new predictions are required.

The complete hierarchical formulation of our graph-based binary classification model can be succinctly summarized as shown hereafter. In addition, the directed acyclic graph (Figure 6 in the supplementary material [Baladandayuthapani et al. (2014)]) shows a graphical representation of our model

where the circles indicate parameters and the squares observed random variables:

$$\begin{aligned}
\mathbf{Y} &= [\mathbf{Y}^t, \mathbf{Y}^u] \sim \mathbf{Z}N(\boldsymbol{\mu}^{(1)}, \boldsymbol{\Omega}^{-1(1)}) + (1 - \mathbf{Z})N(\boldsymbol{\mu}^{(2)}, \boldsymbol{\Omega}^{-1(2)}), \\
\mathbf{Z} &= [\mathbf{Z}^t, \mathbf{Z}^u], \\
Z_o^u &\sim \text{Bernoulli}(h_o), \\
h_o &\sim \text{Beta}(\eta, \zeta), \\
\boldsymbol{\mu}^{(k)} &\sim N(\boldsymbol{\mu}_0^{(k)}, \mathbf{B}_0^{-1(k)}), \\
\boldsymbol{\Omega}^{(k)} &= \mathbf{S}^{(k)}(\mathbf{A}^{(k)} \odot \mathbf{R}^{(k)})\mathbf{S}^{(k)}, \\
\mathbf{A}^{(k)}, \boldsymbol{\lambda}, \mathbf{R}^{(k)} | \mathbf{q}^{(k)}, \boldsymbol{\pi} &\sim \prod_{i < j} \text{Uniform}(u_{ij}, v_{ij}) \text{Bernoulli}(q_{ij}^{(k)}) \\
&\quad \times \text{Bernoulli}(\pi_{ij}) I(\mathbf{C}^{(k)} \in \mathbb{C}_p), \\
q_{ij}^{(k)} &\sim \text{Beta}(\alpha_{ij}^{(k)}, \beta_{ij}^{(k)}), \\
\pi_{ij} &\sim \text{Beta}(e_{ij}, f_{ij}), \quad i \neq j, \\
S_i^{(k)} &\sim IG(g, h),
\end{aligned}$$

where $k = 1, 2$ corresponds to the two classes, $i, j = 1, \dots, p$, and $o = 1, \dots, n_u$ corresponds to the size of the test/validation sample. The full conditional distributions for MCMC sampling of the model parameters and random variables are provided in Appendix B in the supplementary material [Baladandayuthapani et al. (2014)].

2.3. Incorporating prior pathway information and hyperparameter settings. As we mentioned before, there exists a huge amount of literature (prior knowledge) describing the functional behaviors of proteins, as characterized in metabolic, signaling and other regulation pathways. We formally incorporate this a priori knowledge in our model through the hyperparameter settings on the prior specification of q_{ij} , the probability that the edge between protein (i, j) will be selected. In particular, we impose an informative prior on $\pi(q_{ij}) \sim \text{Beta}(a_{ij}, b_{ij})$ and set the hyperparameters a_{ij} and b_{ij} such that the distribution has a higher mean to reflect our prior knowledge of the presence of an edge. For example, one could set the following:

- prior on q_{ij} as $\text{Beta}(2, 10)$ with mean 0.17, if there is biological evidence that the edge does not play an important role in the pathway;
- prior on q_{ij} as $\text{Beta}(10, 2)$ with mean 0.83, if there is biological evidence that the edge plays an important role in the pathway;
- prior on q_{ij} as $\text{Beta}(2, 2)$ with mean 0.5, if no prior information is available.

The prior information incorporated in the q_{ij} 's from online databases is assumed to represent normal conditions only. Information on relations between proteins that is affected by an intervention and/or mutation can be elicited by expert opinion (e.g., from a biologist). Information on the edges of graphs that is perturbed by a mutation can be incorporated formally through our prior on π_{ij} , the probability that controls the differential/conserved edge between two different conditions. We specify informative priors in a manner analogous to that of q_{ij} (as shown above) in cases where such information exists by setting e_{ij}, f_{ij} similarly to a_{ij} and b_{ij} . Finally, for the hyperparameters of the variance components, we obtain a vague inverse gamma prior by setting $(g, h) = 1$ and set the hyperparameters for the beta prior on h_o to be diffuse using $(\eta, \zeta) = 2$.

3. FDR-based determination of significant networks. Once we apply the MCMC methods, we are left with posterior samples of the model parameters that we can use to perform Bayesian inference. Our objective is twofold: to detect the “best” network/pathway based on the significance of the edges and also to detect differential networks between treatment classes. Given p proteins, our network consists of $p(p - 1)/2$ unique edges, which could be large even for a moderate number of proteins. Therefore, we need a mechanism that will control for these large-scale comparisons, discover edges that are significant and also detect differential edges between classes. We accomplish this in a statistically coherent manner using false discovery rate (FDR)-based thresholding to find significant networks and also to differentiate networks across samples.

The MCMC samples explore the distribution of possible network configurations suggested by the data, with each configuration leading to a different topology of the network based on the model parameters. Some edges that are strongly supported by the data may appear in most of the MCMC samples, whereas others with less evidence may appear less often. There are different ways to summarize this information in the samples. One could choose the most likely (posterior mode) network configuration and conduct conditional inference on this particular network topology. The benefit of this approach would be the yielding of a single set of defined edges, but the drawback is that the most likely configuration may still appear only in a very small proportion of MCMC samples. Alternatively, one could use all of the MCMC samples and, applying Bayesian model averaging (BMA) [Hoeting et al. (1999)], mix the inference over the various configurations visited by the sampler. This approach better accounts for the uncertainty in the data, leads to estimators of the precision matrix with the smallest mean squared error and should lead to better predictive performance in class predictions [Raftery, Madigan and Hoeting (1997)]. We will use this Bayesian model averaging approach.

From our MCMC iterations, suppose we have M posterior samples of the corresponding parameter set $\{A_{ij}^{(m)}, m = 1, \dots, M\}$, for which the selection indicator of the ij th edge is in the model. Suppose further that the model averaged set of posterior probabilities is set \mathcal{P} , the ij th element of which $\mathcal{P}_{ij} = M^{-1} \sum_m A_{ij}^{(m)}$ and is a $p \times p$ -dimensional matrix. Note that $1 - \mathcal{P}_{ij}$ can be considered Bayesian q -values, or estimates of the local false discovery rate [Storey and Tibshirani (2003), Newton et al. (2004)], as they measure the probability of a false positive if the ij th edge is called a “discovery” or is significant. Given a desired global FDR bound $\alpha \in (0, 1)$, we can determine a threshold ϕ_α with which to flag a set of edges $\mathcal{X}_\phi = \{(i, j) : \mathcal{P}_{ij} \geq \phi_\alpha\}$ as significant edges.

The significance threshold ϕ_α can be determined based on classical Bayesian utility considerations such as those described in Müller et al. (2004) and based on the elicited relative costs of false-positive and false-negative errors or can be set to control the average Bayesian FDR, as in Morris et al. (2008). The latter is the process we follow here. For example, suppose we are interested in finding the value ϕ_α that controls the overall average FDR at some level α , meaning that we expect that only $100\alpha\%$ of the edges that are declared significant are in fact false positives. Let $\text{vec}(\mathcal{P}) = [\mathcal{P}_t; t = 1, \dots, p(p-1)/2]$ be the vectorized version of the set \mathcal{P} containing the unique posterior probabilities of the edges, stacked columnwise. We first sort \mathcal{P}_t in descending order to yield $\mathcal{P}_{(t)}, t = 1, \dots, p(p-1)/2$. Then $\phi_\alpha = \mathcal{P}_{(\xi)}$, where $\xi = \max\{j^* : j^{*-1} \sum_{j=1}^{j^*} \mathcal{P}_{(t)} \leq \alpha\}$. The set of regions $\mathcal{X}_{\phi_\alpha}$ then can be claimed to be significant edges based on an average Bayesian FDR of α .

This FDR-based thresholding procedure can be extended to find differential networks between different populations (tumor classes/subtypes), for example, to identify edges that are significantly different between tumor types. To this end, we use the corresponding parameter set $\{\lambda_{ij}^{(m)}, m = 1, \dots, M\}$, which is the selection indicator of the differential edge between the ij th covariates in the model. The model averaged set of posterior probabilities is set \mathcal{P}^d , the ij th element of which $\mathcal{P}_{ij}^d = M^{-1} \sum_m \lambda_{ij}^{(m)}$. We use this same procedure to arrive at a set of differential edges $\mathcal{X}_\phi = \{(i, j) : \mathcal{P}_{ij}^d \geq \phi_\alpha\}$ with ϕ_α chosen to control the Bayesian FDR at level α . We use a similar procedure on the parameter set $\{1 - \lambda_{ij}^{(m)}, m = 1, \dots, M\}$, to arrive at a set of common edges $\mathcal{X}_\phi = \{(i, j) : \mathcal{P}_{ij}^c \geq \phi_\alpha\}$ with ϕ_α chosen to control the Bayesian FDR at level α .

4. Data analysis.

4.1. *Classification of breast and ovarian cancer cell lines.* Breast and ovarian cancer are two of the leading causes of cancer-related deaths in

women [Jemal et al. (2009)]. Both of these diseases are frequently affected by mutations in kinase signaling cascades, particularly those involving components of the PI3K-AKT pathway [Mills et al. (2003), Hennessy et al. (2008), Yuan and Cantley (2008), Bast Jr., Hennessy and Mills (2009)]. The PI3K-AKT pathway is one of the most important signaling networks in carcinogenesis [Vivanco and Sawyers (2002)] and is affected more than any other signaling pathway by activating mutations in cancer tissues [Yuan and Cantley (2008)]. Aggressive drug development efforts have targeted this critical oncogenic pathway, and inhibitors of multiple different components of the PI3K-AKT pathway have been developed and are in various stages of preclinical and clinical testing [Hennessy et al. (2005), Courtney, Corcoran and Engelman (2010)].

We applied our methodology to identify differences in the regulation of the PI3K-AKT signaling network in breast and ovarian cancers. For this analysis, we used expression data of $p = 50$ protein markers in signaling pathways from an RPPA analysis of human breast ($n_1 = 51$) and ovarian ($n_2 = 31$) cancer cell lines grown under normal tissue culture conditions [Stemke-Hale et al. (2008)]. We used the known connections in the PI3K-AKT pathway suggested by previous studies and expert opinion as a priori information in our model, as stated in Section 2.3.

The significant networks based on a Bayesian FDR cutoff of $\alpha = 0.1$ for breast and ovarian cancer samples are shown in Figure 2(a) and (b), respectively. The red edges indicate a negative association (regulation) and the green edges indicate a positive interaction between the proteins. The edges are represented by lines of varying degrees of thickness based on the strength of the association (correlation), with higher weights having thicker edges and lower weights having thinner edges. In order to identify biological similarities and differences between the breast and ovarian cancer cell lines, we compared the results of our network analyses of the two cancer types. Plotted in Figure 3(a) are the conserved (common) edges between the two cancer types. The differential network between the two cancer types, controlling for a Bayesian FDR cutoff of $\alpha = 0.1$, is shown in Figure 3(b).

A number of protein-protein relationships demonstrated significant similarity between the two cancer types. For example, both breast cancer and ovarian cancer cell lines exhibited a marked negative association between the levels of PTEN and phosphorylated AKT (Akt.pT308). This relationship was expected due to the critical regulation of 3-phosphatidylinositols by the lipid phosphatase activity of PTEN, and has previously been demonstrated as a significant interaction in multiple tumor types [Davies et al. (1998, 1999, 2009), Stemke-Hale et al. (2008), Vasudevan et al. (2009), Park et al. (2010)]. Although this concordance was expected, our analysis also identified a large network of differential protein interactions among the breast and ovarian cancer cell lines [Figure 3(b)]. In this figure, the

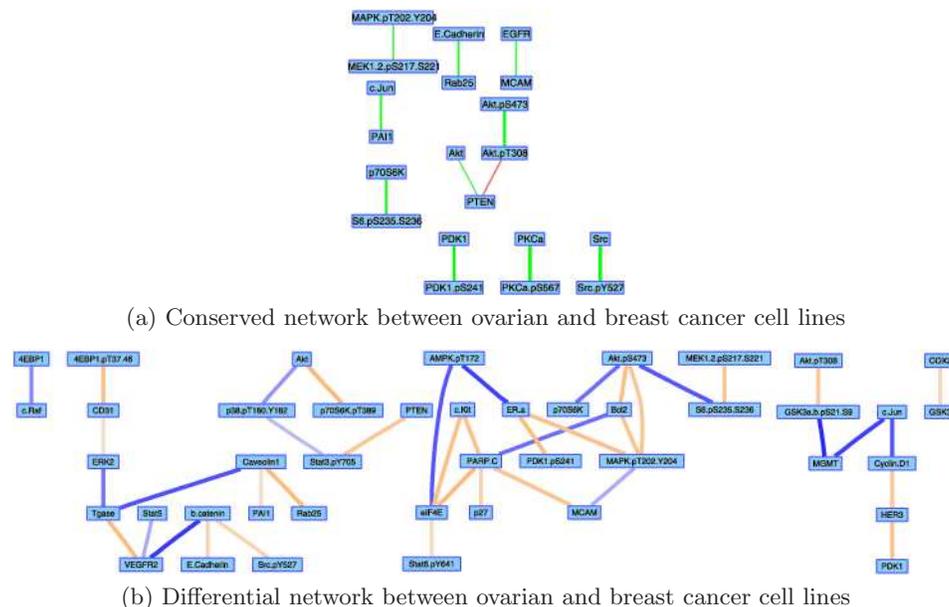


FIG. 3. Conserved and differential networks for the proteins in the PI3K-AKT kinase pathway between breast and ovarian cancer cell lines computed using a Bayesian FDR set to 0.10. In the conserved network (top panel), the red (green) lines between the proteins indicate a negative (positive) correlation between the proteins. In the differential network (bottom row), the blue lines between the proteins indicate a relationship that was significant in the ovarian cancer cell lines but not in the breast cancer cell lines; the orange lines between the proteins indicate a relationship in the breast cancer cell lines but not in the ovarian cancer cell lines. The thickness of the edges corresponds to the strength of the associations, with stronger associations having greater thickness.

edges in blue indicate relationships between proteins that were present in the ovarian cancer cell lines but not in the breast cancer cell lines using our FDR cutoff, and the orange edges indicate relationships present in the breast cancer cell lines but not in the ovarian cancer cell lines. In addition, the thickness of the edges corresponds to the strength of the association. Notable differential connections in this analysis include the association of phosphorylated AKT (Akt.pS473) with BCL-2 (Bcl2) and phosphorylated MAPK (MAPK.pT202.Y204) in breast cancer. Both of these, BCL-2(Bcl2) and phosphorylated (activated) MAPK (MAPK.pT202.Y204), may contribute to tumor proliferation and survival, and are therapeutic targets with available inhibitors. The association of different proteins with the expression of the estrogen receptor, phosphorylated PDK1 (PDK1.pS241) and MAPK (MAPK.pT202.Y204) in breast cancer and phosphorylated AMPK (AMPK.pT172) in ovarian cancer, may also have therapeutic implications, as the estrogen-receptor blockade is a treatment used in both advanced breast and ovarian cancer.

TABLE 1

Misclassification error rates for different classifiers for ovarian and breast cancer data sets. The methods compared here are SVM (network-based support vector machine), LDA (linear discriminant analysis), KNN (K-nearest neighbor), DQDA (diagonal quadratic discriminant analysis), DLDA (diagonal linear discriminant analysis), NBC (naive Bayes classifier) and BGBC (Bayesian graph-based classifier), which is the method proposed in this paper with and without incorporating prior information, denoted by BGBC (prior) and BGBC (w/o prior), respectively. The mean and the standard deviation are values of the misclassification percentage over 100 random splits of the data

	SVM	KNN	LDA	DLDA	DQDA	NBC	BGBC	BGBC w/o prior
Mean	8.03	15.14	25.48	12.07	13.74	13.37	6.59	10.88
SD	5.44	6.82	10.63	5.829	6.70	6.96	4.06	6.31

We used this network information to build a classifier to distinguish between breast cancer and ovarian cancer samples as explained in Section 2. We assessed the performance of the classifiers using cross-validation techniques. In particular, we generated 100 random selections of training and test data sets with 66% and 33% splits of the data, respectively. We fit our Bayesian graph-based classifier (BGBC) and compared our method to six other methods: the network-based support vector machine (SVM) [Zhu, Shen and Pan (2009)], K -nearest neighbor (KNN), linear discriminant analysis (LDA), diagonal linear discriminant analysis (DLDA), diagonal quadratic discriminant analysis (DQDA) and naive Bayes classifier (NBC) [John and Langley (1995)] methods. We implemented the network-based SVM using the R package “pathclass.” The network structure was specified to be the common network for the two classes obtained from the BGBC algorithm, as this method does not explicitly estimate the network. All other input parameters were set at the default settings for the network-based SVM function. We implemented all the other methods using the corresponding MATLAB functions.

The average misclassification errors (along with standard errors) across all splits for all the methods on the test set are shown in Table 1. The BGBC method had much lower misclassification rates compared to the other methods (the other methods ignore the underlying network structure of the proteins). We believe that this improved precision is due to the fact that the mean expression profiles of the breast and ovarian cancer cell lines are very similar so there is not enough information in the mean to classify the two cases. Hence, means-based classifiers, especially KNN and LDA (both of which use identity and diagonal covariances), underperform as compared to our method. The results of the DQDA method could be a bit closer to those of the BGBC method, but the former method ignores the cross-connections, that is, network information, and hence results in a higher misclassification

rate. The QDA could not be performed because the estimation of different covariance matrices for different classes is an ill-posed problem for $n < p$. We also tested the performance of BGBC using prior information and without using prior information in estimating the networks. The last two columns of Table 1 show that incorporating prior information improves our classification performance. Furthermore, the inclusion of prior information leads to sparser networks (as shown in Figure 7 in the supplementary material [Baladandayuthapani et al. (2014)]), as the prior information provides information about important and unimportant relationships, which aids our classification model.

We further note that nonlinear (quadratic) boundaries are obtained by using network information (since we model the covariance matrix), whereas linear boundaries are obtained by ignoring the network information (linear/diagonal discriminant based approaches). The classification boundary (Figure 8 in the supplementary material [Baladandayuthapani et al. (2014)]) exemplifies our intuition and approach. We have a $p(= 50)$ -dimensional quadratic classification boundary based on the GGM. In order to visualize this, we projected the boundary and the data onto two randomly selected dimensions/covariates. Two of those projections are shown in the figure, which confirm our intuition of how nonlinear boundary is more effective than a linear boundary in classification.

4.2. Effects of tissue culture conditions on network topology. Cell lines derived from tumors are a powerful research tool, as they allow for detailed characterization and functional testing. Genetic studies support the concept that cell lines generally mirror the changes that are detected in tumors, particularly at the DNA and RNA levels [Neve et al. (2006)]. However, the activation status of proteins can be impacted by the use of different environmental conditions in the culturing of cells. A key scientific question in the analysis of protein networks in cancer cell lines is the variability of network topologies due to differing tissue culture conditions. In order to assess if different network connectivity is observed under varying culture conditions, we used three different tissue culture conditions to grow the 31 ovarian cancer cell lines used in the previous analysis.

For condition “A,” the cells were grown in tissue culture media that was supplemented with growth factors in the form of fetal calf serum (5% of the total volume), which is a standard condition for the culturing of cancer cells. For condition “B,” the cells were harvested after being cultured in the absence of growth factors (serum) for 24 hours. For condition “C,” cells were grown in the absence of growth factors for 24 hours, then they were stimulated acutely (20 minutes) with growth factors (5% fetal calf serum). Proteins were harvested from each cell line for each tissue culture condition. The experimental procedure used for the isolation and RPPA analysis

of proteins from the cancer cells growing under normal, serum-replete tissue culture conditions has been described previously [Davies et al. (2009), Park et al. (2010)]. Protein isolation, processing and RPPA analysis were performed using the same methodology for all three conditions.

The RPPA data for each condition were then analyzed for protein–protein interactions using the GGM method. The topology maps for the ovarian cancer cells for the A, B and C tissue culture conditions are shown in Figure 12(a), (b) and (c) (provided in the supplementary material [Baladandayuthapani et al. (2014)]), respectively. We then performed comparisons of the results based on each of the three conditions in order to identify protein topology networks that were similar and different between each of the tissue culture conditions. As conditions A (media replete with growth factor) and B (media starved of growth factor) both represented steady-state tissue culture conditions, we initially compared these protein networks using a Bayesian FDR of 10%. The networks that are shared between the two conditions are shown in Figure 4(a); the differential associations are presented in Figure 4(d). We detected 21 significant protein interactions that were common for conditions A and B, and 4 interactions that were different. Thus, the overwhelming majority of protein–protein associations that were observed were maintained regardless of the presence or absence of growth factors (serum) in the tissue culture media. We then compared the significant relationships identified for condition B (media starved of growth factor) versus condition C (media starved, then acutely stimulated with growth factor). This comparison showed increased discordance of results, as we detected 20 associations that were common for conditions B and C [Figure 4(b)], but 11 associations that differed significantly [Figure 4(e)]. Similarly, the comparison of networks between the A and C conditions identified 22 shared protein interactions [Figure 4(c)] and 12 differential interactions [Figure 4(f)]. Of the differential interactions noted for the comparisons of conditions B versus C and A versus C, only 2 were observed in both comparisons (c-KIT and P38; VEGFR2 and MAPK.pT202.Y204). Neither of these 2 relationships was among the differential protein interactions in the analysis of condition A versus condition B. Of the 4 relationships that differed in the comparison of condition A versus condition B, 3 of the relationships were also identified as differing significantly when comparing condition B versus condition C (eIF4E and P38.pT180.Y182; c-Kit and PARP.cleaved; PARP.cleaved and ER.alpha), and the fourth differed significantly for the comparison of condition A versus condition C (AMPK.pT172 and eIF4E). This analysis suggests that protein–protein relationships are largely maintained under steady-state tissue culture conditions. However, these interactions may differ significantly in the setting of acute growth factor stimulation. We have included the explicit comparisons of our inferred networks with the prior PI3K-AKT pathway in Figures 13–16 in the supplementary material [Baladandayuthapani et

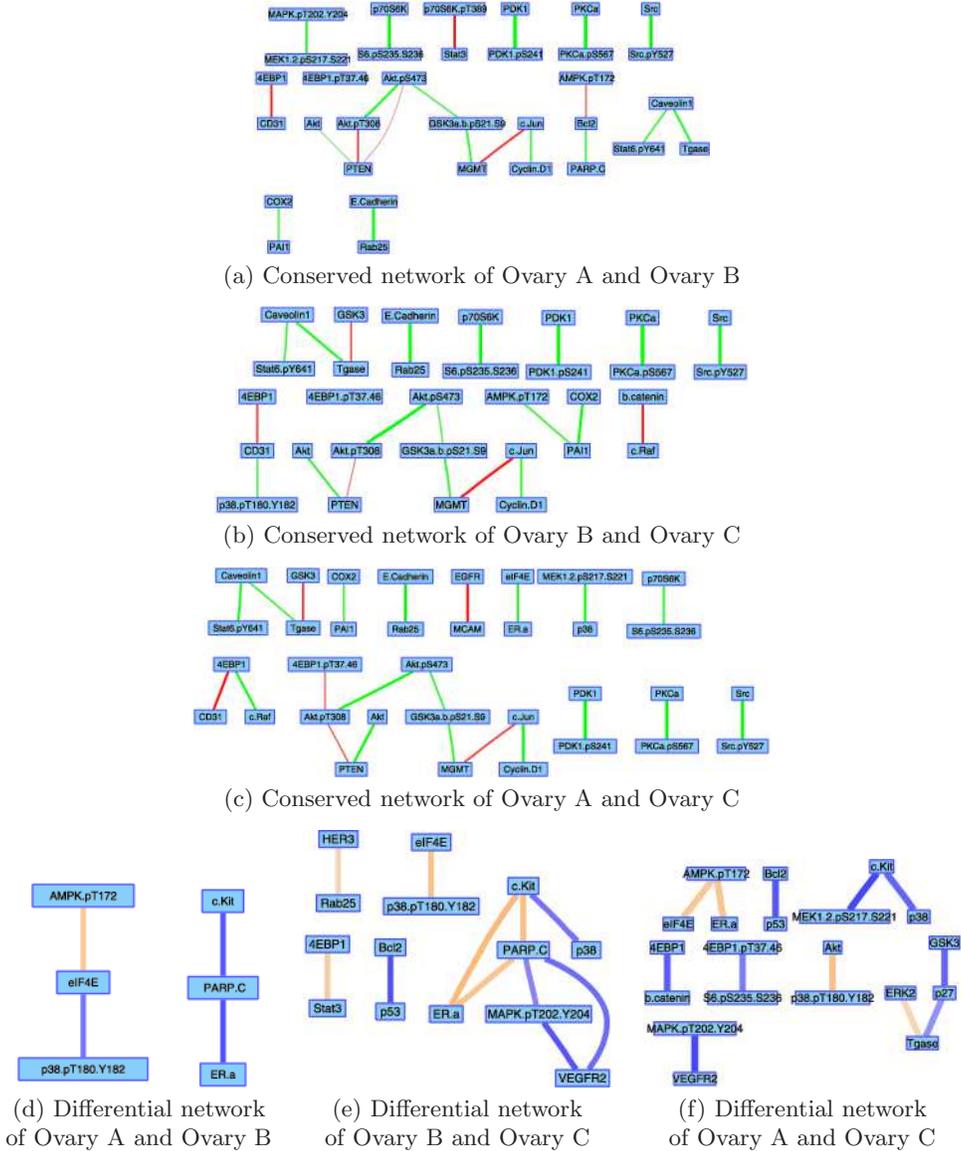


FIG. 4. Conserved and differential networks for the proteins in the PI3K-AKT kinase pathway between ovarian cancer cell lines grown in three different tissue culture conditions: A, B and C (see main text) computed using a Bayesian FDR set to 0.10. In the conserved network [(a)–(c)], the red (green) lines between the proteins indicate a negative (positive) correlation between the proteins. In the differential network [(d)–(f)], the blue lines between the proteins indicate a relationship that was significant in the ovarian cancer cell lines but not in the breast cancer cell lines; the orange lines between the proteins indicate a relationship in the breast cancer cell lines but not in the ovarian cancer cell lines. The thickness of the edges corresponds to the strength of the associations, with stronger associations having greater thickness.

al. (2014)]. The posterior means of the covariance matrices corresponding to the networks are also now plotted as heat maps in Figures 17–20 in the supplementary material [Baladandayuthapani et al. (2014)]. The exact posterior mean estimates are also provided as excel files downloadable from the corresponding authors’ website at http://odin.mdacc.tmc.edu/~vbaladan/Veera_Home_Page/Software_files/Covariance_Matrices.xlsx.

5. Discussion and conclusions. We present methodology to model sparse graphical models in the presence of class variables in high-dimensional settings, with a particular focus on protein signaling networks. Our methods allow for borrowing strength between classes to assess differential and common networks between the classes of cancer/tumor conditions. In addition, our method allows for the effective use of prior information about signaling pathways that is already available to us from various sources to help in decoding the complex protein networks. Improved understanding of the differential networks can be crucial for biologists when designing their experiments, allowing them to concentrate on the most important factors that distinguish tumor types. Such information may also help to narrow the drug targets for specific types of cancer. Knowledge of the common networks can be used to develop a drug for two different types of cancer that targets proteins that are active in both types. Data on the differential edges may be used as a good screening analysis, allowing researchers to eliminate unimportant proteins and concentrate on effective proteins when designing advanced patient-based translational experiments.

In this article we focused on undirected graphical models and not on directed (casual) networks. Directed graphical models, such as Bayesian networks and directed acyclic graphs (DAGs), have explicit causal modeling goals that require further modeling assumptions. In our formulation, we provide a natural and useful technical step in the identification of high posterior probability undirected graphical models, assuming a random sampling paradigm. In addition, our models infer network topologies that assume a steady-state network. Some of the protein networks may be dependent on causal relations between the nodes, which would require us to model data over time to infer the complete dynamics of the network. We leave this task for future consideration.

With regard to computation time, our MCMC chains are fairly fast for high-dimensional data sets such as those we considered, with a 5000-iteration run taking about 15 minutes. The source code, in MATLAB (The Mathworks, Inc., Natick, MA), takes advantage of several matrix optimizations available in that language environment. The computationally-involved step is the imposition of a positive definiteness on the correlation matrix. Optimizations to the code have been made by porting some functions into C. The software is available by contacting the first author.

Our main motivation for this work was to provide a constructive framework to conduct classification using sparse graphical methods that incorporate prior information. We assume parametric structures (likelihood/priors) throughout for ease of interpretation and computation, and our results indicate that this performs reasonably well on both real and simulated data sets. Extending to nonparametric settings would be an excellent avenue of future research that we would wish to undertake.

Acknowledgments. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Cancer Institute or the National Institutes of Health.

SUPPLEMENTARY MATERIAL

Supplement to “Bayesian sparse graphical models for classification with application to protein expression data” (DOI: [10.1214/14-AOAS722SUPP](https://doi.org/10.1214/14-AOAS722SUPP); .pdf). The supplementary material includes Appendix A: Positive definiteness constraint, Appendix B: Full conditional distributions and Appendix C: Simulations.

REFERENCES

- BALADANDAYUTHAPANI, V., TALLURI, R., JI, Y., COOMBS, K. R., LU, Y., HENNESSY, B. T., DAVIES, M. A. and MALLICK, B. K. (2014). Supplement to “Bayesian sparse graphical models for classification with application to protein expression data”. DOI:[10.1214/14-AOAS722SUPP](https://doi.org/10.1214/14-AOAS722SUPP).
- BARNARD, J., McCULLOCH, R. and MENG, X.-L. (2000). Modeling covariance matrices in terms of standard deviations and correlations, with application to shrinkage. *Statist. Sinica* **10** 1281–1311. [MR1804544](#)
- BAST C. R. JR., HENNESSY, B. and MILLS, G. B. (2009). The biology of ovarian cancer: New opportunities for translation. *Nat. Rev. Cancer* **9** 415–428.
- BICKEL, P. J. and LEVINA, E. (2008). Regularized estimation of large covariance matrices. *Ann. Statist.* **36** 199–227. [MR2387969](#)
- BLOWER, P. E., VERDUCCI, J. S., LIN, S., ZHOU, J., CHUNG, J.-H., DAI, Z., LIU, C.-G., REINHOLD, W., LORENZI, P. L., KALDJIAN, E. P., CROCE, C. M., WEINSTEIN, J. N. and SADEE, W. (2007). MicroRNA expression profiles for the NCI-60 cancer cell panel. *Mol. Cancer Ther.* **6** 1483–1491.
- CARVALHO, C. M. and SCOTT, J. G. (2009). Objective Bayesian model selection in Gaussian graphical models. *Biometrika* **96** 497–512. [MR2538753](#)
- CHAUDHURI, S., DRTON, M. and RICHARDSON, T. S. (2007). Estimation of a covariance matrix with zeros. *Biometrika* **94** 199–216. [MR2307904](#)
- COURTNEY, K. D., CORCORAN, R. B. and ENGELMAN, J. A. (2010). The PI3K pathway as drug target in human cancer. *J. Clin. Oncol.* **28** 1075–1083.
- COX, D. R. and WERMUTH, N. (2002). On some models for multivariate binary variables parallel in complexity with the multivariate Gaussian distribution. *Biometrika* **89** 462–469. [MR1913973](#)
- DAVIES, M., HENNESSY, B. and MILLS, G. B. (2006). Point mutations of protein kinases and individualised cancer therapy. *Expert Opin. Pharmacother.* **7** 2243–2261.

- DAVIES, M. A., LU, Y., SANO, T., FANG, X., TANG, P., LAPUSHIN, R., KOUL, D., BOOKSTEIN, R., STOKOE, D., YUNG, W. K., MILLS, G. B. and STECK, P. A. (1998). Adenoviral transgene expression of MMAC/PTEN in human glioma cells inhibits Akt activation and induces anoikis. *Cancer Res.* **58** 5285–5290.
- DAVIES, M. A., KOUL, D., DHESI, H., BERMAN, R., MCDONNELL, T. J., MCCONKEY, D., YUNG, W. K. and STECK, P. A. (1999). Regulation of Akt/PKB activity, cellular growth, and apoptosis in prostate carcinoma cells by MMAC/PTEN. *Cancer Res.* **59** 2551–2556.
- DAVIES, M. A., STEMKE-HALE, K., LIN, E., TELLEZ, C., DENG, W., GOPAL, Y. N., WOODMAN, S. E., CALDERONE, T. C., JU, Z., LAZAR, A. J., PRIETO, V. G., ALDAPE, K., MILLS, G. B. and GERSHENWALD, J. E. (2009). Integrated molecular and clinical analysis of AKT activation in metastatic melanoma. *Clin. Cancer Res.* **15** 7538–7546.
- DEMPSTER, A. P. (1972). Covariance Selection. *Biometrics* **28** 157–175.
- DOBRA, A., HANS, C., JONES, B., NEVINS, J. R., YAO, G. and WEST, M. (2004). Sparse graphical models for exploring gene expression data. *J. Multivariate Anal.* **90** 196–212. [MR2064941](#)
- EHRICH, M., TURNER, J., GIBBS, P., LIPTON, L., GIOVANNETTI, M., CANTOR, C. and VAN DEN BOOM, D. (2008). Cytosine methylation profiling of cancer cell lines. *Proc. Natl. Acad. Sci. USA* **105** 4844–4849.
- FAN, Y., JIN, J. and YAO, Z. (2013). Optimal classification in sparse Gaussian graphic model. *Ann. Statist.* **41** 2537–2571. [MR3161437](#)
- FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2008). Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* **9** 432–441.
- GAUR, A., JEWELL, D. A., LIANG, Y., RIDZON, D., MOORE, J. H., CHEN, C., AMBROS, V. R. and ISRAEL, M. A. (2007). Characterization of microRNA expression levels and their biological correlates in human cancer cell lines. *Cancer Res.* **67** 2456–2468.
- GIUDICI, P. and GREEN, P. J. (1999). Decomposable graphical Gaussian model determination. *Biometrika* **86** 785–801. [MR1741977](#)
- HALABAN, R., ZHANG, W., BACCHIOCCHI, A., CHENG, E., PARISI, F., ARIYAN, S., KRAUTHAMMER, M., MCCUSKER, J. P., KLUGER, Y. and SZNOL, M. (2010). PLX4032, a selective BRAF V600E kinase inhibitor, activates the ERK pathway and enhances cell migration and proliferation of BRAF WT melanoma cells. *Pigment Cell & Melanoma Research* **23** 190–200.
- HENNESSY, B. T., SMITH, D. L., RAM, P. T., LU, Y. and MILLS, G. B. (2005). Exploiting the PI3K/AKT pathway for cancer drug discovery. *Nat. Rev., Drug Discov.* **4** 988–1004.
- HENNESSY, B. T., MURPH, M., NANJUNDAN, M., CAREY, M., AUERSPERG, N., ALMEIDA, J., COOMBES, K. R., LIU, J., LU, Y., GRAY, J. W. and MILLS, G. B. (2008). Ovarian cancer: Linking genomics to new target discovery and molecular markers—the way ahead. *Adv. Exp. Med. Biol.* **617** 23–40.
- HENNESSY, B. T., LU, Y., GONZALEZ-ANGULO, A. M., CAREY, M. S., MYHRE, S., JU, Z., DAVIES, M. A., LIU, W., COOMBES, K., MERIC-BERNSTAM, F. et al. (2010). A technical assessment of the utility of reverse phase protein arrays for the study of the functional proteome in nonmicrodissected human breast cancers. *Clinical Proteomics* **6** 129–151.
- HOETING, J. A., MADIGAN, D., RAFTERY, A. E. and VOLINSKY, C. T. (1999). Bayesian model averaging: A tutorial. *Statist. Sci.* **14** 382–401. [MR1765176](#)
- JEMAL, A., SIEGEL, R., WARD, E., HAO, Y., XU, J. and THUN, M. J. (2009). Cancer statistics, 2009. *CA Cancer J. Clin.* **59** 225–249.

- JOHN, G. H. and LANGLEY, P. (1995). Estimating continuous distributions in Bayesian classifiers. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence* 338–345. Morgan Kaufmann, San Francisco, CA.
- KELLEY, R. and IDEKER, T. (2005). Systematic interpretation of genetic interactions using protein networks. *Nat. Biotechnol.* **23** 561–566.
- LAURITZEN, S. L. (1996). *Graphical Models*. Clarendon, Oxford.
- LIECHTY, J. C., LIECHTY, M. W. and MÜLLER, P. (2004). Bayesian correlation estimation. *Biometrika* **91** 1–14. [MR2050456](#)
- MEINSHAUSEN, N. and BÜHLMANN, P. (2006). High-dimensional graphs and variable selection with the lasso. *Ann. Statist.* **34** 1436–1462. [MR2278363](#)
- MIGUEL HERNÁNDEZ-LOBATO, J., HERNÁNDEZ-LOBATO, D. and SUÁREZ, A. (2011). Network-based sparse Bayesian classification. *Pattern Recogn.* **44** 886–900.
- MILLS, G. B., KOHN, E., LU, Y., EDER, A., FANG, X., WANG, H., BAST, R. C., GRAY, J., JAFFE, R. and HORTOBAGYI, G. (2003). Linking molecular diagnostics to molecular therapeutics: Targeting the PI3K pathway in breast cancer. *Semin. Oncol.* **30** 93–104.
- MIRZOEVA, O. K., DAS, D., HEISER, L. M., BHATTACHARYA, S., SIWAK, D., GENDELMAN, R., BAYANI, N., WANG, N. J., NEVE, R. M., GUAN, Y., HU, Z., KNIGHT, Z., FEILER, H. S., GASCARD, P., PARVIN, B., SPELLMAN, P. T., SHOKAT, K. M., WYROBEK, A. J., BISSELL, M. J., MCCORMICK, F., KUO, W.-L., MILLS, G. B., GRAY, J. W. and KORN, W. M. (2009). Basal subtype and MAPK/ERK kinase (MEK)-phosphoinositide 3-kinase feedback signaling determine susceptibility of breast cancer cells to MEK inhibition. *Cancer Res.* **69** 565–572.
- MONNI, S. and LI, H. (2010). Bayesian methods for network-structured genomics data. *UPenn Biostatistics Working Papers* **34**.
- MORRIS, J. S., BROWN, P. J., HERRICK, R. C., BAGGERLY, K. A. and COOMBES, K. R. (2008). Bayesian analysis of mass spectrometry proteomic data using wavelet-based functional mixed models. *Biometrics* **64** 479–489. [MR2432418](#)
- MÜLLER, P., PARMIGIANI, G., ROBERT, C. and ROUSSEAU, J. (2004). Optimal sample size for multiple testing: The case of gene expression microarrays. *J. Amer. Statist. Assoc.* **99** 990–1001. [MR2109489](#)
- NEELEY, E. S., KORNBLAU, S. M., COOMBES, K. R. and BAGGERLY, K. A. (2009). Variable slope normalization of reverse phase protein arrays. *Bioinformatics* **25** 1384–1389.
- NEVE, R. M., CHIN, K., FRIDLAND, J., YEH, J., BAEHNER, F. L., FEVR, T., CLARK, L., BAYANI, N., COPPE, J. P., TONG, F., SPEED, T., SPELLMAN, P. T., DEVRIES, S., LAPUK, A., WANG, N. J., KUO, W. L., STILWELL, J. L., PINKEL, D., ALBERTSON, D. G., WALDMAN, F. M., MCCORMICK, F., DICKSON, R. B., JOHNSON, M. D., LIPPMAN, M., ETHIER, S., GAZDAR, A. and GRAY, J. W. (2006). A collection of breast cancer cell lines for the study of functionally distinct cancer subtypes. *Cancer Cell* **10** 515–527.
- NEWTON, M. A., NOUEIRY, A., SARKAR, D. and AHLQUIST, P. (2004). Detecting differential gene expression with a semiparametric hierarchical mixture method. *Biostatistics* **5** 155–176.
- NISHIZUKA, S., CHEN, S.-T., GWADRY, F. G., ALEXANDER, J., MAJOR, S. M., SCHERF, U., REINHOLD, W. C., WALTHAM, M., CHARBONEAU, L., YOUNG, L., BUSSEY, K. J., KIM, S., LABABIDI, S., LEE, J. K., PITTALUGA, S., SCUDIERO, D. A., SAUSVILLE, E. A., MUNSON, P. J., PETRICOIN, E. F. I., LIOTTA, L. A., HEWITT, S. M., RAFFELD, M. and WEINSTEIN, J. N. (2003). Diagnostic markers that distinguish colon and ovarian adenocarcinomas: Identification by genomic, proteomic, and tissue array profiling. *Cancer Res.* **63** 5243–5250.

- PARK, E. S., RABINOVSKY, R., CAREY, M., HENNESSY, B. T., AGARWAL, R., LIU, W., JU, Z., DENG, W., LU, Y., WOO, H. G., KIM, S.-B., CHEONG, J.-H., GARRAWAY, L. A., WEINSTEIN, J. N., MILLS, G. B., LEE, J.-S. and DAVIES, M. A. (2010). Integrative analysis of proteomic signatures, mutations, and drug responsiveness in the NCI 60 cancer cell line set. *Mol. Cancer Ther.* **9** 257–267.
- PAWELETZ, C. P., CHARBONEAU, L., BICHSEL, V. E., SIMONE, N. L., CHEN, T., GILLESPIE, J. W., EMMERT-BUCK, M. R., ROTH, M. J., PETRICOIN, E. F. and LIOTTA, L. A. (2001). Reverse phase protein microarrays which capture disease progression show activation of pro-survival pathways at the cancer invasion front. *Oncogene* **20** 1981–1989.
- RAFTERY, A. E., MADIGAN, D. and HOETING, J. A. (1997). Bayesian model averaging for linear regression models. *J. Amer. Statist. Assoc.* **92** 179–191. [MR1436107](#)
- RAPAPORT, F., ZINOVYEV, A., DUTREIX, M., BARILLOT, E. and VERT, J.-P. (2007). Classification of microarray data using gene networks. *BMC Bioinformatics* **8** 35.
- RODRÍGUEZ, A., LENKOSKI, A. and DOBRA, A. (2011). Sparse covariance estimation in heterogeneous samples. *Electron. J. Stat.* **5** 981–1014. [MR2836767](#)
- ROVERATO, A. (2002). Hyper inverse Wishart distribution for nondecomposable graphs and its application to Bayesian inference for Gaussian graphical models. *Scand. J. Stat.* **29** 391–411. [MR1925566](#)
- SHANKAVARAM, U. T., REINHOLD, W. C., NISHIZUKA, S., MAJOR, S., MORITA, D., CHARY, K. K., REIMERS, M. A., SCHERF, U., KAHN, A., DOLGINOW, D., COSSMAN, J., KALDJIAN, E. P., SCUDIERO, D. A., PETRICOIN, E., LIOTTA, L., LEE, J. K. and WEINSTEIN, J. N. (2007). Transcript and protein expression profiles of the NCI-60 cancer cell panel: An integromic microarray study. *Mol. Cancer Ther.* **6** 820–832.
- STEMKE-HALE, K., GONZALEZ-ANGULO, A. M., LLUCH, A., NEVE, R. M., KUO, W.-L., DAVIES, M., CAREY, M., HU, Z., GUAN, Y., SAHIN, A., SYMMANS, W. F., PUSZTAI, L., NOLDEN, L. K., HORLINGS, H., BERN, K., HUNG, M.-C., VAN DE VIJVER, M. J., VALERO, V., GRAY, J. W., BERNARDS, R., MILLS, G. B. and HENNESSY, B. T. (2008). An integrative genomic and proteomic analysis of PIK3CA, PTEN, and AKT mutations in breast cancer. *Cancer Res.* **68** 6084–6091.
- STOREY, J. D. and TIBSHIRANI, R. (2003). Statistical significance for genomewide studies. *Proc. Natl. Acad. Sci. USA* **100** 9440–9445. [MR1994856](#)
- TABCHY, A., HENNESSY, B. T., GONZALEZ-ANGULO, A. M., BERNSTAM, F. M., LU, Y. and MILLS, G. B. (2011). Quantitative proteomic analysis in breast cancer. *Drugs of Today (Barcelona, Spain: 1998)* **47** 169–182.
- TABUS, I., HATEGAN, A., MIRCEAN, C., RISSANEN, J., SHMULEVICH, I. and WEI ZHANG ASTOLA, J. A. (2006). Nonlinear modeling of protein expressions in protein arrays. *IEEE Trans. Signal Process.* **54** 2394–2407.
- TIBES, R., QIU, Y., LU, Y., HENNESSY, B., ANDREEFF, M., MILLS, G. B. and KORNBLAU, S. M. (2006). Reverse phase protein array: Validation of a novel proteomic technology and utility for analysis of primary leukemia specimens and hematopoietic stem cells. *Mol. Cancer Ther.* **5** 2512–2521.
- VARAMBALLY, S., YU, J., LAXMAN, B., RHODES, D. R., MEHRA, R., TOMLINS, S. A., SHAH, R. B., CHANDRAN, U., MONZON, F. A., BECICH, M. J., WEI, J. T., PIANTA, K. J., GHOSH, D., RUBIN, M. A. and CHINNAIYAN, A. M. (2005). Integrative genomic and proteomic analysis of prostate cancer reveals signatures of metastatic progression. *Cancer Cell* **8** 393–406.
- VASUDEVAN, K. M., BARBIE, D. A., DAVIES, M. A., RABINOVSKY, R., MCNEAR, C. J., KIM, J. J., HENNESSY, B. T., TSENG, H., POCHANARD, P., KIM, S. Y., DUNN, I. F., SCHINZEL, A. C., SANDY, P., HOERSCH, S., SHENG, Q., GUPTA, P. B., BOEHM, J. S., REILING, J. H., SILVER, S., LU, Y., STEMKE-HALE, K., DUTTA, B., JOY, C.,

- SAHIN, A. A., GONZALEZ-ANGULO, A. M., LLUCH, A., RAMEH, L. E., JACKS, T., ROOT, D. E., LANDER, E. S., MILLS, G. B., HAHN, W. C., SELLERS, W. R. and GARRAWAY, L. A. (2009). AKT-independent signaling downstream of oncogenic PIK3CA mutations in human cancer. *Cancer Cell* **16** 21–32.
- VERT, J. P. and KANEHISA, M. (2003). Extracting active pathways from gene expression data. *Bioinformatics* **19** ii238–244.
- VIVANCO, I. and SAWYERS, C. L. (2002). The phosphatidylinositol 3-kinase AKT pathway in human cancer. *Nat. Rev. Cancer* **2** 489–501.
- WHITTAKER, J. (1990). *Graphical Models in Applied Multivariate Statistics*. Wiley, Chichester. [MR1112133](#)
- WONG, F., CARTER, C. K. and KOHN, R. (2003). Efficient estimation of covariance selection models. *Biometrika* **90** 809–830. [MR2024759](#)
- YUAN, T. L. and CANTLEY, L. C. (2008). PI3K pathway alterations in cancer: Variations on a theme. *Oncogene* **27** 5497–5510.
- YUAN, M. and LIN, Y. (2007). Model selection and estimation in the Gaussian graphical model. *Biometrika* **94** 19–35. [MR2367824](#)
- ZHANG, L., WEI, Q., MAO, L., LIU, W., MILLS, G. B. and COOMBES, K. (2009). Serial dilution curve: A new method for analysis of reverse phase protein array data. *Bioinformatics* **25** 650–654.
- ZHU, Y., SHEN, X. and PAN, W. (2009). Network-based support vector machine for classification of microarray samples. *BMC Bioinformatics* **10** S21.

V. BALADANDAYUTHAPANI
 R. TALLURI
 DEPARTMENT OF BIostatISTICS
 THE UNIVERSITY OF TEXAS
 M.D. ANDERSON CANCER CENTER
 HOUSTON, TEXAS 77030
 USA
 E-MAIL: veera@mdanderson.org
rtalluri@mdanderson.org

K. R. COOMBES
 DEPARTMENT OF BIOMEDICAL INFORMATICS
 THE OHIO STATE UNIVERSITY
 WEXNER MEDICAL CENTER
 COLUMBUS, OHIO 77030
 USA
 E-MAIL: kcoombes@mdanderson.org

B. T. HENNESSY
 BEAUMONT HOSPITAL
 DUBLIN
 IRELAND
 E-MAIL: bryanhennesy74@gmail.com

Y. JI
 NORTHSHORE UNIVERSITY HEALTHSYSTEM
 1001 UNIVERSITY PLACE
 EVANSTON, ILLINOIS 60201
 USA
 E-MAIL: jiyuan@uchicago.edu

Y. LU
 DEPARTMENT OF SYSTEMS BIOLOGY
 THE UNIVERSITY OF TEXAS
 M.D. ANDERSON CANCER CENTER
 HOUSTON, TEXAS 77030
 USA
 E-MAIL: yilinglu@mdanderson.org

M. A. DAVIES
 DEPARTMENT OF MELANOMA
 MEDICAL ONCOLOGY
 THE UNIVERSITY OF TEXAS
 M.D. ANDERSON CANCER CENTER
 HOUSTON, TEXAS 77030
 USA
 E-MAIL: madavies@mdanderson.org

B. K. MALLICK
 DEPARTMENT OF STATISTICS
 TEXAS A&M UNIVERSITY
 COLLEGE STATION, TEXAS 77843
 USA
 E-MAIL: bmallick@stat.tamu.edu