

A statistical framework for pathway and gene identification from integrative analysis



Quefeng Li^{a,b,*}, Menggang Yu^c, Sijian Wang^c

^a Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27517, USA

^b Statistical and Applied Mathematical Sciences Institute, Research Triangle Park, NC 27709, USA

^c Department of Biostatistics and Medical Informatics, University of Wisconsin at Madison, Madison, WI 53792, USA

ARTICLE INFO

Article history:

Received 15 June 2016

Available online 21 January 2017

AMS subject classification:

92B15

Keywords:

Gene and pathway

High dimensional analysis

Integrative analysis

Variable selection

ABSTRACT

In the era of big data, integrative analyses that pool data from different sources are now extensively conducted in order to improve performance. Among many interesting applications, genomics research is an area where integrative methods become popular tools to identify prognostic biomarkers for various diseases. In this paper, we propose such a framework for pathway and gene identification. Our method employs a hierarchical decomposition on genes' effects followed by a proper regularization to identify important pathways and genes across multiple studies. Asymptotic theories are provided to show that our method is both pathway and gene selection consistent. More importantly, we explicitly show that pathway selection consistency needs milder statistical conditions than gene selection consistency, as it would allow false positives and negatives at the gene selection level. Finite-sample performance of our method is shown to be superior than other ad hoc methods in various simulation studies. We further apply our method to analyze five cardiovascular disease studies. Our method is intrinsically a general method on group-wise and element-wise selections from integrative analysis, which can have other applications beyond genomic research.

© 2016 Elsevier Inc. All rights reserved.

1. Introduction

High-throughput sequencing technology has become a prevalent tool in biological and medical research. A unique characteristic of the sequencing data is that they are high-dimensional-low-sample-size (HDLSS), i.e., the number of biomarkers is substantially greater than the number of participants in the studies. On the other hand, the decisive biomarkers regulating the phenotypes are usually sparse compared with the total number of biomarkers along the whole genome and their effects are often weak; making the results from individual studies unremarkable and hard to be reproduced [19]. For this reason, joint analysis of multiple genomic data has been used widely and proven to be essential for identifying decisive biomarkers. A good example is the discovery of the risk loci for type 2 diabetes [22,30].

There are two types of joint analysis. One is the classical meta-analysis that aggregates summary statistics from individual datasets to obtain an overall score, based on which statistical significance across all studies is assessed. A comprehensive review of meta-analysis and its applications in genomic studies can be found in [24]. The other is the integrative analysis using individual patient data (IPD) from multiple studies [14–16]. In the era of big data, IPD becomes more accessible from

* Corresponding author at: Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC 27517, USA.
E-mail address: quefeng@email.unc.edu (Q. Li).

large genomic consortia such as Gene Expression Omnibus (GEO) and The Cancer Genome Atlas (TCGA). The integrative analysis links phenotypes to gene expression via a joint model of multiple studies and employs variable selection techniques to select decisive genes. It explores the relationship between phenotypes and genes in a more direct way and aims to identify genes that can better predict the phenotypes. Compared with the traditional meta-analysis, it is more appealing as the raw data contain more information than the summary statistics.

Groups of genes, commonly referred to as pathways, are often involved in the same biological processes. The pathway information has been well documented in publicly available databases such as KEGG [20] and Gene Ontology [2]. Since pathways play more defined biological roles than individual genes, pathway analysis has been advocated for better comprehension of the biological mechanism of the disease process. [13] provided a review of pathway analysis approaches developed in the recent decade, including popular methods such as Gene Set Enrichment Analysis [23] and Over-Representation Analysis [10,12]. Some recent advance in pathway analysis, e.g., a Bayesian method proposed in [3], further considers within-pathway selections of genes. In addition to pathways, this method identifies a list of genes under selected pathways that are associated with the response. However, all the above methods only focus on the analysis of a single dataset.

To the best of our knowledge, there are no existing methods that conduct pathway and within-pathway selections based on joint analysis of multiple datasets. In this paper, we propose such a new integrative analysis that utilizes a novel decomposition on genes' effects to reflect the relationship between pathways and genes. This decomposition, together with sparsity-induced penalties, enables both pathway and gene selections. A main difficulty of pathway selection is that the effects of pathways are unobserved. We show that once a reasonable identifiability condition is assumed, whether a pathway has nonzero effect is definitive. Hence, the pathway selection is sensible.

Under our framework, the employment of pathway information offers substantial benefit in statistical properties of the proposed method. We *explicitly* show that pathway selection requires much weaker conditions to achieve desired asymptotic properties. The main reason is that correct selection of pathways can allow some false positives and negatives at the gene selection level (see Section 4). In particular, for the pathway selection, we are able to avoid the irrepresentable condition that has been shown to be necessary for variable selection consistency in HDLSS settings, when either a convex or folded-concave penalty is used [6,31]. The reason is that our penalty's subgradient is infinitely large at the origin, leading to the removal of unimportant pathways without the irrepresentable condition. In addition, we show that correct pathway selection requires weaker conditions on the minimal signal strength than correct gene selection.

The rest of the paper is organized as follows. In Section 2, we introduce our method. In Section 3, we discuss the estimation and computation of our method. In Section 4, we provide theoretical results regarding pathway and gene selection consistency, where conditions needed for the two types of consistency are explicitly compared. In Section 5, we discuss how our method deals with the issue of pathway overlapping. In Section 6, we examine the finite-sample performance of our method through simulation studies. Examples of both overlapping and non-overlapping pathways are simulated. In Section 7, our method is applied to an integrative analysis of five cardiovascular studies. The paper is completed with a discussion in Section 8. All technical proofs are relegated to the [Appendix](#).

2. A hierarchical decomposition for pathway selection

We consider M independent studies with n_m subjects in study m . Let y_{mi} be a binary phenotype of the i th subject in the m th study and $\mathbf{x}_{mi} = (x_{mi,1}, \dots, x_{mi,d})^\top$ be the corresponding expression of d genes. Our proposed methods can be applied to data with any type of outcomes. We focus on the binary phenotype in this paper, because our motivational studies (see Section 7) have binary outcomes/phenotypes. It is natural to assume the following logistic regression model:

$$\log \frac{\Pr(y_{mi} = 1 | \mathbf{x}_{mi})}{\Pr(y_{mi} = 0 | \mathbf{x}_{mi})} = \alpha_m + \sum_{k=1}^K \sum_{j=1}^{G_k} x_{mi,kj} \beta_{kjm}, \quad k = 1, \dots, K; j = 1, \dots, G_k,$$

where α_m is the intercept, β_{kjm} is the effect of the j th gene under the k th pathway in the m th study, and they both are allowed to vary with m . K is the number of pathways and G_k is the number of genes under the k th pathway. We assume pathways do not have overlapping genes until Section 5 in which we then discuss the case of overlapping pathways.

We decompose β_{kjm} into three components as

$$\beta_{kjm} = p_k g_{kj} \zeta_{kjm}. \quad (1)$$

The parameter p_k is the effect of the k th pathway among all studies; g_{kj} is the effect of the j th gene under the k th pathway; ζ_{kjm} is the gene's effect in the m th study. Such a decomposition is inspired by the hierarchical LASSO [32], which targets for group selection within a single dataset. In decomposition (1), the exact values of p_k , g_{kj} and ζ_{kjm} are not identifiable (since $c p_k \cdot g_{kj} \zeta_{kjm} / c = p_k g_{kj} \zeta_{kjm}$ for a nonzero constant c). However, as will be shown in [Proposition 1](#), it is definitive whether they equal to zero once a reasonable identifiability condition is assumed.

A natural definition of an important pathway across multiple studies is that at least one of its genes has nonzero effect in at least one study. More explicitly, the set of important pathways is defined as

$$\mathcal{P} = \{k : \text{there exists at least one } (j, m) \text{ such that } \beta_{kjm}^* \neq 0\}, \quad (2)$$

where β_{kjm}^* denotes the true effect.

Next, we show that this definition is reasonable as it essentially requires an important pathway to have nonzero pathway effect, i.e., $p_k^* \neq 0$. In fact, the true effect β_{kjm}^* can always be decomposed hierarchically as $\beta_{kjm}^* = p_k^* g_{kj}^* \zeta_{kjm}^*$, where $(p_k^*, g_{kj}^*, \zeta_{kjm}^*)$ is a counterpart of $(p_k, g_{kj}, \zeta_{kjm})$. One possible decomposition is

$$p_k^* = \sum_{j=1}^{G_k} \sum_{m=1}^M |\beta_{kjm}^*|, \quad g_{kj}^* = \left(\sum_{m=1}^M |\beta_{kjm}^*| \right) / \left(\sum_{j=1}^{G_k} \sum_{m=1}^M |\beta_{kjm}^*| \right), \quad \zeta_{kjm}^* = \beta_{kjm}^* / \left(\sum_{m=1}^M |\beta_{kjm}^*| \right),$$

where $0/0$ is defined as 1. Even though the decomposition is not unique and therefore the exact values of $(p_k^*, g_{kj}^*, \zeta_{kjm}^*)$ are not definitive, we show that $I(p_k^* \neq 0)$ is definitive once a reasonable identifiability condition is imposed. In this sense, we are clear whether a pathway has a nonzero pathway effect. Moreover, we show that \mathcal{P} defined in (2) is the same as the support of \mathbf{p}^* (the set of nonzero elements of \mathbf{p}^*), where $\mathbf{p}^* = (p_1^*, \dots, p_K^*)^\top$. In other words, we require an important pathway to have nonzero pathway effect.

Proposition 1. *It holds that $\mathcal{P} = \{k : p_k^* \neq 0\}$, under the following condition*

(C1) *For any fixed (k, j) , if $\beta_{kjm}^* = 0$ for all $m = 1, \dots, M$, then $g_{kj}^* = 0$. In addition, for any fixed k , if $\beta_{kjm}^* = 0$ for all $j = 1, \dots, G_k$ and $m = 1, \dots, M$, then $p_k^* = 0$.*

Condition (C1) requires that if a gene’s effects are 0 in every study, its total gene effect should be 0. In addition, if all its member genes’ effects are 0 in every study, a pathway should have zero effect. Without this condition, p_k^* can be any number if $\beta_{kjm}^* = 0$ for all $j \in \{1, \dots, G_k\}$ and $m \in \{1, \dots, M\}$, as long as $g_{kj}^* = 0$ and $\zeta_{kjm}^* = 0$. In other words, a pathway’s effect is still undefined even though all its genes have zero effects in all studies. Hence, Condition (C1) is imposed to avoid the identifiability issue under the above circumstance.

3. Estimation and computation

Let $\ell_m(\alpha_m, \mathbf{p}, \mathbf{g}, \boldsymbol{\zeta}_m)$ be the log-likelihood of the m th study (divided by n_m) such that

$$\ell_m(\alpha_m, \mathbf{p}, \mathbf{g}, \boldsymbol{\zeta}_m) = \frac{1}{n_m} \sum_{i=1}^{n_m} \left[y_{mi} \left(\alpha_m + \sum_{k=1}^K \sum_{j=1}^{G_k} x_{mi,kj} p_k g_{kj} \zeta_{kjm} \right) - \log \left\{ 1 + \exp \left(\alpha_m + \sum_{k=1}^K \sum_{j=1}^{G_k} x_{mi,kj} p_k g_{kj} \zeta_{kjm} \right) \right\} \right],$$

where $\mathbf{p} = (p_1, \dots, p_K)^\top$, $\mathbf{g} = (g_{11}, \dots, g_{KG_K})^\top$, and $\boldsymbol{\zeta}_m = (\zeta_{11m}, \dots, \zeta_{KG_K m})^\top$. We impose L_1 penalties on $\mathbf{p}, \mathbf{g}, \boldsymbol{\zeta}_m$ and solve the following optimization problem:

$$(\hat{\boldsymbol{\alpha}}, \hat{\mathbf{p}}, \hat{\mathbf{g}}, \hat{\boldsymbol{\zeta}}) = \underset{\boldsymbol{\alpha}, \mathbf{p}, \mathbf{g}, \boldsymbol{\zeta}}{\operatorname{argmin}} \sum_{m=1}^M -\ell_m(\alpha_m, \mathbf{p}, \mathbf{g}, \boldsymbol{\zeta}_m) + \sum_{k=1}^K |p_k| + \sum_{k=1}^K \sum_{j=1}^{G_k} |g_{kj}| + \chi_n \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}|, \tag{3}$$

where χ_n is a positive tuning parameter, $\hat{\boldsymbol{\alpha}} = (\hat{\alpha}_1, \dots, \hat{\alpha}_M)^\top$ and $\hat{\boldsymbol{\zeta}} = (\hat{\zeta}_1^\top, \dots, \hat{\zeta}_M^\top)^\top$. The coefficient β_{kjm} is then estimated by $\hat{\beta}_{kjm} = \hat{p}_k \hat{g}_{kj} \hat{\zeta}_{kjm}$. Nonzero $\hat{\beta}_{kjm}$ ’s are considered as selected. Selection properties will be studied in Section 4.

Remark 1. In problem (3), the tuning parameters in front of penalties on \mathbf{p} and \mathbf{g} can always be normalized to be 1, due to the fact that $cp_k \cdot dg_{kj} \cdot (cd)^{-1} \zeta_{kjm} = p_k g_{kj} \zeta_{kjm}$, for nonzero constants c and d .

Similar as (2), we define pathways selected by our method as

$$\hat{\mathcal{P}} = \{k : \text{there exists at least one } (j, m) \text{ such that } \hat{\beta}_{kjm} \neq 0\}. \tag{4}$$

In Proposition 2, we show that $\hat{\mathcal{P}}$ is the same as the support of $\hat{\mathbf{p}}$. Theorem 1 in Section 4 shows that asymptotically $\hat{\mathcal{P}} = \mathcal{P}$ using definitions in (2) and (4). Therefore, it implies that $\{k : \hat{p}_k \neq 0\} = \{k : p_k^* \neq 0\}$. In other words, the solution $\hat{\mathbf{p}}$ given by our method can asymptotically recover the support of \mathbf{p}^* .

Proposition 2. $\hat{\mathcal{P}} = \{k : \hat{p}_k \neq 0\}$.

To solve (3), we propose fitting $\mathbf{p}, \mathbf{g}, \boldsymbol{\zeta}$ and $\boldsymbol{\alpha}$ iteratively. We fix three parameters and solve for the other and iterate between the steps until the algorithm converges. At each step, (3) is a convex problem, which can be efficiently solved by the coordinate descent algorithm [8,29]. The value of the objective function decreases after each iteration, hence the convergence is guaranteed. Let $\mathbf{X}_m \in \mathcal{R}^{n_m \times d}$ be the matrix containing data from the m th study. Denote $\mathbf{H} \in \mathcal{R}^{d \times K}$ as the indicator matrix of the membership of genes in different pathways such that its (j, k) th element $h_{jk} = 1$ if the j th gene belongs to the k th pathway and $h_{jk} = 0$ otherwise. Let $\hat{\boldsymbol{\alpha}}^{(s)}, \hat{\mathbf{p}}^{(s)}, \hat{\mathbf{g}}^{(s)}$, and $\hat{\boldsymbol{\zeta}}^{(s)}$ be the solutions of $\boldsymbol{\alpha}, \mathbf{p}, \mathbf{g}$ and $\boldsymbol{\zeta}$ at the s th step. We elaborate the algorithm as follows.

The algorithm of solving (3):

Step 1. For each dataset, standardize columns to have zero mean and unit variance. Initialize $\hat{\mathbf{g}}^{(0)}$ and $\mathbf{H}\hat{\mathbf{p}}^{(0)}$, e.g., let $\hat{\mathbf{g}}^{(0)} = \mathbf{1}$ and $\mathbf{H}\hat{\mathbf{p}}^{(0)} = \mathbf{1}$.

Step 2. In the s th iteration, let $\tilde{\mathbf{X}}_m = \mathbf{X}_m \cdot \hat{\mathbf{g}}^{(s-1)} \cdot \mathbf{H}\hat{\mathbf{p}}^{(s-1)}$, where \cdot represents product by columns. Denote

$$\ell(\boldsymbol{\zeta}) = \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \left[y_{mi} \left(\alpha_m^{(s-1)} + \sum_{k=1}^K \sum_{j=1}^{G_k} \tilde{x}_{mi,kj} \zeta_{kjm} \right) - \log \left\{ 1 + \exp \left(\alpha_m^{(s-1)} + \sum_{k=1}^K \sum_{j=1}^{G_k} \tilde{x}_{mi,kj} \zeta_{kjm} \right) \right\} \right].$$

Solve

$$\hat{\boldsymbol{\zeta}}^{(s)} = \underset{\boldsymbol{\zeta}}{\operatorname{argmin}} -\ell(\boldsymbol{\zeta}) + \chi_n \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}|.$$

Step 3. Update $\tilde{\mathbf{X}}_m$ by letting $\tilde{\mathbf{X}}_m = \mathbf{X}_m \cdot \hat{\boldsymbol{\zeta}}_m^{(s)} \cdot \mathbf{H}\hat{\mathbf{p}}^{(s-1)}$, where $\hat{\boldsymbol{\zeta}}_m^{(s)}$ is the subvector of $\hat{\boldsymbol{\zeta}}^{(s)}$ for the solutions in the m th dataset. Let

$$\ell(\mathbf{g}) = \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \left[y_{mi} \left(\alpha_m^{(s-1)} + \sum_{k=1}^K \sum_{j=1}^{G_k} \tilde{x}_{mi,kj} g_{kj} \right) - \log \left\{ 1 + \exp \left(\alpha_m^{(s-1)} + \sum_{k=1}^K \sum_{j=1}^{G_k} \tilde{x}_{mi,kj} g_{kj} \right) \right\} \right].$$

Solve

$$\hat{\mathbf{g}}^{(s)} = \underset{\mathbf{g}}{\operatorname{argmin}} -\ell(\mathbf{g}) + \sum_{k=1}^K \sum_{j=1}^{G_k} |g_{kj}|.$$

Step 4. Let $\mathbf{Z}_m = (\mathbf{X}_m \cdot \hat{\mathbf{g}}^{(s)} \cdot \hat{\boldsymbol{\zeta}}_m^{(s)}) \mathbf{H}$ and

$$\ell(\mathbf{p}) = \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \left[y_{mi} \left(\alpha_m^{(s-1)} + \sum_{k=1}^K z_{mk} p_k \right) - \log \left\{ 1 + \exp \left(\alpha_m^{(s-1)} + \sum_{k=1}^K z_{mk} p_k \right) \right\} \right],$$

where z_{mk} is the (m, k) th element of \mathbf{Z}_m . Solve

$$\hat{\mathbf{p}}^{(s)} = \underset{\mathbf{p}}{\operatorname{argmin}} -\ell(\mathbf{p}) + \sum_{k=1}^K |p_k|.$$

Step 5. Update $\tilde{\mathbf{X}}_m$ by letting $\tilde{\mathbf{X}}_m = \mathbf{X}_m \cdot \mathbf{H}\hat{\mathbf{p}}^{(s)} \cdot \hat{\mathbf{g}}^{(s)} \cdot \hat{\boldsymbol{\zeta}}_m^{(s)}$. Let

$$\ell(\boldsymbol{\alpha}) = \sum_{m=1}^M \frac{1}{n_m} \sum_{i=1}^{n_m} \left[y_{mi} \left(\alpha_m + \sum_{t=1}^d \tilde{x}_{mi,t} \right) - \log \left\{ 1 + \exp \left(\alpha_m + \sum_{t=1}^d \tilde{x}_{mi,t} \right) \right\} \right].$$

Solve $\hat{\boldsymbol{\alpha}}^{(s)} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \ell(\boldsymbol{\alpha})$.

Step 6. Let $\hat{\beta}_{kjm}^{(s)} = \hat{p}_k^{(s)} \hat{g}_{kj}^{(s)} \hat{\zeta}_{kjm}^{(s)}$. Return to step 2 and iterate, unless the difference between the values of the objective functions in (3) at two consecutive steps is less than some predefined threshold.

Next, we show that (3) has an equivalent form in terms of the original coefficient β_{kjm} .

Lemma 1. The problem (3) is equivalent to

$$(\hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\beta}}) = \underset{\boldsymbol{\alpha}, \boldsymbol{\beta}}{\operatorname{argmin}} \sum_{m=1}^M -\ell_m(\alpha_m, \boldsymbol{\beta}_m) + 3\chi_n^{1/6} \sum_{k=1}^K \left\{ \sum_{j=1}^{G_k} \left(\sum_{m=1}^M |\beta_{kjm}| \right)^{1/2} \right\}^{1/3}. \quad (5)$$

The equivalent form (5) describes our method in a different perspective. It shows that our method penalizes the negative log-likelihood with different penalties at different levels of the coefficients. Regarding $\sum_{m=1}^M |\beta_{kjm}|$ as the overall effect of the j th gene under the k th pathway across M studies, it penalizes the gene's overall effect with an $L_{1/2}$ -penalty. At the pathway level, an $L_{1/3}$ -penalty is imposed. Due to the singularities of $L_{1/2}$ and $L_{1/3}$ penalties at the origin point, our method can select pathways as well as its important members. In addition, the L_1 -penalty on the genes' effects in individual studies enables across-study selection for each gene and adapts to the heterogeneity among the studies. Due to the equivalence, (5) will be used to establish the theoretical results in Section 4.

It is easy to see that (5) is a non-convex problem thus hard to solve directly. However, due to the equivalence, it can be transferred into a sequence of convex optimization problems by using (3). This is another benefit brought by the hierarchical decomposition.

4. Theoretical properties

We present two results regarding the pathway and gene selection accuracy of our method. The results are considered for fixed M , but d and K can diverge with n , as long as $\log d = o(n)$. To simplify the presentation, we assume each dataset has n observations and the intercept terms $\{\alpha_{0m}^*\}_{m=1}^M$ are all 0.

We begin with introducing some notation. For any vector $\mathbf{a} \in \mathcal{R}^d$, let $\|\mathbf{a}\|_1 = \sum_{j=1}^d |a_j|$ be the L_1 -norm of \mathbf{a} and $\|\mathbf{a}\|_\infty = \max_j |a_j|$ be the sup-norm of \mathbf{a} . Define the support of \mathbf{a} as $\{j : a_j \neq 0\}$. For a symmetric matrix $\mathbf{A} \in \mathcal{R}^{d \times d}$, let $\|\mathbf{A}\|_\infty = \max_i \sum_{j=1}^d |a_{ij}|$ be the matrix sup-norm of \mathbf{A} . For two matrices \mathbf{A} and \mathbf{B} , denote $\mathbf{A} \circ \mathbf{B}$ as the Hadamard (componentwise) product of \mathbf{A} and \mathbf{B} . In addition, denote \mathbf{a}_S as the sub-vector of \mathbf{a} with indices in S and \mathbf{A}_S as columns of \mathbf{A} with indices in S . In particular, we denote $kj\cdot$ as the set of $\{kj1, \dots, kjM\}$. For a set S , we denote $|S|$ as its cardinality.

4.1. Pathway selection consistency

According to (2), to ensure that an important pathway is selected, it suffices to select one of its members in at least one study. In fact, we show that the member with the *largest* overall effect can be selected with large probability, where the overall effect of a gene is defined as the L_1 -norm of its effects in all studies. To this end, it only requires the *largest* overall effect not to be small, which is much weaker than the typical assumption on the minimal signal strength. On the other hand, as will be seen in the proof of Theorem 1, the unique property of the penalty term in (5) (its subgradient is infinitely large at the origin point) ensures that unimportant pathways are ruled out without additional assumptions. This is not like the L_1 or other folded-concave penalties, where irrerepresentable conditions are needed [6,31].

To present the result, let $K^* = M|\mathcal{P}|$, the number of important pathways multiplied by M , and

$$J = \left\{ (k, j, m) : k \in \mathcal{P}, j = \underset{j \in G_k}{\operatorname{argmax}} \|\boldsymbol{\beta}_{kj}^*\|_1, m = 1, \dots, M \right\},$$

the collection of $|\mathcal{P}|$ genes' effects in M studies, whose overall effect is the largest in their own pathways. Denote $B = \min_{k \in \mathcal{P}, j: kj \in J} \|\boldsymbol{\beta}_{kj}^*\|_1/2$, the minimal overall effects of these $|\mathcal{P}|$ genes. We write B and K^* in terms of orders of n as $B \asymp n^{-\alpha_B}$ and $K^* \asymp n^{-\alpha_{K^*}}$.

We present (5) in a matrix form as

$$\underset{\boldsymbol{\beta} \in \mathcal{R}^{Md}}{\operatorname{argmin}} \frac{1}{n} \{ -\mathbf{Y}^\top \mathbf{X} \boldsymbol{\beta} + \mathbf{e}^\top \mathbf{f}(\mathbf{X} \boldsymbol{\beta}) \} + \lambda_n \rho(\boldsymbol{\beta}), \tag{6}$$

where $\mathbf{X} \in \mathcal{R}^{Mn \times Md}$ is the block-diagonal design matrix, whose m th block contains data from the m th study, $\mathbf{Y} = (y_{11}, y_{12}, \dots, y_{Mn})^\top$, $\mathbf{f}(\boldsymbol{\theta}) = (f(\theta_1), \dots, f(\theta_{Mn}))^\top$ with $f(\theta_i) = \log\{1 + \exp(\theta_i)\}$, $\boldsymbol{\mu}(\boldsymbol{\theta}) = (f'(\theta_1), \dots, f'(\theta_{Mn}))^\top$, $\boldsymbol{\Sigma}(\boldsymbol{\theta}) = \operatorname{diag}(f''(\theta_1), \dots, f''(\theta_{Mn}))$, \mathbf{e} is a vector with all elements equal to 1, the penalty term $\rho(\boldsymbol{\beta}) = \sum_{k=1}^K \{ \sum_{j=1}^{G_k} \|\boldsymbol{\beta}_{kj}\|_1 \}^{1/2}$ and $\lambda_n = 3\chi_n^{1/6}$.

Theorem 1. Assume the following conditions hold.

(C2) $0 < \alpha_B < 1/2$;

(C3) $0 < \alpha_{K^*} < 1/2$;

(C4) $\|\{ \mathbf{X}_j^\top \boldsymbol{\Sigma}(\mathbf{X} \boldsymbol{\beta}^*) \mathbf{X}_j \}^{-1}\|_\infty = O(n^{-1})$;

(C5) $\max_{\delta \in \mathcal{N}} \max_{k,j,m} \lambda_{\max}(\mathbf{X}_j^\top \operatorname{diag}\{ |\mathbf{X}_{kjm}| \circ |\boldsymbol{\mu}''(\mathbf{X}_j \delta)| \} \mathbf{X}_j) = O(n)$, where $\mathcal{N} = \{ \delta \in \mathcal{R}^{K^*} : \|\delta - \boldsymbol{\beta}^*\|_\infty \leq B \}$.

If we choose λ_n such that $\lambda_n = o(B^{5/6} n^{-1/2})$, $\lambda_n \kappa_n = o(\tau_n)$, where

$$\kappa_n = \max_{\delta \in \mathcal{N}} (5/36) \|\delta\|_1^{-11/6},$$

and $\tau_n = \min_{\delta \in \mathcal{N}} \lambda_{\min}(n^{-1} \mathbf{X}_j^\top \boldsymbol{\Sigma}(\mathbf{X}_j \delta) \mathbf{X}_j)$, then for sufficiently large n , with probability greater than $1 - 2K^*/n$, there exists a solution $\hat{\boldsymbol{\beta}}$ to (5), such that

(a) $\hat{\mathcal{P}} = \mathcal{P}$;

(b) $\|\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j^*\|_\infty \leq cn^{-\gamma}$, where c is a generic positive constant and γ is some number that is arbitrarily smaller than $1/2$.

Remark 2. Theorem 1 implies that our method is asymptotically pathway selection consistent and the corresponding estimation of $\boldsymbol{\beta}_j^*$ is uniformly consistent. In general, conditions needed for pathway selection consistency is weaker than that for gene selection consistency (see Theorem 2). For example, condition (C2) only requires the minimum of the *largest* overall effects under important pathways not to be small, compared with a stronger condition (C6) needed for gene selection consistency. Moreover, Theorem 1 does not need the irrerepresentable condition [6,31] to reach pathway selection consistency, while this condition (C9) is needed to achieve gene selection consistency. The intuition is that, pathway selection consistency is easier to be achieved than gene selection consistency. For important pathways, correct pathway selection can tolerate false positive and negative selections of their genes, as long as at least one gene under the important pathway is selected in at least one study. For unimportant pathways, our method is able to remove them without an irrerepresentable condition due to the fact that the subgradient of our penalty at the origin is infinitely large.

4.2. Gene selection consistency

If stronger conditions are satisfied, we show that our method is further gene selection consistent in *every study* in the sense that the support of $\hat{\beta}$ is the same as that of β^* . The genes' effects are often heterogeneous among various studies due to many practical factors, such as various types of biospecimens, inconsistent experimental environments, etc. For this reason, we allow a gene's effect to be zero in some studies but nonzero in others. The nonzero effects are also allowed to vary. The challenge of achieving gene selection consistency is, for genes have zero/nonzero effects in different studies, to identify the studies where the gene has zero effect. The gene selection consistency is stronger than pathway selection consistency, since the former implies the latter.

To present the result, we divide elements of β^* into three subsets:

$$I = \{(k, j, m) : \beta_{kjm}^* \neq 0\}; \quad II = \{(k, j, m) : \beta_{kjm}^* = 0, \|\beta_{kj\cdot}^*\|_1 \neq 0\};$$

$$III = \{(k, j, m) : \beta_{kjm}^* = 0, \|\beta_{kj\cdot}^*\|_1 = 0\}.$$

Denote $s = |I|$, the cardinality of set I and $b = (1/2) \min\{|\beta_{kjm}^*| : (k, j, m) \in I\}$, the minimal nonzero element of β^* ,

$$h_n = \min_{(k,j,m) \in I} \left(\sum_{j=1}^{G_k} \|\beta_{kj\cdot}^*\|_1^{1/2} \right)^{2/3} \|\beta_{kj\cdot}^*\|_1^{1/2}, \quad u_n = \max_{(k,j,m) \in I} \left(\sum_{j=1}^{G_k} \|\beta_{kj\cdot}^*\|_1^{1/2} \right)^{2/3} \|\beta_{kj\cdot}^*\|_1^{1/2}.$$

The magnitudes of b, s, d are quantified in the order of n as $b \asymp n^{-\alpha_b}$, $s \asymp n^{\alpha_s}$, and $\log d \asymp n^{\alpha_d}$, where α_s, α_u and α_d are positive numbers.

Theorem 2. Under the following conditions

(C6) $0 < \alpha_b < 1/2$;

(C7) $0 < \alpha_s < 1/2$ and $0 < \alpha_d < 1$;

(C8) $\|\{X_I^T \Sigma(X\beta^*)X_I\}^{-1}\|_\infty = O(n^{-1})$;

(C9) $\|X_{II}^T \Sigma(X\beta^*)X_I \{X_I^T \Sigma(X\beta^*)X_I\}^{-1}\|_\infty < h_n/(10u_n)$;

(C10) $\max_{\delta \in \mathcal{N}_0} \max_{k,j,m} \lambda_{\max}(X_I^T \text{diag}\{|\mathbf{X}_{kjm}| \circ |\mu''(\mathbf{X}_I \delta)|\} X_I) = O(n)$, $\mathcal{N}_0 = \{\delta \in \mathcal{R}^s : \|\delta - \beta_I^*\|_\infty \leq b\}$.

If we choose the penalty λ_n such that $(n^{-1} \log d)^{1/2} = o(\lambda_n)$, $sn^{-2/3} = o(\lambda_n)$, $\lambda_n = o(h_n n^{-1/2})$, and $\lambda_n \kappa_{0n} = o(\tau_{0n})$, where

$$\kappa_{0n} = \max_{\delta \in \mathcal{N}_0, (k,j,m) \in I} \frac{1}{18} \left(\sum_{j=1}^{G_k} \|\delta_{kj\cdot}\|_1^{1/2} \right)^{-5/3} \|\delta_{kj\cdot}\|_1^{-1} + \frac{1}{12} \left(\sum_{j=1}^{G_k} \|\delta_{kj\cdot}\|_1^{1/2} \right)^{-2/3} \|\delta_{kj\cdot}\|_1^{-3/2}$$

and $\tau_{0n} = \min_{\delta \in \mathcal{N}_0} \lambda_{\min}(n^{-1} X_I^T \Sigma(X_I \delta) X_I)$, then for sufficiently large n , with probability greater than $1 - 2\{s/n + (Md - s)/d^2\}$, there exists a solution $\hat{\beta}$ to (5), such that

(c) $\hat{\beta}_{II \cup III} = \mathbf{0}$;

(d) $\|\hat{\beta}_I - \beta_I^*\|_\infty \leq cn^{-\gamma}$, where c is a generic positive constant and γ is some number that is arbitrarily smaller than $1/2$.

5. Pathway overlap

So far, we have assumed that pathways do not have overlapping genes. In practice, pathways can share some common genes and this makes the analysis and interpretation much more challenging than the non-overlap case. In general, one can always classify any pathway into the following three cases according to whether its important member genes are shared with some other pathways. For a pathway that has at least one exclusive important gene (a gene that only belongs to that pathway and its effect is nonzero in at least one study), it is important according to definition in (2). For a pathway whose genes all have zero effects, it is unimportant. For a pathway whose important genes are all shared with other pathways, its importance is unclear. The reason is that the shared genes' contribution in different pathways is undetermined even if its total effect is known. Hence, if a pathway's importance cannot be judged by its exclusive genes, its importance is not always identifiable.

To quantify the distribution of the shared gene's effects in different pathways, we use the following assumption. Suppose $X_{k_1 j_1 m} = X_{k_2 j_2 m} = \dots = X_{k_T j_T m}$ for all $1 \leq m \leq M$, i.e., the gene belongs to T pathways. We still adopt the decomposition (1) and assume

(C11) $g_{k_1 j_1} = \dots = g_{k_T j_T} = g$ and $\zeta_{k_1 j_1 m} = \dots = \zeta_{k_T j_T m} = \zeta_m$ for all $1 \leq m \leq M$.

Under condition (C11), this shared gene's effect can be decomposed as $(p_{k_1} + \dots + p_{k_T})g\zeta_m$. Its contribution in pathway k_ℓ is proportional to the pathway effect p_{k_ℓ} . Once condition (C11) is assumed, our algorithm in Section 2 still works for the overlapping case and gives a pathway effect estimator \hat{p} . We select pathways by referring to the nonzero elements of \hat{p} . As for gene selection, when a shared gene is selected, our algorithm tells under which pathways the gene is selected (by referring to nonzero elements of \hat{p}).

In reality, the overlapping scheme of pathways can be complex. The asymptotic properties are difficult to establish under this setting. We leave it as a topic for future research. Nevertheless, we investigate the empirical performance of our method under this setting by simulation studies, where overlapping pathways for the first two definitive cases are simulated.

6. Simulation studies

We simulate scenarios of non-overlapping and overlapping pathways to inspect our method’s finite-sample performance.

6.1. Non-overlapping pathways

We simulate $M = 10$ studies with sample size $n_m = 50$ in each study. The expression of $d = 100$ genes from $K = 20$ pathways are simulated, with 5 genes in each pathway. The gene expression of the i th subject in the m th study \mathbf{x}_{mi} is i.i.d. from $\mathcal{N}_{100}(\mathbf{0}, \mathbf{I})$. The corresponding phenotype y_{mi} is generated from the logistic model

$$\Pr(y_{mi} = 1 | \mathbf{x}_{mi}) = \frac{\exp(\mathbf{x}_{mi}^\top \boldsymbol{\beta}_m^*)}{1 + \exp(\mathbf{x}_{mi}^\top \boldsymbol{\beta}_m^*)},$$

where $\boldsymbol{\beta}_m^* = (\beta_{1m}^*, \beta_{2m}^*, \dots, \beta_{dm}^*)^\top$.

We let $\beta_{kjm}^* = a_k b_{kj} c_{kjm}$, where a_k is i.i.d. from $N(v_k, 0.5^2)$ with $v_1 = 8, v_2 = 8, v_3 = -4, v_4 = -4, v_5 = -8$, and $a_k = 0$ for $k > 5$; b_{kj} is i.i.d. from Bernoulli(π_g); and c_{kjm} is i.i.d. from Bernoulli(π_m). That is, only the first five pathways are important. Each gene under the important pathway has probability π_g to be important. For an important gene, it has probability π_m to be important in the m th study. The probabilities π_g and π_m control the heterogeneity among genes and studies, respectively. We choose two values (0.3 and 0.9) for both π_g and π_m to represent different levels of heterogeneity. Simulations were repeated for 100 runs.

Our method is compared with two ad-hoc selection methods:

Separate Group LASSO (seGLASSO): it runs study-by-study selections by treating genes under each pathway as a group and imposes a Group LASSO penalty therein, i.e., in each study it solves the problem

$$\operatorname{argmin}_{\boldsymbol{\beta}} -\ell_m(\boldsymbol{\beta}_m) + \lambda_m \sum_{k=1}^{20} \left(\sum_{j=1}^5 \beta_{kjm}^2 \right)^{1/2},$$

where λ_m is a positive tuning parameter.

Stack Group LASSO (stGLASSO): it stacks data in all studies and penalizes the total negative log-likelihood with a Group LASSO penalty, i.e., it solves the problem

$$\operatorname{argmin}_{\boldsymbol{\beta}} \sum_{m=1}^{10} -\ell_m(\boldsymbol{\beta}_m) + \lambda \sum_{k=1}^{20} \left(\sum_{j=1}^5 \sum_{m=1}^{10} \beta_{kjm}^2 \right)^{1/2},$$

where λ is a positive tuning parameter.

The optimal tuning parameters in the three methods are all selected by minimizing the Bayesian Information Criterion.

The three competitors’ performance is assessed by their selection capability of pathways and genes. Two measurements will be presented: (a) Sensitivity: the proportion of important pathways/genes being selected; (b) Specificity: the proportion of unimportant pathways/genes not being selected. In particular, the set of important pathways is $\{1, \dots, 5\}$. The set of important genes is defined as $\{(k, j, m) : \beta_{kjm}^* \neq 0\}$.

Fig. 1 gives the boxplots for pathway and gene selection among the 100 simulations. The overall performance is measured by both sensitivity and specificity. Tables 1 and 2 give the means of sensitivity and specificity. These results clearly show that the pathway selection by our method is the best in all settings. The gene selection is also the best when studies are heterogeneous ($\pi_m = 0.3$) and comparable to stGLASSO when studies are homogeneous ($\pi_m = 0.9$). The seGLASSO performs badly as it uses one data at a time and does not borrow strength across studies. The performance of stGLASSO changes drastically in different settings. When π_m is small, it performs badly as it does not take study heterogeneity into account. Even when π_m is large, its performance of pathway selection is still worse than ours. In conclusion, our method performs consistently well for various settings in terms of both pathway selection and gene selection accuracy.

The computational speed of our estimator is comparable to that of seGLASSO and stGLASSO. From the algorithm in Section 3, each iteration of our method solves a convex optimization, similar to seGLASSO and stGLASSO. In the simulations, we noticed that our method usually took very few iterations to convergence.

6.2. Overlapping pathways

We simulate three more examples for pathways with overlapping genes. We simulate $M = 5$ studies with sample size $n_m = 30$ in each study, the number of pathways $K = 21$ and the number of genes $d = 100$. The gene expression values \mathbf{x}_{mi} are generated in the same way as Section 6.1. The true effects of genes are generated by letting $\beta_{kjm}^* = p_k g_{kj}^* s_{kjm}^*$ and the following three cases are considered. Table 3 gives the genes’ effects for the first four pathways. The unlisted genes have zero effects.

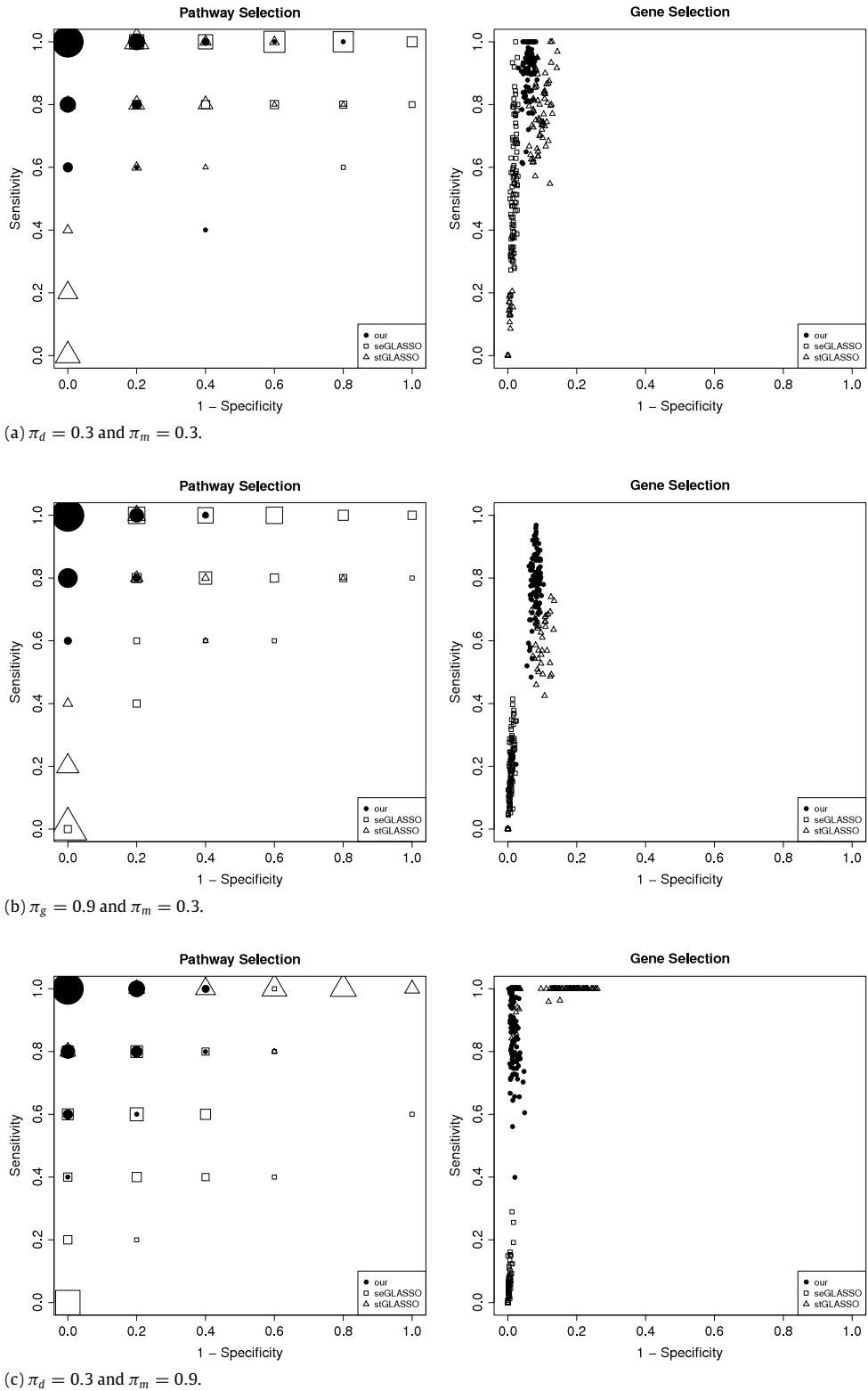
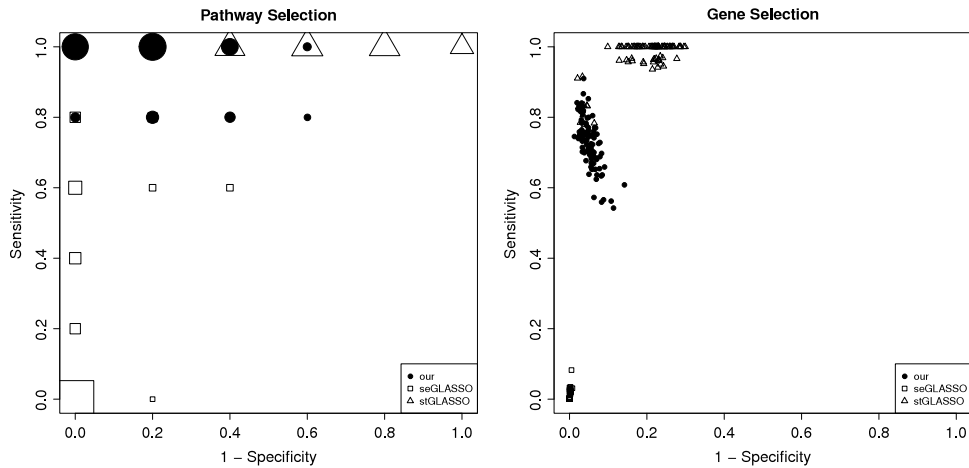


Fig. 1. ROC type of scatter plots for pathway and gene selections (the size of a point indicates its frequency among 100 simulations).



(d) $\pi_d = 0.9$ and $\pi_m = 0.9$.

Fig. 1. (continued)

Table 1
Pathway selection accuracy by the three methods.

		$\pi_m = 0.3$		$\pi_m = 0.9$	
		Sensitivity	Specificity	Sensitivity	Specificity
$\pi_g = 0.3$	Our	0.93	0.93	0.93	0.93
	seGLASSO	0.96	0.42	0.39	0.86
	stGLASSO	0.61	0.86	0.98	0.56
$\pi_g = 0.9$	Our	0.94	0.96	0.85	0.91
	seGLASSO	0.88	0.60	0.20	0.97
	stGLASSO	0.36	0.94	1.00	0.41

Table 2
Gene selection accuracy by the three methods.

		$\pi_m = 0.3$		$\pi_m = 0.9$	
		Sensitivity	Specificity	Sensitivity	Specificity
$\pi_g = 0.3$	Our	0.92	0.94	0.84	0.98
	seGLASSO	0.55	0.98	0.04	0.99
	stGLASSO	0.52	0.93	0.99	0.86
$\pi_g = 0.9$	Our	0.78	0.92	0.60	0.97
	seGLASSO	0.19	0.99	0.01	0.99
	stGLASSO	0.23	0.96	0.97	0.80

Example 1. $p_k^* = 1$ for $1 \leq k \leq 4$ and $p_k^* = 0$ otherwise; $g_1^* = (5, 5, 2.5, 2.5, 0)^\top$, $g_2^* = (2.5, 2.5, 5, 2.5, 2.5)^\top$, $g_3^* = (2.5, 2.5, 5, 5, 5)^\top$, $g_4^* = (-8, -8, -8, -8, -8)^\top$, and $g_k^* = (0, 0, 0, 0, 0)^\top$ otherwise; $\zeta_{kjm}^* = 1$ for $1 \leq k \leq 4$, $1 \leq j \leq 5$, $1 \leq m \leq 5$ and $\zeta_{kjm}^* = 0$ otherwise.

Example 2. $p_k^* = 1$ for $k = 1, 2, 4$ and $p_k^* = 0$ otherwise; $g_1^* = (5, 5, 2.5, 2.5, 0)^\top$, $g_2^* = (2.5, 2.5, 5, 0, 0)^\top$, $g_4^* = (-8, -8, -8, -8, -8)^\top$ and $g_k^* = (0, 0, 0, 0, 0)^\top$ otherwise; $\zeta_{kjm}^* = 1$ for $1 \leq k \leq 4$, $1 \leq j \leq 5$, $1 \leq m \leq 5$ and $\zeta_{kjm}^* = 0$ otherwise.

Example 3. $p_k^* = 1$ for $1 \leq k \leq 4$ and $p_k^* = 0$ otherwise; $g_1^* = (5, 5, 2, 2, 0)^\top$, $g_2^* = (2, 2, 5, 0, 1)^\top$, $g_3^* = (2, 2, 0, -2, 5, 5)^\top$, $g_4^* = (-1, -1, 0, -8, -8, -8)^\top$ and $g_k^* = (0, 0, 0, 0, 0)^\top$ otherwise; $\zeta_{kjm}^* = 1$ for $1 \leq k \leq 4$, $1 \leq j \leq 5$, $1 \leq m \leq 5$ and $\zeta_{kjm}^* = 0$ otherwise.

In **Example 1**, the first four pathways are important and three of them share some common genes. In **Example 2**, only pathways 1, 2, and 4 are important. Even though pathway 3 is unimportant, it shares two genes (gene 6 and 7) with an important pathway (pathway 2). We design such an example to see how well our method could remove an unimportant pathway, which overlaps with an important one. In **Example 3**, genes 3 and 4 are shared by all the four important pathways. **Fig. 2** gives the proportion of each pathway being selected among 100 simulations. **Table 4** gives the means of sensitivity/specificity of pathway and gene selection.

Fig. 2 shows that our method can correctly distinguish important and unimportant pathways. For the challenging case of **Example 2**, our method correctly selects pathway 2 and removes pathway 3, even though they share two common genes.

Table 3

The effects of genes in the first four pathways for the three simulated examples.

Example 1	Genes														
Pathway 1	5	5	2.5	2.5	0	-	-	-	-	-	-	-	-	-	-
Pathway 2	-	-	2.5	2.5	5	2.5	2.5	-	-	-	-	-	-	-	-
Pathway 3	-	-	-	-	-	2.5	2.5	5	5	5	-	-	-	-	-
Pathway 4	-	-	-	-	-	-	-	-	-	-	-8	-8	-8	-8	-8
Overall effect	5	5	5	5	5	5	5	5	5	5	-8	-8	-8	-8	-8
Example 2	Genes														
Pathway 1	5	5	2.5	2.5	0	-	-	-	-	-	-	-	-	-	-
Pathway 2	-	-	2.5	2.5	5	0	0	-	-	-	-	-	-	-	-
Pathway 3	-	-	-	-	-	0	0	0	0	0	-	-	-	-	-
Pathway 4	-	-	-	-	-	-	-	-	-	-	-8	-8	-8	-8	-8
Overall effect	5	5	5	5	5	0	0	0	0	0	-8	-8	-8	-8	-8
Example 3	Genes														
Pathway 1	5	5	2	2	0	-	-	-	-	-	-	-	-	-	-
Pathway 2	-	-	2	2	5	0	-	1	-	-	-	-	-	-	-
Pathway 3	-	-	2	2	-	-	0	-	-2	5	5	-	-	-	-
Pathway 4	-	-	-1	-1	-	-	-	-	-	-	-	0	-8	-8	-8
Overall effect	5	5	5	5	5	0	0	1	-2	5	5	0	-8	-8	-8

Table 4

Performance of pathway/gene selection of our method for overlapping pathways.

	Pathway selection		Gene selection	
	Sensitivity	Specificity	Sensitivity	Specificity
Example 1	0.97	0.86	0.41	0.98
Example 2	0.97	0.85	0.53	0.98
Example 3	0.98	0.85	0.44	0.98

For the other challenging case of [Example 3](#), our method correctly selects the first four pathways, even though two genes are shared by all of them.

7. An integrative analysis of five cardiovascular disease studies

For further illustration, we applied our method to an integrative analysis of five cardiovascular disease (CVD) studies. These studies were aimed to identify biomarkers that are associated with immune response in the development of atherosclerosis. The phenotypes in these studies are binary. The case groups are patients showing certain atherosclerotic syndrome, e.g., having had ischemic strokes. The control groups are healthy people. All subjects were sequenced using microarrays. The raw data can be found on Gene Expression Omnibus with accession names “GSE12288”, “GSE26561”, “GSE20129”, “GSE22255”, and “GSE28829”. [Table 5](#) presents more details of the five studies. However, a careful inspection of their original findings revealed that these studies identified completely different sets of genes. Therefore, the underlying genomic mechanism is still largely unknown. It also indicated that the gene selections from case-by-case studies are hard to be reproduced. This motivated us to incorporate external pathway information and integrate the datasets. We cross-referenced the genes with the pathway information listed on the Kyoto Encyclopedia of Genes and Genomes (KEGG). In total, 4156 genes in 210 pathways from KEGG were involved in the analysis. We applied our method to select pathways as well as their important member genes.

Eleven pathways were identified. [Table 6](#) gives the selected pathways and genes. [Table 7](#) gives the selected genes in each individual study. Among the selected pathways, “Antigen process and presentation”, “Hedgehog signaling”, “Osteoclast differentiation” and “Phagosome” were known to control key modules for activating human body’s immune system [[11,21,26,28](#)]. These pathways can certainly be triggered by inflammation. Epidemiological and clinical studies have shown strong and consistent relationships between markers of inflammation and risk of cardiovascular events [[27](#)]. Moreover, “Ubiquitin mediated proteolysis”, “mRNA surveillance” and “Protein processing in endoplasmic reticulum” play important roles in a broad array of basic cellular processes [[4,9](#)]. “Axon guidance” was previously identified to regulate molecules for the angiogenic growth of blood vessels [[1](#)]. The “peroxisome proliferator-activated receptor” (PPAR) was shown to be an important regulator of cardiac metabolism [[7](#)]. Elevated circulating levels of “cytokines and/or cytokine receptors” was also known to predict adverse outcomes in patients with heart failure [[5](#)]. These key pathways were identified by our method.

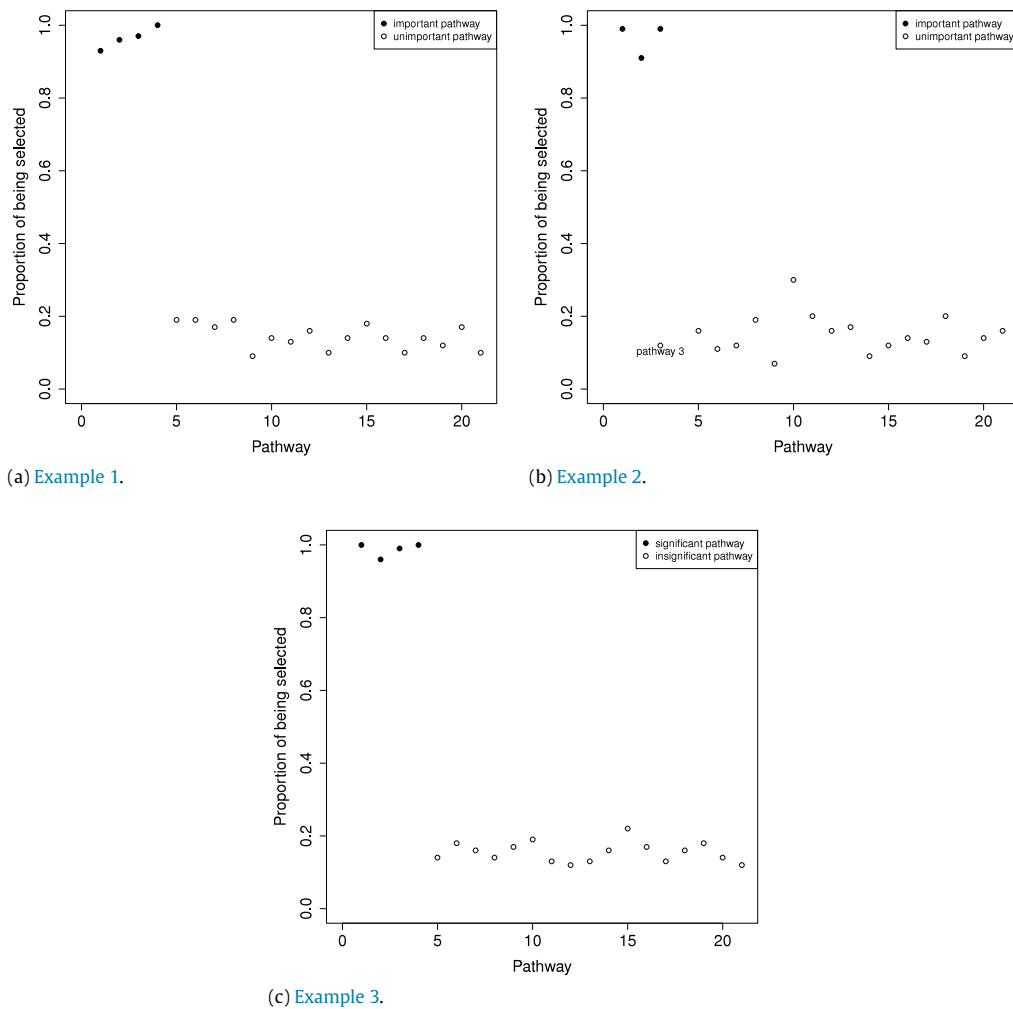


Fig. 2. Proportion of each pathway being selected among 100 simulations.

In addition, Huntington's disease patients were known to have higher risk of CVD, presumably as a result of accumulation of toxic amyloid-like inclusions [17]. Our selection of "Huntington's disease" pathway indicates that Huntington's disease and CVD can share some common prognostic biomarkers.

Moreover, we used the pathway enrichment score proposed in [18] to assess how our selection of pathways is enriched. In particular, for each selected pathway, we created a 2×2 contingency table of the number of genes selected in and out of that pathway. Then, we tested whether that pathway is over-represented by the Fisher's exact test. The corresponding p -values were used as the enrichment scores, where smaller p -values indicate stronger evidence that the selected pathway is enriched. We calculated pathway enrichment scores in each individual study as well as over all five studies together. Table 6 gives the results of the enrichment scores. It clearly indicates that the selected pathways become more enriched when the integrative method is applied over multiple studies than their appearance in a single study.

In conclusion, our integrative method identified more key biomarkers than the original individual studies. The resulting pathway selections become more enriched. Both results necessitate our proposal of integrating data and pave the way for further studies of the interactions among the identified biomarkers (see Table 8).

8. Discussion

In this article, we have provided a general framework of enhancing the integrative analysis with pathway information, which selects pathways as well as their functioning genes. The incorporation of pathway information not only improves the biological interpretation of the results but also brings statistical advantage. We adopt a hierarchical decomposition on gene effects and show that whether a pathway has nonzero effect is determined once a reasonable identifiability condition is assumed. We show that our method can consistently select both pathways and genes. The pathway selection requires

Table 5
Summaries of datasets in the five cardiovascular studies.

	GSE12288	GSE16561	GSE20129	GSE22255	GSE28829
Total sample size	222	63	119	40	29
Case group size	110	39	48	20	16
Control group size	112	24	71	20	13
Microarray Platform	Affymetrix U133A	Illumina HumanRef8 V3.0	Illumina HumanRef8 V2.0	Affymetrix U133plus2	Affymetrix U133plus2

Table 6
Pathway and gene selections by our method and the p -values of pathway enrichment.

Names of selected pathways	Name of selected genes
Antigen processing and presentation	CANX, HLA-A, HLA-DRA
Axon guidance	SEMA6C
Cytokine–cytokine receptor interaction	BMP2, CXCR3, EDAR, IFNA4, IL21R, INHBA, INHBE
Hedgehog signaling	BMP2, BMP6, CSNK1D
Huntington's disease	ATP5B, DNAI2, NDUFA4L2
mRNA surveillance	NUDT21
Osteoclast differentiation	ACP5, GAB2, RELB, LILRB4
Phagosome	CANX, HLA-DRA
PPAR signaling	GK, RXRG
Protein processing in endoplasmic reticulum	CANX, CAPN1, DERL2, HSPH1
Ubiquitin mediated proteolysis	ANAPC13, ANAPC2, PPIL2, TRIM37, UBE3C

Table 7
Gene selections by our method in each dataset.

GSE12288	ACP5, ANAPC2, ATP5B, BMP2, BMP6, CANX, CAPN1, CSNK1D, CXCR3, DERL2, EDAR, HLA-A, HSPH1, LILRB4, NDUFA4L2, RELB, SEMA6C, TRIM37
GSE16561	ANAPC13, ATP5B, CANX, CSNK1D, DNAI2, RXRG
GSE20219	ACP5, ANAPC13, BMP6, CAPN1, CXCR3, EDAR, GAB2, GK, HLA-DRA, IFNA4, IL21R, INHBA, INHBE, NUDT21, UBE3C
GSE22255	ANAPC13, ATP5B, CANX, EDAR, PPIL2, RXRG
GSE28829	ACP5, ANAPC13, ANAPC2, PPIL2

Table 8
Pathway enrichment scores in each dataset and over all five datasets. The scores are not calculated when a pathway is not selected in a particular dataset.

	GSE12288	GSE16561	GSE20219	GSE22255	GSE28829	Overall
Antigen processing and presentation	0.02	–	0.19	–	–	0.008
Axon guidance	0.41	–	–	–	–	0.59
Cytokine–cytokine receptor interaction	0.08	–	0.01	0.46	–	0.001
Hedgehog signaling pathway	0.02	0.12	0.17	–	–	0.006
Huntington's disease	0.02	0.05	–	0.32	–	0.11
mRNA surveillance pathway	–	–	0.22	–	–	0.403
Osteoclast differentiation	0.11	–	0.08	–	0.30	0.017
Phagosome	0.45	–	–	0.29	–	0.08
PPAR signaling pathway	–	0.15	0.22	0.15	–	0.089
Protein processing in endoplasmic reticulum	0.02	0.30	0.42	–	–	0.021
Ubiquitin mediated proteolysis	0.09	0.26	0.07	0.03	0.003	0.001

weaker minimal signal strength condition which allows some false positives and negatives at the gene selection level, thus avoids the restrictive irrepresentable condition. Such advantages have been explicitly quantified in our theoretical results.

Our method can also be adapted to certain new scientific information. For example, for some pathways, their genes tend to be all upregulated or downregulated. Such information can be easily built into our method by requiring $g_{kj} \geq 0$ and $\zeta_{kjm} \geq 0$ for such a pathway k . These constraints can be added to Steps 2 and 3 of our algorithm at minimum extra computational cost. Once they are incorporated in the algorithm, our solution guarantees that the effects of genes under such pathways will be concordant.

The L_1 penalty on the decomposed parameters imposed in (3) is not the only choice. If the studies being integrated are more homogeneous, the L_1 penalty on ζ can be replaced by the L_2 penalty. In this way, our method will select genes in an “all-in-or-all-out” fashion, meaning that a gene is either effective in all studies or not effective at all.

Acknowledgments

The authors would like to thank the editor-in-chief, the associate editor and two reviewers for reviewing the manuscript and provide valuable comments which have led to a substantial improvement from the original version. Li's work is partially

supported by NSF Grant DMS-1127914 and NIH Grants R01GM047845 and R01AI029168. Wang’s work is partially supported by NIH Grant 5 R01 HG007377-02.

Appendix. Proofs

Proof of Proposition 1. First, if $p_k^* = 0$, then $\beta_{kjm}^* = 0$. Hence, it is obvious that $\{k : p_k^* \neq 0\}^c \subset \mathcal{P}^c$. On the other hand, by the definition of \mathcal{P} and the identifiability condition (Condition 1), $\mathcal{P}^c \subset \{k : p_k^* = 0\}$.

Proof of Proposition 2. Since $\hat{\beta}_{kjm} = \hat{p}_k \hat{g}_{kj} \hat{\zeta}_{kjm}$, $\hat{p}_k = 0$ implies $\hat{\beta}_{kjm} = 0$ for all $j = 1, \dots, G_k$ and $m = 1, \dots, M$. Hence, $\{k : \hat{p}_k \neq 0\}^c \subset \hat{\mathcal{P}}^c$. On the other hand, if $\hat{\beta}_{kjm} = 0$ for all j and m , since $(\hat{p}_k, \hat{g}_{kj}, \hat{\zeta}_{kjm})$ minimizes (3), it must hold that $\hat{p}_k = \hat{g}_{kj} = \hat{\zeta}_{kjm} = 0$, otherwise $(0, 0, 0)$ will give a smaller value of the objective function in (3).

Proof of Lemma 1. Referring to the main body of the paper, under decomposition (2), (3) is equivalent to

$$\begin{aligned} \operatorname{argmin}_{\alpha, \mathbf{p}, \mathbf{g}, \zeta} & - \sum_{m=1}^M \ell_m(\alpha_m, \boldsymbol{\beta}_m) + \sum_{k=1}^K |p_k| + \sum_{k=1}^K \sum_{j=1}^{G_k} |g_{kj}| + \chi_n \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}|, \\ \text{s.t. } & \beta_{kjm} = p_k g_{kj} \zeta_{kjm}. \end{aligned}$$

Since the term $-\sum_{m=1}^M \ell_m(\alpha_m, \boldsymbol{\beta}_m)$ is irrelevant to $(\mathbf{p}, \mathbf{g}, \zeta)$, the above problem is further equivalent to

$$\begin{aligned} \min_{\mathbf{p}, \mathbf{g}, \zeta} & \sum_{k=1}^K |p_k| + \sum_{k=1}^K \sum_{j=1}^{G_k} |g_{kj}| + \chi_n \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}|, \\ \text{s.t. } & \beta_{kjm} = p_k g_{kj} \zeta_{kjm}. \end{aligned} \tag{A.1}$$

Next, we show that (A.1) is equivalent to

$$\begin{aligned} \min_{\mathbf{p}, \mathbf{g}, \zeta} & \sum_{k=1}^K |p_k| + \sum_{k=1}^K \sum_{j=1}^{G_k} |g_{kj}| + \chi_n \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}|, \\ \text{s.t. } & |\beta_{kjm}| = |p_k g_{kj} \zeta_{kjm}|. \end{aligned} \tag{A.2}$$

In fact, it holds trivially that any solution of (A.1) is also a solution of (A.2). On the contrary, if $(\hat{p}_k, \hat{g}_{kj}, \hat{\zeta}_{kjm})$ is a solution of (A.2), then

$$(\bar{p}_k, \bar{g}_{kj}, \bar{\zeta}_{kjm}) = \begin{cases} (\hat{p}_k, \hat{g}_{kj}, \hat{\zeta}_{kjm}) & \text{if } \beta_{kjm} \hat{p}_k \hat{g}_{kj} \hat{\zeta}_{kjm} \geq 0 \\ (\hat{p}_k, \hat{g}_{kj}, -\hat{\zeta}_{kjm}) & \text{if } \beta_{kjm} \hat{p}_k \hat{g}_{kj} \hat{\zeta}_{kjm} < 0 \end{cases}$$

is a solution of (A.1). Hence (A.1) and (A.2) are equivalent in the sense that there is a one-to-one correspondence between their solutions. By Lagrange multiplier, (A.2) is equivalent to minimizing

$$\mathcal{H}(\mathbf{p}, \mathbf{g}, \zeta) := \sum_{k=1}^K |p_k| + \sum_{k=1}^K \sum_{j=1}^{G_k} |g_{kj}| + \chi_n \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}| + \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M \alpha_{kjm} (|\beta_{kjm}| - |p_k g_{kj} \zeta_{kjm}|).$$

Then, we have

$$\frac{\partial \mathcal{H}}{\partial |p_k|} = 1 - \sum_{j=1}^{G_k} \sum_{m=1}^M \alpha_{kjm} |g_{kj} \zeta_{kjm}| = 0, \tag{A.3}$$

$$\frac{\partial \mathcal{H}}{\partial |g_{kj}|} = 1 - \sum_{m=1}^M \alpha_{kjm} |p_k \zeta_{kjm}| = 0, \tag{A.4}$$

$$\frac{\partial \mathcal{H}}{\partial |\zeta_{kjm}|} = \chi_n - \alpha_{kjm} |p_k g_{kj}| = 0, \tag{A.5}$$

(A.5) implies that $\alpha_{kjm} = \chi_n / |p_k g_{kj}|$, which together with (A.4) give $|g_{kj}| = \chi_n \sum_{m=1}^M |\zeta_{kjm}|$. Then, it follows from (A.3) that

$$|p_k| = \sum_{j=1}^{G_k} |g_{kj}| = \chi_n \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}|. \tag{A.6}$$

Since $|\beta_{kjm}| = |p_k g_{kj} \zeta_{kjm}|$, we have $\chi_n \sum_{m=1}^M |\beta_{kjm}| = |p_k| |g_{kj}|^2$. Hence,

$$|g_{kj}| = \left(|p_k|^{-1} \chi_n \sum_{m=1}^M |\beta_{kjm}| \right)^{1/2}.$$

This together with $|p_k| = \sum_{j=1}^{G_k} |g_{kj}|$ gives that $|p_k| = \chi_n^{1/6} \{ \sum_{j=1}^{G_k} (\sum_{m=1}^M |\beta_{kjm}|)^{1/2} \}^{1/3}$. Therefore,

$$\sum_{j=1}^{G_k} |g_{kj}| = \chi_n \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}| = \chi_n^{1/6} \left\{ \sum_{j=1}^{G_k} \left(\sum_{m=1}^M |\beta_{kjm}| \right)^{1/2} \right\}^{1/3}.$$

This together with (A.6) shows that

$$\sum_{k=1}^K |p_k| + \sum_{k=1}^K \sum_{j=1}^{G_k} |g_{kj}| + \chi_n \sum_{k=1}^K \sum_{j=1}^{G_k} \sum_{m=1}^M |\zeta_{kjm}| = 3 \chi_n^{1/6} \sum_{k=1}^K \left\{ \sum_{j=1}^{G_k} \left(\sum_{m=1}^M |\beta_{kjm}| \right)^{1/2} \right\}^{1/3}.$$

Hence, (3) and (5) are equivalent.

Proof of Theorem 1. We show that there exists a solution $\hat{\beta}$ to (5) such that its restriction on set J , i.e., $\hat{\beta}_J$, satisfies (b), and all other elements are 0. Then, by definition, this solution also satisfies (a).

By optimization theory, the vector $\hat{\beta}$ that satisfies the following Karush–Kuhn–Tucker (KKT) conditions is a solution to (6).

$$\mathbf{X}_J^\top \mathbf{Y} - \mathbf{X}_J^\top \boldsymbol{\mu}(\mathbf{X}\hat{\beta}) = n\lambda_n \nabla \rho(\hat{\beta}_J), \tag{A.7}$$

$$\mathbf{X}_{J^c}^\top \mathbf{Y} - \mathbf{X}_{J^c}^\top \boldsymbol{\mu}(\mathbf{X}\hat{\beta}) \in n\lambda_n \partial \rho(\hat{\beta}_{J^c}), \tag{A.8}$$

$$\lambda_{\min}(\mathbf{X}_J^\top \boldsymbol{\Sigma}(\mathbf{X}\hat{\beta})\mathbf{X}_J) > n\lambda_n \kappa(\hat{\beta}_J). \tag{A.9}$$

In (A.7), $\nabla \rho(\hat{\beta}_{kjm}) = (1/6) \|\hat{\beta}_{kj}\|_1^{-5/6} \text{sign}(\hat{\beta}_{kjm})$ for $(k, j, m) \in J$. In (A.8), the subgradient $\partial \hat{\beta}_{kjm} = (-\|\hat{\beta}_{kj}\|_1^{-5/6}/6, \|\hat{\beta}_{kj}\|_1^{-5/6}/6)$ for $(k, j, m) \in J^c$. In (A.9),

$$\kappa(\hat{\beta}_J) = \max_{(k,j,m) \in J} (5/36) \|\hat{\beta}_{kj}\|_1^{-11/6}.$$

Let $\boldsymbol{\xi} = \mathbf{X}^\top \mathbf{Y} - \mathbf{X}^\top \boldsymbol{\mu}(\mathbf{X}\beta^*)$. Consider the event $E = \{\|\boldsymbol{\xi}_J\|_\infty \leq (2^{-1}n \log n)^{1/2}\}$. Since $y_{mi} \in \{0, 1\}$ and columns of \mathbf{X} have been standardized to have norm $n^{1/2}$, by Proposition 4(a) of [6],

$$\Pr(|\xi_{kjm}| \geq t) \leq 2 \exp(-2t^2/n), \tag{A.10}$$

where ξ_{kjm} is the (k, j, m) th element of $\boldsymbol{\xi}$. Then, by the union bound,

$$\Pr(E) \geq 1 - \sum_{(k,j,m) \in J} P\{\xi_{kjm} \geq (2^{-1}n \log n)^{1/2}\} \geq 1 - 2K^*/n. \tag{A.11}$$

In the event E , we show two results: [1] within the hypercube $\mathcal{M} := \{\beta : \|\beta - \beta_J\|_\infty \leq cn^{-\gamma}\}$, there exists $\hat{\beta}_J \in \mathcal{R}^{K^*}$ that satisfies (A.7) and (A.9); [2] $\hat{\beta} = (\hat{\beta}_J, \mathbf{0})^\top$ satisfies (A.8). These two results together with the KKT conditions and (A.11) complete the proof.

[1] Let $\boldsymbol{\eta} = n\lambda_n \nabla \rho(\beta)$, where the element $\eta_{kjm} = (1/6)n\lambda_n \|\beta_{kj}\|_1^{-5/6} \text{sign}(\beta_{kjm})$. Under Condition 2, for sufficiently large n , we have

$$\|\beta_{kj}^*\|_1 \geq 2B > 2cn^{-\gamma} \geq 2\|\beta_{kj}\|_1 - \|\beta_{kj}^*\|_1, \tag{A.12}$$

for all $\beta \in \mathcal{M}$. In the second inequality, we use the fact that γ is arbitrarily close to $1/2$ so that $B > cn^{-\gamma}$. Therefore, $\|\beta_{kj}\|_1 \geq \|\beta_{kj}^*\|_1/2$. It further implies that $\|\boldsymbol{\eta}\|_\infty \leq B^{-5/6}n\lambda_n$. Define

$$\Psi(\beta) = \mathbf{X}_J^\top \{\boldsymbol{\mu}(\mathbf{X}_J\beta) - \boldsymbol{\mu}(\mathbf{X}_J\beta_J^*)\} - (\boldsymbol{\xi}_J - \boldsymbol{\eta}). \tag{A.13}$$

We show that $\Psi(\beta) = 0$ has a solution $\hat{\beta}_J$ within \mathcal{M} . Then, $\hat{\beta}_J$ also solves (A.7). By Taylor expansion,

$$\mathbf{X}_J^\top \{\boldsymbol{\mu}(\mathbf{X}_J\beta) - \boldsymbol{\mu}(\mathbf{X}_J\beta_J^*)\} = \mathbf{X}_J^\top \boldsymbol{\Sigma}(\mathbf{X}\beta^*)\mathbf{X}_J(\beta - \beta_J^*) + \mathbf{r},$$

where the (k, j, m) th element of \mathbf{r} has $r_{kjm} = \frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\beta}_j^*)^\top \mathbf{R}(\boldsymbol{\delta}^{(kjm)}) (\boldsymbol{\beta} - \boldsymbol{\beta}_j^*)$, $\mathbf{R}(\boldsymbol{\delta}^{(kjm)}) = \mathbf{X}_j^\top \{\text{diag}(\mathbf{X}_{kjm} \circ \boldsymbol{\mu}''(\mathbf{X}_j \boldsymbol{\delta}^{(kjm)}))\} \mathbf{X}_j$, $\boldsymbol{\delta}^{(kjm)}$ is a vector lying on the line segment joining $\boldsymbol{\beta}$ and $\boldsymbol{\beta}_j^*$. By Conditions 2, 5, $\|\mathbf{r}\|_\infty = O(K^* n^{1-2\gamma}) = o(n^{1-\gamma})$, where in the last equality we use the fact that $K^* = o(n^\gamma)$ as γ is arbitrarily close to $1/2$.

Let

$$\bar{\Psi}(\boldsymbol{\beta}) = \{\mathbf{X}_j^\top \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}^*) \mathbf{X}_j\}^{-1} \Psi(\boldsymbol{\beta}) = \boldsymbol{\beta} - \boldsymbol{\beta}_j^* + \mathbf{v},$$

where $\mathbf{v} = -\{\mathbf{X}_j^\top \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}^*) \mathbf{X}_j\}^{-1} (\boldsymbol{\xi}_j - \boldsymbol{\eta} - \mathbf{r})$. Then, we have

$$\begin{aligned} \|\mathbf{v}\|_\infty &\leq \|\{\mathbf{X}_j^\top \boldsymbol{\Sigma}(\mathbf{X}\boldsymbol{\beta}^*) \mathbf{X}_j\}^{-1}\|_\infty (\|\boldsymbol{\xi}_j - \boldsymbol{\eta}\|_\infty + \|\mathbf{r}\|_\infty) \\ &= O(n^{-1} (\|\boldsymbol{\xi}_j\|_\infty + \|\boldsymbol{\eta}\|_\infty + \|\mathbf{r}\|_\infty)) \\ &= O((n^{-1} \log n)^{1/2}) + o(n^{-1/2}) + o(n^{-\gamma}) = o(n^{-\gamma}). \end{aligned} \tag{A.14}$$

Hence, for sufficiently large n , if $(\boldsymbol{\beta} - \boldsymbol{\beta}_j^*)_{kjm} = n^{-\gamma}$, we have $\{\bar{\Psi}(\boldsymbol{\beta})\}_{kjm} \geq n^{-\gamma} - \|\mathbf{v}\|_\infty \geq 0$, and if $(\boldsymbol{\beta} - \boldsymbol{\beta}_j^*)_{kjm} = -n^{-\gamma}$, we have $\{\bar{\Psi}(\boldsymbol{\beta})\}_{kjm} \leq -n^{-\gamma} + \|\mathbf{v}\|_\infty \leq 0$. Since the function $\bar{\Psi}(\boldsymbol{\beta})$ is continuous in \mathcal{M} , an application of Miranda's existence theorem (see, e.g., [25]) shows that equation $\bar{\Psi}(\boldsymbol{\beta}) = \mathbf{0}$ has a solution $\hat{\boldsymbol{\beta}}_j$ in \mathcal{M} . Then, $\hat{\boldsymbol{\beta}}_j$ also solves $\Psi(\boldsymbol{\beta}) = 0$ and further solves (A.7).

Since $B > cn^{-\gamma} \geq \|\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}_j^*\|_\infty$, we have $\|\hat{\boldsymbol{\beta}}_{kj}\|_1 \neq \mathbf{0}$ for $(k, j, \cdot) \in J$. Therefore, for any $k \in \mathcal{P}$, there exists at least one (j, m) such that $\hat{\boldsymbol{\beta}}_{kjm} \neq \mathbf{0}$. By definition, $\mathcal{P} \subset \hat{\mathcal{P}}$.

[2] Next, we show that $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}_j, \mathbf{0})^\top$ satisfies (A.8), therefore solves (5). Thus, the solution $\hat{\boldsymbol{\beta}}$ admits the property that $\hat{\mathcal{P}} = \mathcal{P}$. To show that $\hat{\boldsymbol{\beta}}$ satisfies (A.8), recall that the subgradient $\partial \hat{\boldsymbol{\beta}}_{kjm}$ is the interval $(-\|\hat{\boldsymbol{\beta}}_{kj}\|_1^{-5/6}/6, \|\hat{\boldsymbol{\beta}}_{kj}\|_1^{-5/6}/6)$. When $\hat{\boldsymbol{\beta}}_{jc} = \mathbf{0}$, this interval is $(-\infty, +\infty)$. Hence, (A.8) always holds when $\hat{\boldsymbol{\beta}}_{jc} = \mathbf{0}$.

Proof of Theorem 2. Again, by the KKT conditions, any vector $\hat{\boldsymbol{\beta}}$ satisfies (A.15)–(A.17) is a solution to (5).

$$\mathbf{X}_I^\top \mathbf{Y} - \mathbf{X}_I^\top \boldsymbol{\mu}(\mathbf{X}\hat{\boldsymbol{\beta}}) = n\lambda_n \nabla \rho(\hat{\boldsymbol{\beta}}_I), \tag{A.15}$$

$$\mathbf{X}_{II}^\top \mathbf{Y} - \mathbf{X}_{II}^\top \boldsymbol{\mu}(\mathbf{X}\hat{\boldsymbol{\beta}}) = n\lambda_n \partial \rho(\hat{\boldsymbol{\beta}}_{II}), \tag{A.16}$$

$$\lambda_{\min}(\mathbf{X}_I^\top \boldsymbol{\Sigma}(\mathbf{X}\hat{\boldsymbol{\beta}}) \mathbf{X}_I) > n\lambda_n \kappa(\hat{\boldsymbol{\beta}}_I), \tag{A.17}$$

where $\nabla \rho(\hat{\boldsymbol{\beta}}_{kjm}) = \tau_{kj}(\hat{\boldsymbol{\beta}}_{kj}) \text{sign}(\hat{\boldsymbol{\beta}}_{kjm})/6$ for $(k, j, m) \in I$, $\partial \rho(\hat{\boldsymbol{\beta}}_{kjm}) \in (-\tau_{kj}(\hat{\boldsymbol{\beta}}_{kj})/6, \tau_{kj}(\hat{\boldsymbol{\beta}}_{kj})/6)$, for $(k, j, m) \in II$,

$$\tau_{kj}(\hat{\boldsymbol{\beta}}_{kj}) = \left(\sum_{j=1}^{G_k} \|\hat{\boldsymbol{\beta}}_{kj}\|_1^{\frac{1}{2}} \right)^{-\frac{2}{3}} \|\hat{\boldsymbol{\beta}}_{kj}\|_1^{-\frac{1}{2}}, \tag{A.18}$$

$$\kappa(\hat{\boldsymbol{\beta}}_I) = \max_{(k,j,m) \in I} \frac{1}{18} \left(\sum_{j=1}^{G_k} \|\hat{\boldsymbol{\beta}}_{kj}\|_1^{\frac{1}{2}} \right)^{-\frac{5}{3}} \|\hat{\boldsymbol{\beta}}_{kj}\|_1^{-1} + \frac{1}{12} \left(\sum_{j=1}^{G_k} \|\hat{\boldsymbol{\beta}}_{kj}\|_1^{\frac{1}{2}} \right)^{-\frac{2}{3}} \|\hat{\boldsymbol{\beta}}_{kj}\|_1^{-\frac{3}{2}}. \tag{A.19}$$

Let $\boldsymbol{\xi} = \mathbf{X}^\top \mathbf{Y} - \mathbf{X}^\top \boldsymbol{\mu}(\mathbf{X}\boldsymbol{\beta}^*)$. Consider events

$$E_1 = \{\|\boldsymbol{\xi}_I\|_\infty \leq (2^{-1}n \log n)^{1/2}\} \quad \text{and} \quad E_2 = \{\|\boldsymbol{\xi}_{I^c}\|_\infty \leq (n \log p)^{1/2}\}.$$

By (A.10) and the union bound, we have

$$\begin{aligned} \Pr(E_1 \cap E_2) &\geq 1 - \sum_{(k,j,m) \in I} P\{|\xi_{kjm}| \geq (2^{-1}n \log n)^{1/2}\} - \sum_{(k,j,m) \in I^c} P\{|\xi_{kjm}| \geq (n \log p)^{1/2}\} \\ &\geq 1 - 2\{s/n + (Md - s)/d^2\}. \end{aligned}$$

In event $E_1 \cap E_2$, we show two results: [1] within the set $\mathcal{M}_0 := \{\boldsymbol{\beta} \in \mathcal{R}^S : \|\boldsymbol{\beta} - \boldsymbol{\beta}_I^*\|_\infty \leq cn^{-\gamma}\}$, there exists a vector $\hat{\boldsymbol{\beta}}_I \in \mathcal{R}^S$ satisfying (A.15) and (A.17); [2] $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}_I, \mathbf{0})^\top$ satisfies (A.16). These two steps together with the KKT conditions complete the proof.

[1] Let $\boldsymbol{\eta} = n\lambda_n \nabla \rho(\boldsymbol{\beta})$, where $\eta_{kjm} = n\lambda_n \tau_{kj}(\boldsymbol{\beta}_{kj}) \text{sign}(\boldsymbol{\beta}_{kjm})/6$. Under Condition 7, Similarly as in (A.12) we have $\|\boldsymbol{\beta}_{kj}^*\|_1 > 2\|\boldsymbol{\beta}_{kj}^* - \boldsymbol{\beta}_{kj}\|_1$ and further $\|\boldsymbol{\beta}_{kj}\|_1 \geq \|\boldsymbol{\beta}_{kj}^*\|_1/2$. Therefore, $6|\eta_{kjm}| = n\lambda_n \tau_{kj}(\boldsymbol{\beta}_{kj}) \leq n\lambda_n \tau_{kj}(\boldsymbol{\beta}_{kj}^*)$. Hence, $\|\boldsymbol{\eta}\|_\infty \leq n\lambda_n h_n^{-1}$. It further implies that $\|\boldsymbol{\xi}_I - \boldsymbol{\eta}\|_\infty \leq \|\boldsymbol{\xi}_I\|_\infty + \|\boldsymbol{\eta}\|_\infty \leq (2^{-1}n \log n)^{1/2} + n\lambda_n h_n^{-1}$.

Define

$$\Psi(\boldsymbol{\beta}) = \mathbf{X}_I^\top \{\boldsymbol{\mu}(\mathbf{X}_I \boldsymbol{\beta}) - \boldsymbol{\mu}(\mathbf{X}_I \boldsymbol{\beta}_I^*)\} - (\boldsymbol{\xi}_I - \boldsymbol{\eta}). \tag{A.20}$$

Again, we show that $\Psi(\beta) = \mathbf{0}$ has a solution within \mathcal{M}_0 . By a Taylor expansion, we have,

$$\mathbf{X}_l^\top \{\mu(\mathbf{X}_l \beta) - \mu(\mathbf{X}_l \beta_l^*)\} = \mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l (\beta - \beta_l^*) + \mathbf{r},$$

where the Lagrange remainder $\mathbf{r} = (r_{kjm})^\top$ such that $r_{kjm} = (\beta - \beta_l^*)^\top \mathbf{R}(\tilde{\delta}^{(kjm)})(\beta - \beta_l^*)/2$, $\mathbf{R}(\tilde{\delta}^{(kjm)}) = \mathbf{X}_l^\top \{\text{diag}(\mathbf{X}_{kjm} \circ \mu''(\mathbf{X}_l \tilde{\delta}^{(kjm)}))\} \mathbf{X}_l$ and $\tilde{\delta}^{(kjm)}$ is a vector lying on the line segment joining β and β_l^* . By Condition 10, we have

$$\|\mathbf{r}\|_\infty = O(sn^{1-2\gamma}). \tag{A.21}$$

Let

$$\tilde{\Psi}(\beta) = \{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1} \Psi(\beta) = \beta - \beta_l^* + \mathbf{v},$$

where $\mathbf{v} = -\{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1} (\xi_l - \eta - \mathbf{r})$. Then, it follows from Conditions 7, 8, and $\lambda_n = o(h_n n^{-\gamma})$ that

$$\begin{aligned} \|\mathbf{v}\|_\infty &\leq \|\{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1}\|_\infty (\|\xi_l - \eta\|_\infty + \|\mathbf{r}\|_\infty) \\ &= O((n^{-1} \log n)^{1/2} + \lambda_n h_n^{-1} + sn^{-2\gamma}) = o(n^{-\gamma}). \end{aligned} \tag{A.22}$$

By the same argument as in the proof of Theorem 1, there exists a vector $\hat{\beta}_l$ within \mathcal{M}_0 such that $\tilde{\Psi}(\hat{\beta}_l) = \mathbf{0}$. Hence, $\hat{\beta}_l$ also solves (A.15). On the other hand, by the stated choice of λ_n , $\hat{\beta}_l$ satisfies (A.17) for sufficiently large n .

[2] Let $\hat{\beta} = (\hat{\beta}_l, \mathbf{0})^\top$. Next, we prove that $\hat{\beta}$ satisfies (A.16) for the stated choice of λ_n . Indeed, (A.16) requires that

$$|\mathbf{X}_{kjm}^\top \mathbf{Y} - \mathbf{X}_{kjm}^\top \mu(\mathbf{X} \hat{\beta})| < \frac{1}{6} n \lambda_n \tau_{kj}(\hat{\beta}_{kj}), \tag{A.23}$$

for any $(k, j, m) \in \mathcal{H}$. As proved in [1], $\|\beta_{kj}^*\|_1 > 2\|\beta_{kj}^* - \hat{\beta}_{kj}\|_1$, hence $\|\hat{\beta}_{kj}\|_1 < 1.5\|\beta_{kj}^*\|_1$. Then, we have $\tau_{kj}(\hat{\beta}_{kj}) \geq 0.6\tau_{kj}(\beta_{kj}^*)$. It implies that

$$\min_{k,j} \tau_{kj}(\hat{\beta}_{kj}) \geq \min_{k,j} 0.6\tau_{kj}(\beta_{kj}^*) \geq 0.6u_n^{-1}.$$

To prove (A.16), by (A.23) and the above arguments, it suffices to show that

$$\|\mathbf{X}_l^\top \mathbf{Y} - \mathbf{X}_l^\top \mu(\mathbf{X} \hat{\beta})\|_\infty < 0.1n\lambda_n u_n^{-1}. \tag{A.24}$$

Observe that

$$\mathbf{X}_l^\top \mathbf{Y} - \mathbf{X}_l^\top \mu(\mathbf{X} \hat{\beta}) = \mathbf{X}_l^\top \{\mathbf{Y} - \mu(\mathbf{X} \beta^*)\} + \mathbf{X}_l^\top \{\mu(\mathbf{X} \beta^*) - \mu(\mathbf{X} \hat{\beta})\}. \tag{A.25}$$

In event E_2 , $\|\mathbf{X}_l^\top \{\mathbf{Y} - \mu(\mathbf{X} \beta^*)\}\|_\infty = O((n \log d)^{1/2})$. Then, it follows from $(n^{-1} \log d)^{1/2} = o(\lambda_n)$ that

$$(n\lambda_n)^{-1} \|\mathbf{X}_l^\top \{\mathbf{Y} - \mu(\mathbf{X} \beta^*)\}\|_\infty = o(1). \tag{A.26}$$

For the second term on the right hand side of (A.25),

$$\mathbf{X}_l^\top \{\mu(\mathbf{X} \hat{\beta}) - \mu(\mathbf{X} \beta^*)\} = \mathbf{X}_l^\top \{\mu(\mathbf{X}_l \hat{\beta}_l) - \mu(\mathbf{X}_l \beta_l^*)\} = \mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l (\hat{\beta}_l - \beta_l^*) + \mathbf{w},$$

where $\mathbf{w} = (w_{kjm})^\top$ such that $w_{kjm} = (\hat{\beta}_l - \beta_l^*)^\top \mathbf{R}(\tilde{\delta}^{(kjm)})(\hat{\beta}_l - \beta_l^*)/2$, where $\tilde{\delta}^{(kjm)}$ is a vector lying on the line segment joining $\hat{\beta}_l$ and β_l^* . Similarly as in (A.21),

$$\|\mathbf{w}\|_\infty = O(sn^{1-2\gamma}). \tag{A.27}$$

Since $\hat{\beta}_l$ solves $\tilde{\Psi}(\delta) = \mathbf{0}$, we have $\hat{\beta}_l - \beta_l^* = \{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1} (\xi_l - \eta - \mathbf{r})$. Therefore,

$$\begin{aligned} &(n\lambda_n)^{-1} \mathbf{X}_l^\top \{\mu(\mathbf{X}_l \hat{\beta}_l) - \mu(\mathbf{X}_l \beta_l^*)\} \\ &\leq (n\lambda_n)^{-1} \|\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l \{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1}\|_\infty \cdot (\|\xi_l - \eta\|_\infty + \|\mathbf{r}\|_\infty) + (n\lambda_n)^{-1} \|\mathbf{w}\|_\infty \\ &\leq (n\lambda_n)^{-1} \|\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l \{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1}\|_\infty \|\eta\|_\infty + (n\lambda_n)^{-1} O(\|\xi_l\|_\infty + \|\mathbf{r}\|_\infty + \|\mathbf{w}\|_\infty), \end{aligned}$$

because by Condition 9, $\|\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l \{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1}\|_\infty < 1$. It follows from (A.21) and $sn^{-2\gamma} = o(\lambda_n)$ that $(n\lambda_n)^{-1} O(\|\xi_l\|_\infty + \|\mathbf{r}\|_\infty) = o(1)$. Meanwhile, by (A.27) and the choice of λ_n , $(n\lambda_n)^{-1} \|\mathbf{w}\|_\infty = o(1)$. Using Condition 9 and $\|\eta\|_\infty \leq n\lambda_n h_n^{-1}$, we have

$$(n\lambda_n)^{-1} \|\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l \{\mathbf{X}_l^\top \Sigma(\mathbf{X} \beta^*) \mathbf{X}_l\}^{-1}\|_\infty \cdot \|\eta\|_\infty < 0.1u_n^{-1}.$$

Therefore, (A.24) holds.

References

- [1] R.H. Adams, A. Eichmann, Axon guidance molecules in vascular patterning, *Cold Spring Harbor. Perspect. Biol.* 2 (5) (2010) a001875.
- [2] M. Ashburner, C. Ball, J. Blake, D. Botstein, H. Butler, J. Cherry, A. Davis, K. Dolinski, S. Dwight, J. Eppig, et al., Gene Ontology: tool for the unification of biology, *Nature Genet.* 25 (1) (2000) 25–29.
- [3] P. Carbonetto, M. Stephens, Integrated enrichment analysis of variants and pathways in genome-wide association studies indicates central role for IL-2 signaling genes in type 1 diabetes, and cytokine signaling genes in Crohn's disease, *PLoS Genet.* 9 (10) (2013) e1003770.
- [4] L. Carrier, S. Schlossarek, M.S. Willis, T. Eschenhagen, The ubiquitin-proteasome system and nonsense-mediated mRNA decay in hypertrophic cardiomyopathy, *Cardiovasc. Res.* 85 (2) (2009) 330–338.
- [5] A. Deswal, N.J. Petersen, A.M. Feldman, J.B. Young, B.G. White, D.L. Mann, Cytokines and cytokine receptors in advanced heart failure an analysis of the cytokine database from the Vesnarinone Trial (VEST), *Circulation* 103 (16) (2001) 2055–2059.
- [6] J. Fan, J. Lv, Nonconcave penalized likelihood with NP-dimensionality, *IEEE Trans. Inform. Theory* 57 (8) (2011) 5467–5484.
- [7] F. Finck, The PPAR regulatory system in cardiac physiology and disease, *Cardiovasc. Res.* 73 (2) (2007) 269–277.
- [8] J. Friedman, T. Hastie, R. Tibshirani, Regularization paths for generalized linear models via coordinate descent, *J. Stat. Softw.* 33 (1) (2010) 1–22.
- [9] J. Groenendyk, P.K. Sreenivasaiiah, D.H. Kim, L.B. Agellon, M. Michalak, Biology of endoplasmic reticulum stress in the heart, *Circ. Res.* 107 (10) (2010) 1185–1197.
- [10] D.W. Huang, B.T. Sherman, R.A. Lempicki, Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists, *Nucleic Acids Res.* 1 (2009) 1–13.
- [11] J.C. Kagan, A. Iwasaki, Phagosome as the organelle linking innate and adaptive immunity, *Traffic* 13 (8) (2012) 1053–1061.
- [12] P. Khatri, S. Drăghici, Ontological analysis of gene expression data: current tools, limitations, and open problems, *Bioinformatics* 21 (18) (2005) 3587–3595.
- [13] P. Khatri, M. Sirota, A.J. Butte, Ten years of pathway analysis: current approaches and outstanding challenges, *PLoS Comput. Biol.* 8 (2) (2012) e1002375.
- [14] Q. Li, S. Wang, C.-C. Huang, M. Yu, J. Shao, Meta-analysis based variable selection for gene expression data, *Biometrics* 70 (4) (2014) 872–880.
- [15] J. Liu, J. Huang, S. Ma, Integrative analysis of multiple cancer genomic datasets under the heterogeneity model, *Stat. Med.* 32 (20) (2013) 3509–3521.
- [16] J. Liu, J. Huang, Y. Zhang, Q. Lan, N. Rothman, T. Zheng, S. Ma, Integrative analysis of prognosis data on multiple cancer subtypes, *Biometrics* 70 (3) (2014) 480–488.
- [17] G.C. Melkani, A.S. Trujillo, R. Ramos, R. Bodmer, S.I. Bernstein, K. Ocorr, Huntington's disease induced cardiac amyloidosis is reversed by modulating protein folding and oxidative stress pathways in the drosophila heart, *PLoS Genet.* 9 (12) (2013) e1004024-17.
- [18] M.A. Newton, F.A. Quintana, J.A. Den Boon, S. Sengupta, P. Ahlquist, Random-set methods identify distinct aspects of the enrichment signal in gene-set analysis, *Ann. Appl. Stat.* (2007) 85–106.
- [19] E. Ntzani, J. Ioannidis, et al., Predictive ability of DNA microarrays for cancer outcomes and correlates: an empirical assessment, *Lancet* 362 (9394) (2003) 1439–1444.
- [20] H. Ogata, S. Goto, K. Sato, W. Fujibuchi, H. Bono, M. Kanehisa, KEGG: Kyoto encyclopedia of genes and genomes, *Nucleic Acids Res.* 28 (1) (2000) 27–30.
- [21] F. Rochais, K. Mesbah, R.G. Kelly, Signaling pathways controlling second heart field development, *Circ. Res.* 104 (8) (2009) 933–942.
- [22] R. Saxena, B.F. Voight, V. Lyssenko, N.P. Burt, P.I. de Bakker, H. Chen, J.J. Roix, S. Kathiresan, J.N. Hirschhorn, M.J. Daly, et al., Genome-wide association analysis identifies loci for type 2 diabetes and triglyceride levels, *Science* 316 (5829) (2007) 1331–1336.
- [23] A. Subramanian, P. Tamayo, V.K. Mootha, S. Mukherjee, B.L. Ebert, M.A. Gillette, A. Paulovich, S.L. Pomeroy, T.R. Golub, E.S. Lander, et al., Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles, *PNAS* 102 (43) (2005) 15545–15550.
- [24] G.C. Tseng, D. Ghosh, E. Feingold, Comprehensive literature review and statistical considerations for microarray meta-analysis, *Nucleic Acids Res.* 40 (9) (2012) 3785–3799.
- [25] M.N. Vrahatis, A short proof and a generalization of Mirandas existence theorem, *Proc. Amer. Math. Soc.* 107 (3) (1989) 701–703.
- [26] J.M. Vyas, A.G. Van der Veen, H.L. Ploegh, The known unknowns of antigen processing and presentation, *Nat. Rev. Immunol.* 8 (8) (2008) 607–618.
- [27] J.T. Willerson, P.M. Ridker, Inflammation as a cardiovascular risk factor, *Circulation* 109 (21 suppl 1) (2004) II-2–II-10.
- [28] Y. Wu, M.B. Humphrey, M.C. Nakamura, Osteoclasts—the innate immune cells of the bone, *Autoimmunity* 41 (3) (2008) 183–194.
- [29] T. Wu, K. Lange, Coordinate descent algorithms for lasso penalized regression, *Ann. Appl. Stat.* 2 (1) (2008) 224–244.
- [30] E. Zeggini, M.N. Weedon, C.M. Lindgren, T.M. Frayling, K.S. Elliott, H. Lango, N.J. Timpson, J.R. Perry, N.W. Rayner, R.M. Freathy, et al., Replication of genome-wide association signals in UK samples reveals risk loci for type 2 diabetes, *Science* 316 (5829) (2007) 1336–1341.
- [31] P. Zhao, B. Yu, On model selection consistency of lasso, *J. Mach. Learn. Res.* 7 (2006) 2541–2563.
- [32] N. Zhou, J. Zhu, Group variable selection via a hierarchical lasso and its oracle property, *Stat. Interface* 3 (2010) 557–574.