

Variable Selection with Multiply-Imputed Datasets: Choosing Between Stacked and Grouped Methods

Jiacong Du, Jonathan Boss, Peisong Han, Lauren J. Beesley, Michael Kleinsasser, Stephen A. Goutman, Stuart Batterman, Eva L. Feldman & Bhramar Mukherjee

To cite this article: Jiacong Du, Jonathan Boss, Peisong Han, Lauren J. Beesley, Michael Kleinsasser, Stephen A. Goutman, Stuart Batterman, Eva L. Feldman & Bhramar Mukherjee (2022) Variable Selection with Multiply-Imputed Datasets: Choosing Between Stacked and Grouped Methods, Journal of Computational and Graphical Statistics, 31:4, 1063-1075, DOI: [10.1080/10618600.2022.2035739](https://doi.org/10.1080/10618600.2022.2035739)

To link to this article: <https://doi.org/10.1080/10618600.2022.2035739>



View supplementary material [↗](#)



Published online: 28 Mar 2022.



Submit your article to this journal [↗](#)



Article views: 1233



View related articles [↗](#)



View Crossmark data [↗](#)



Citing articles: 18 View citing articles [↗](#)



Variable Selection with Multiply-Imputed Datasets: Choosing Between Stacked and Grouped Methods

Jiacong Du^a, Jonathan Boss^a, Peisong Han^a, Lauren J. Beesley^{a, }, Michael Kleinsasser^a, Stephen A. Goutman^b, Stuart Batterman^c, Eva L. Feldman^b, and Bhramar Mukherjee^a

^aDepartment of Biostatistics, University of Michigan, Ann Arbor, MI; ^bDepartment of Neurology, University of Michigan, Ann Arbor, MI; ^cDepartment of Environmental Health Science, University of Michigan, Ann Arbor, MI

ABSTRACT

Penalized regression methods are used in many biomedical applications for variable selection and simultaneous coefficient estimation. However, missing data complicates the implementation of these methods, particularly when missingness is handled using multiple imputation. Applying a variable selection algorithm on each imputed dataset will likely lead to different sets of selected predictors. This article considers a general class of penalized objective functions which, by construction, force selection of the same variables across imputed datasets. By pooling objective functions across imputations, optimization is then performed jointly over all imputed datasets rather than separately for each dataset. We consider two objective function formulations that exist in the literature, which we will refer to as “stacked” and “grouped” objective functions. Building on existing work, we (i) derive and implement efficient cyclic coordinate descent and majorization-minimization optimization algorithms for continuous and binary outcome data, (ii) incorporate adaptive shrinkage penalties, (iii) compare these methods through simulation, and (iv) develop an R package *miselect*. Simulations demonstrate that the “stacked” approaches are more computationally efficient and have better estimation and selection properties. We apply these methods to data from the University of Michigan ALS Patients Biorepository aiming to identify the association between environmental pollutants and ALS risk. Supplementary materials for this article are available online.

ARTICLE HISTORY

Received June 2020
Revised January 2022

KEYWORDS

Elastic net; Group LASSO; Majorization-minimization; Missing data; Multiple imputation; Pooled objective function

1. Introduction

Variable selection in the presence of missing data is a common problem in applied statistics. Existing statistical methods (Table S1) to deal with both missing data and variable selection fall into two broad classes depending on whether the missing data is handled prior to or at the same time as variable selection. While one can, in principle, handle missingness and selection simultaneously (Yang et al. 2005; Johnson et al. 2008; Garcia et al. 2010), in this article we focus on a different scenario, where missing data have already been handled via multiple imputation and we wish to do variable selection post-imputation. In practice, we often find ourselves in the latter scenario for many reasons. Some examples include: (i) existing data sources or analytical pipelines provide the analyst with already imputed datasets (Schafer 2001), (ii) data privacy settings where the missingness mechanism depends on sensitive information that is not needed for downstream analyses, and (iii) practical situations where the analyst wants to use an existing imputation software package to generate the multiply imputed datasets to be used across various analyses. That is, in this article we assume that the imputation model is correctly specified, imputed datasets are given to us, and we focus purely on variable selection techniques post-imputation. In particular, the challenge is to guarantee uniform selection of variables across multiply imputed datasets

and to generate point estimates for the regression coefficients corresponding to the selected variables.

To address this issue, there are some existing methods in the literature proposed for variable selection post-imputation. The most common approach is to define an ad hoc rule for harmonizing selection across imputed datasets (Wood et al. 2008; Lachenbruch 2011; Ghosh et al. 2015). Wood et al. (2008) outlined three rules for finalizing selection after identifying the active set of predictors in each imputed dataset. These rules consider a variable to be selected if it is selected in at least one imputed dataset, if it is selected in all imputed datasets, or if it is selected in at least half of the imputed datasets. Another related approach is to apply variable selection methods on bootstrapped datasets (Heymans et al. 2007; Long and Johnson 2015; Liu et al. 2016). One such method, proposed by Long and Johnson (2015), imputes missing values in the bootstrapped samples from the original data and then applies penalized regression on each bootstrapped imputed dataset. Proportions of each covariate being selected over all bootstrap imputed datasets are provided and the final active set is determined by thresholding the proportions. The challenge with such thresholding approaches is that there are no clear guidelines on how to choose the threshold. An alternative class of methods are penalized regression estimators whose objective functions

involve a joint optimization over all imputed datasets. Wan et al. (2015) proposed a penalized regression estimator which is equivalent to fitting the penalized regression model on the stacked imputed datasets. The stacked method can also be interpreted in terms of maximizing an objective function pooled over imputed datasets, where regression coefficients are assumed to be the same across datasets. Chen and Wang (2013) proposed the multiple imputation-least absolute shrinkage and selection operator (MI-LASSO) which utilizes a group LASSO penalty to group regression coefficients of the same variable across imputed datasets together, thus, arriving at uniform variable selection across imputed datasets. This grouped approach is less restrictive than the stacked approach, in that the parameter values are allowed to differ across imputed datasets. The stacked and grouped approaches are appealing because they handle variable selection and regression coefficient estimation simultaneously, and no ad hoc rules are required to determine the final active set.

Despite the advantages of the stacked and grouped approaches, automated implementation of these two methods (Chen and Wang 2013; Wan et al. 2015) has been limited in practice. There is a lack of easy-to-use software for both approaches, and for variable selection with multiple imputed datasets in general. More specifically, off-the-shelf R packages such as *glmnet* (Hastie and Qian 2014) and *gcdnet* (Yang et al. 2017) are inadequate for optimizing the stacked objective function proposed in Wan et al. (2015). Namely, the R package *glmnet* does not exclusively apply adaptive weights to the L_1 penalty, as in adaptive elastic net, and the *gcdnet* package does not implement observation weights. As for the grouped approach with the group LASSO penalty function, *gglasso* (Yang and Zou 2015) cannot be directly used without changing the data structure. The dimension of the new design matrix scales linearly with the number of imputed datasets and introduces many zero entries, which is computationally inefficient. Furthermore, existing approaches are targeted toward continuous outcomes and have not been extended to discrete/categorical outcomes that are abundant in practice. Even for continuous outcomes, the algorithm for MI-LASSO in Chen and Wang (2013) has an important limitation. Specifically, the existing MI-LASSO implementation is an optimization procedure that relies on a local quadratic approximation to the L_1 penalty, which converts the penalized objective function into multiple ridge regressions. Coefficient estimates are always nonzero, requiring a manual threshold for setting small coefficient estimates exactly to zero.

To address the aforementioned limitations, we extend existing work on the stacked and grouped approaches for handling variable selection with multiple imputed datasets to incorporate binary outcomes and adaptive penalties (Section 2). We derive a cyclic coordinate descent algorithm for optimizing the pooled objective functions from the stacked approach and a majorization-minimization (MM) algorithm coupled with block coordinate descent updates to obtain optimizers of the pooled objective functions from the grouped approach (Section 3). Unlike the optimization routine in Chen and Wang (2013), our proposed methods provide exact shrinkage to zero without any ad hoc thresholding. We provide an R package *miselect* available on the Comprehensive R Archive Network (CRAN), allowing users to easily implement the proposed

methods. In our motivating example, we apply these methods to identify persistent organic pollutants (POPs) associated with Amyotrophic Lateral Sclerosis (ALS) susceptibility using data collected from the University of Michigan ALS Patients Biorepository (Section 4). Seventeen of the 23 POPs have between 10% and 60% below their respective detection limits, making complete-case analysis infeasible. Finally, through a simulation study, we compare the performance of the proposed methods in terms of variable selection and estimation accuracy (Section 5).

2. Methods

We will first review and illustrate the issues with identifying an active set of predictors when fitting separate penalized regressions on each individual imputed dataset. We will then present the stacked and grouped approaches for uniform variable selection with imputed datasets. Let \mathbf{X}_d denote the $n \times p$ matrix of predictor variables and \mathbf{Y}_d be the $n \times 1$ vector of responses for the d -th imputed dataset ($d = 1, \dots, D$). Let $\mathbf{X}_{d,i}$ indicate the $p \times 1$ covariate vector for the i th observation in the d th imputed dataset, $Y_{d,i}$ be the response for the i th observation in the d th imputed dataset, and $X_{d,ij}$ denote the j th covariate for the i th observation in the d th imputed dataset.

2.1. Penalized Regression on Individual Datasets

A common approach in practice is to use ad hoc rules for determining the final set of selected variables. Often this proceeds by fitting a penalized regression procedure on each imputed dataset separately:

$$\hat{\boldsymbol{\theta}}_d = \underset{\boldsymbol{\theta}_d}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{i=1}^n \log L(\boldsymbol{\theta}_d | Y_{d,i}, \mathbf{X}_{d,i}) + \lambda P_\alpha(\boldsymbol{\beta}_d) \right\} \quad (1)$$

for $d = 1, \dots, D$ where $\log L$ is the log-likelihood function, $\boldsymbol{\beta}_d$ is the $p \times 1$ vector of regression coefficients for the d -th dataset, and $\boldsymbol{\theta}_d = (\mu_d, \boldsymbol{\beta}_d)$ is the vector of regression parameters including the intercept and coefficients for the d -th imputed dataset. $P_\alpha(\boldsymbol{\beta}_d)$ is the penalty function parameterized by α and $\lambda \in (0, \infty)$ which are tuning parameters that control the relative contribution of the penalty to the overall optimization problem.

In this article, we will focus on four penalty functions: LASSO (Tibshirani 1996), adaptive LASSO (aLASSO) (Zou 2006), elastic net (ENET) (Zou and Hastie 2005) and adaptive elastic net (aENET) (Zou and Zhang 2009). The penalty functions can be expressed as:

1. LASSO: $P_\alpha(\boldsymbol{\beta}_d) = \sum_{j=1}^p |\beta_{d,j}|$
2. aLASSO: $P_\alpha(\boldsymbol{\beta}_d) = \sum_{j=1}^p \hat{a}_{d,j} |\beta_{d,j}|$
3. ENET: $P_\alpha(\boldsymbol{\beta}_d) = \alpha \sum_{j=1}^p |\beta_{d,j}| + (1 - \alpha) \sum_{j=1}^p \beta_{d,j}^2$
4. aENET: $P_\alpha(\boldsymbol{\beta}_d) = \alpha \sum_{j=1}^p \hat{a}_{d,j} |\beta_{d,j}| + (1 - \alpha) \sum_{j=1}^p \beta_{d,j}^2$

$\beta_{d,j}$ is the coefficient for the j th covariate in the d th dataset. The intuitive appeal of adaptive weights is that they allow differential penalization of each covariate based on an initial estimate of the regression coefficient vector $\hat{\boldsymbol{\beta}}_d^0$ (smaller $\hat{\beta}_{d,j}^0$ implies a harsher penalty on $\beta_{d,j}$). Moreover, adaptive penalties are known to

address estimation and selection consistency issues that have been observed for nonadaptive penalties (Zou 2006; Zou and Zhang 2009). Here, the adaptive weight $\hat{a}_{d,j} = (|\hat{\beta}_{d,j}^0| + 1/n)^{-\gamma}$ for some $\gamma > 0$, where $\hat{\beta}_{d,j}^0$ is an initial estimate of $\beta_{d,j}$, is often determined from ordinary least squares (OLS) or maximum likelihood estimation when p is smaller than n . If p is larger than n , $\hat{\beta}_{d,j}^0$ can be obtained using LASSO or ENET depending on the correlation structure of the predictors. For the purposes of this article, we use ENET initial values to calculate the adaptive weights for both aLASSO and aENET. To avoid tuning on γ , we follow Zou and Zhang (2009), which fixes $\gamma = \lceil 2\nu/1 - \nu \rceil + 1$ where $\nu = \log(p)/\log(n)$. The $1/n$ term prevents division by zero for initial regression coefficient estimates that are exactly equal to zero.

For illustrative purposes, suppose that there are three imputed datasets and that LASSO is fit separately on each imputed dataset. Furthermore, assume that \mathbf{X}_d is an $n \times 2$ matrix for all $d = 1, 2, 3$, that is, we are applying LASSO shrinkage to coefficients of only two covariates. Figure 1 visualizes the geometry of the constrained region in this hypothetical scenario. As shown in the figure, the maximum likelihood estimates for each of the imputed datasets are slightly different and, when shrunk to the constrained region, lead to two cases: (i) $\hat{\beta}_{d,1} = 0$ and $\hat{\beta}_{d,2} \neq 0$ and (ii) $\hat{\beta}_{d,1} \neq 0$ and $\hat{\beta}_{d,2} = 0$. This toy example illustrates the fundamental problem of trying to harmonize variable selection across imputed datasets without borrowing information across imputed datasets.

2.2. Stacked Objective Functions

In a stacked objective function, we sum the loss functions for each of the imputed datasets together and jointly optimize the

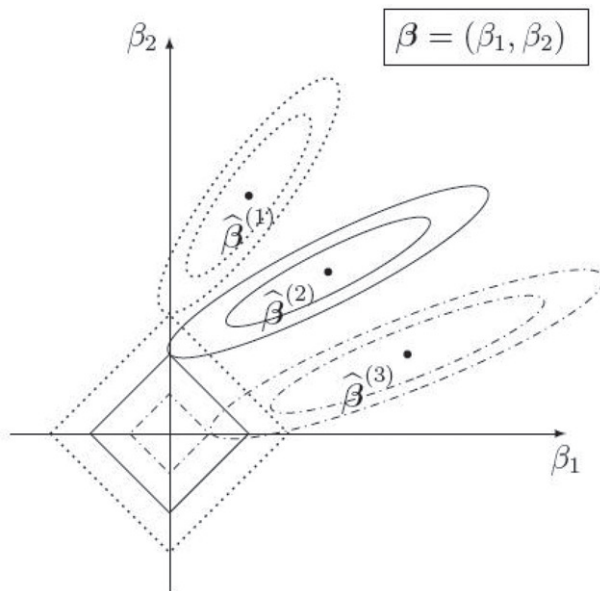


Figure 1. Example of LASSO fit separately on each imputed dataset with two predictor variables. β_1 and β_2 represent the first and second coefficient, respectively. $\hat{\beta}^{(d)}$ denotes the maximum likelihood estimate in the d th imputed dataset without the constraint ($d = 1, 2, 3$). The diamond-shaped areas around the origin represent the resulting constrained regions for the imputed datasets. The contours surrounding the maximum likelihood estimates represent the loss function.

collective objective function:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n \log L(\theta | Y_{d,i}, \mathbf{X}_{d,i}) + \lambda P_{\alpha}(\beta) \right\}. \quad (2)$$

Note that θ is not indexed by d . This implies that optimizing the pooled objective function will result in one estimated parameter vector $\hat{\theta}$, thereby enforcing uniform selection across all imputed datasets.

A nice feature of the stacked objective function is that optimization is straightforward. Namely, optimization of the stacked objective function is equivalent to stacking the imputed datasets and fitting the desired penalized regression algorithm on the stacked imputed datasets with existing software. Therefore, the stacked objective function provides a framework for pooling penalized regression estimates across imputed datasets for a general class of objective functions. However, stacking all imputed datasets can be viewed as artificially increasing the sample size. A common way to address this is to add a weight to each subject, $o_i = 1/D$, so that the total weight for each subject in the stacked dataset sums up to one. An alternative weight specification, proposed in Wan et al. (2015), is $o_i = f_i/D$, where f_i is the number of observed predictors out of the total number of predictors for subject i , which accounts for varying degrees of missing information for each subject. That being said, upweighting subjects with less missingness and downweighting subjects with more missingness can, in some sense, be viewed as making the optimization more like complete-case analysis, which might be problematic for Missing at Random (MAR) and Missing not at Random (MNAR) scenarios. Going forward, we consider both equal weights ($o_i = 1/D$) and the observation weights proposed in Wan et al. (2015) ($o_i = f_i/D$). The weighted stacked objective function can be written as

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i \log L(\theta | Y_{d,i}, \mathbf{X}_{d,i}) + \lambda P_{\alpha}(\beta) \right\}. \quad (3)$$

In (3), we directly sum the log-likelihoods ignoring the correlation for the same observation across all imputed datasets. This is not a problem for point estimation. Averaging over the number of imputed datasets D can be viewed as obtaining the expected log-likelihood conditional on the observed data and then subsequently averaging over all n observations leads to the conditional log-likelihood given the observed data. This strategy of defining an objective function summed across correlated imputed datasets is well-founded in the statistical literature on missing data (Wang and Robins 1998; Robins and Wang 2000). Ignoring the dependence between these datasets will have an impact on estimation of standard errors, but the coefficient estimator will still be consistent (Wang and Robins 1998; Robins and Wang 2000).

From this stacked approach, we extend the penalty functions to LASSO, aLASSO, ENET and aENET. When equal weights are used, the names for different versions of the stacked methods are as follows:

1. Stacked LASSO (SLASSO): $P_{\alpha}(\beta) = \sum_{j=1}^p |\beta_j|$
2. Stacked adaptive LASSO (SaLASSO): $P_{\alpha}(\beta) = \sum_{j=1}^p \hat{a}_j |\beta_j|$

3. Stacked elastic net (SENET): $P_\alpha(\boldsymbol{\beta}) = \alpha \sum_{j=1}^p |\beta_j| + (1 - \alpha) \sum_{j=1}^p \beta_j^2$
4. Stacked adaptive elastic net (SaENET): $P_\alpha(\boldsymbol{\beta}) = \alpha \sum_{j=1}^p \hat{a}_j |\beta_j| + (1 - \alpha) \sum_{j=1}^p \beta_j^2$

where $\hat{a}_j = (|\hat{\beta}_j| + 1/(nD))^{-\gamma}$ and $\hat{\beta}_j$ is estimated through SENET. Following Zou and Zhang (2009), γ is fixed to be $\lceil 2\nu/1 - \nu \rceil + 1$, where $\nu = \log(p)/\log(nD)$.

When $o_i = f_i/D$, the penalized methods are named SLASSO(w), SaLASSO(w), SENET(w), and SaENET(w), respectively. Adaptive weights for the stacked approach with the observation weights are calculated using SENET(w).

2.3. Grouped Objective Functions

An alternative to the stacked objective function is the grouped objective function. This method imposes uniform variable selection across imputed datasets by adding a group LASSO penalty across imputed datasets (Yuan and Lin 2006; Chen and Wang 2013; Geronimi and Saporta 2017). The optimizer of the grouped objective function can be mathematically expressed as

$$(\hat{\boldsymbol{\theta}}_1, \dots, \hat{\boldsymbol{\theta}}_D) = \underset{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_D}{\operatorname{argmin}} \left\{ -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n \log L(\boldsymbol{\theta}_d | Y_{d,i}, \mathbf{X}_{d,i}) + \lambda P(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_D) \right\} \quad (4)$$

where $\lambda \in (0, \infty)$ is a tuning parameter. Chen and Wang (2013) originally formulated a special case of the grouped objective function known as MI-LASSO, where the penalty function is

$$P(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_D) = \sum_{j=1}^p \sqrt{\sum_{d=1}^D \beta_{d,j}^2}.$$

In the grouped objective function, the parameter vector is now indexed by d , meaning that $\boldsymbol{\theta}_1 \neq \dots \neq \boldsymbol{\theta}_D$. Although the $\boldsymbol{\theta}_d$'s are not identical, for any fixed j , the group LASSO penalty jointly shrink all $\beta_{d,j}$'s to zero, that is, $\beta_{1,j} = \dots = \beta_{D,j} = 0$. This allows for uniform selection across imputed datasets but also allows for variability in the nonzero estimated coefficients across imputed datasets. Based on Chen and Wang (2013), we consider the following penalties:

1. Group LASSO (GLASSO): $P(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_D) = \sum_{j=1}^p \sqrt{\sum_{d=1}^D \beta_{d,j}^2}$
2. Group adaptive LASSO (GaLASSO): $P(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_D) = \sum_{j=1}^p \hat{a}_j \sqrt{\sum_{d=1}^D \beta_{d,j}^2}$

where $\hat{a}_j = (\sqrt{\sum_{d=1}^D \hat{\beta}_{d,j}^2} + 1/(nD))^{-\gamma}$. $\hat{\beta}_{d,j}$ is estimated from GLASSO, $\gamma = \lceil 2\nu/1 - \nu \rceil + 1$, and $\nu = \log(pD)/\log(nD)$ (Zou and Zhang 2009).

Remark. We do not extend the grouped LASSO based approaches to grouped ENET based approaches because grouping itself provides constraints similar to the ENET penalty. How to effectively account for correlations among predictors using the grouped penalty remains an open question.

3. Optimization

In this section, we outline the optimization routines for SaENET(w) and GaLASSO with binary outcomes. To optimize the SaENET(w) objective function, we use local quadratic approximation coupled with a cyclic coordinate descent algorithm (Friedman et al. 2010). To obtain an optimizer of the GaLASSO objective function, we use a MM algorithm combined with block coordinate descent updates to handle the group LASSO component of the penalty function. The other objective functions listed in Section 2 are special cases of SaENET(w) and GaLASSO, and therefore, can be optimized with minor simplifications (see Figure 2). Namely, when $\alpha = 1$, SaENET(w) reduces to SaLASSO(w), and SENET(w) reduces to SLASSO(w). Furthermore, when $\hat{a}_j = 1$ for $j = 1, 2, \dots, p$, SaENET(w) reduces to SENET(w) and GaLASSO reduces to GLASSO. Optimizing the stacked objective function with equal weights can be achieved by setting $o_i = 1/D$.

3.1. Optimization of SaENET(w)

Without loss of generality, suppose that the variables in the design matrix are standardized after stacking all the imputed datasets. That is, $\sum_{d=1}^D \sum_{i=1}^n x_{d,ij} = 0$, $\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n x_{d,ij}^2 = 1$, for $j = 1, 2, \dots, p$. Let $\boldsymbol{\eta} = (\boldsymbol{\eta}_1^T, \boldsymbol{\eta}_2^T, \dots, \boldsymbol{\eta}_D^T)^T$ be the linear predictor based on $\boldsymbol{\theta} = (\boldsymbol{\mu}, \boldsymbol{\beta}^T)^T$ such that $\eta_{d,i} = \boldsymbol{\mu} + \mathbf{x}_{d,i}^T \boldsymbol{\beta}$. Then the sum of the weighted loss functions can be expressed as

$$L(\boldsymbol{\eta}) = \frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i \left\{ -y_{d,i} \eta_{d,i} + \log(1 + \exp(\eta_{d,i})) \right\}.$$

Let $\boldsymbol{\eta}^{(t)}$ denote the linear predictor at the t th iteration. Following Friedman et al. (2010), we use a Taylor expansion at $\boldsymbol{\eta}^{(t)}$ to construct a quadratic approximation to the loss function:

$$L_Q(\boldsymbol{\eta} | \boldsymbol{\eta}^{(t)}) = \frac{1}{2n} \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} (\tilde{y}_{d,i} - \boldsymbol{\mu} - \mathbf{x}_{d,i}^T \boldsymbol{\beta})^2 + C(\boldsymbol{\eta}^{(t)})^2,$$

where

$$\tilde{y}_{d,i} = \boldsymbol{\mu}^{(t)} + \mathbf{x}_{d,i}^T \boldsymbol{\beta}^{(t)} + \frac{y_{d,i} - \tilde{p}(x_{d,i})}{\tilde{p}(x_{d,i})(1 - \tilde{p}(x_{d,i}))} \quad (5)$$

is the working response, $w_{d,i} = \tilde{p}(x_{d,i})(1 - \tilde{p}(x_{d,i}))$ is a subject weight that is specific to each imputed dataset, and

$$\begin{aligned} \tilde{p}(x_{d,i}) &= P(Y_{d,i} = 1 | X_{d,i} = x_{d,i}) \\ &= \left(1 + \exp \left(-(\boldsymbol{\mu}^{(t)} + \mathbf{x}_{d,i}^T \boldsymbol{\beta}^{(t)}) \right) \right)^{-1}. \end{aligned}$$

Going forward we will use $O_Q = L_Q(\boldsymbol{\eta} | \boldsymbol{\eta}^{(t)}) + \lambda P_\alpha(\boldsymbol{\beta})$, as shorthand notation for the objective function after quadratic approximation. We then use coordinate descent to solve:

$$\underset{\boldsymbol{\theta}}{\operatorname{argmin}} O_Q = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \left\{ L_Q(\boldsymbol{\eta} | \boldsymbol{\eta}^{(t)}) + \lambda P_\alpha(\boldsymbol{\beta}) \right\}.$$

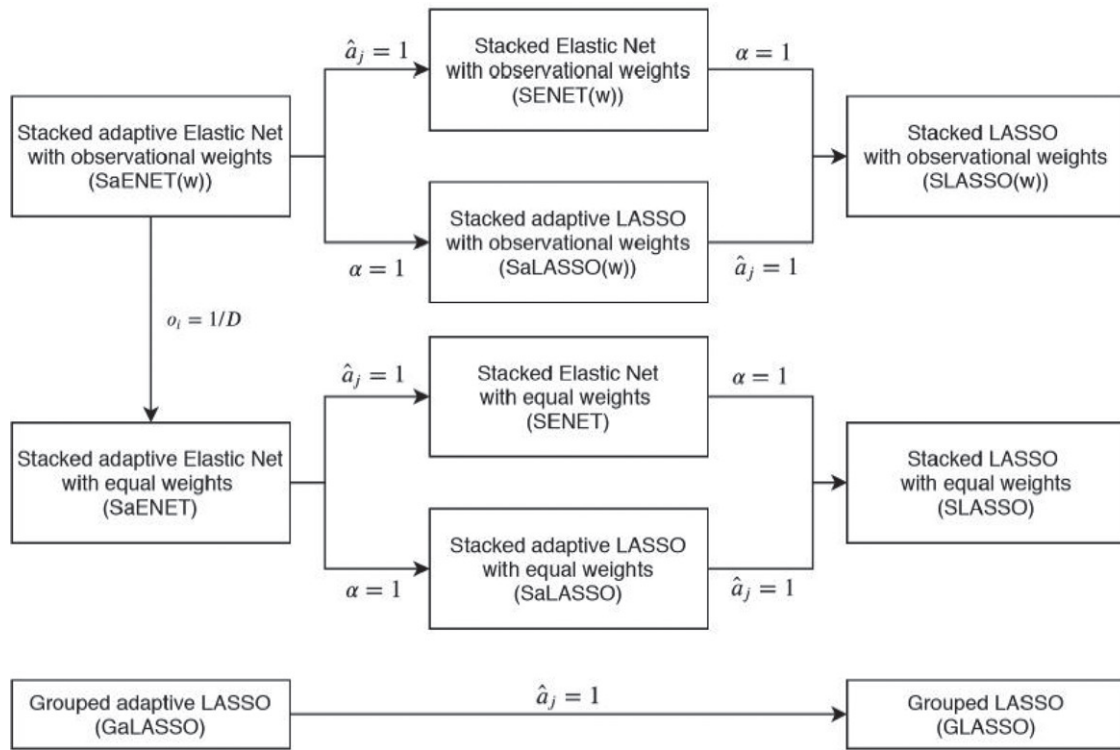


Figure 2. Illustration of how to obtain penalized pooled objective function methods, which are special cases of SaENET(w) and GaLASSO.

The derivative of the approximate objective function with respect to the intercept parameter is:

$$\begin{aligned}\frac{\partial O_Q}{\partial \mu} &= -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} (\tilde{y}_{d,i} - \mu - x_{d,i}^T \boldsymbol{\beta}) \\ &= -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} (\tilde{y}_{d,i} - x_{d,i}^T \boldsymbol{\beta}) + \frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} \mu.\end{aligned}$$

Let $z_0 = \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} (\tilde{y}_{d,i} - x_{d,i}^T \boldsymbol{\beta})$. Then μ can be updated as

$$\mu^{(t+1)} \leftarrow \frac{z_0}{\sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i}}. \quad (6)$$

If $\beta_j > 0$,

$$\begin{aligned}\frac{\partial O_Q}{\partial \beta_j} &= -\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} x_{d,ij} (\tilde{y}_{d,i} - \mu - x_{d,i(-j)}^T \boldsymbol{\beta}_{(-j)}) \\ &\quad + \frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} x_{d,ij}^2 \beta_j + \lambda \alpha \hat{a}_j + 2\lambda(1-\alpha)\beta_j\end{aligned}$$

where $x_{d,i(-j)}$ refers to the value of the covariate vector for the i th observation in the d th imputed dataset after removing the j th covariate, and $\boldsymbol{\beta}_{(-j)}$ refers to the regression coefficient vector $\boldsymbol{\beta}$ without the j th entry. If $\beta_j < 0$, then the derivative of O_Q with respect to β_j is the same as the derivative when $\beta_j > 0$, the only exception being that the $\lambda \alpha \hat{a}_j$ term becomes $-\lambda \alpha \hat{a}_j$. Setting $z_j = \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} x_{d,ij} (\tilde{y}_{d,i} - \mu - x_{d,i(-j)}^T \boldsymbol{\beta}_{(-j)})$, then β_j can be updated as

$$\beta_j^{(t+1)} \leftarrow \frac{S(\frac{1}{n} z_j, \lambda \alpha \hat{a}_j)}{\frac{1}{n} \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} x_{d,ij}^2 + 2\lambda(1-\alpha)} \quad (7)$$

where $S(z, \lambda)$ is the soft-thresholding operator:

$$S(z, \lambda) = \begin{cases} 0 & \text{if } |z| \leq \lambda \\ z - \lambda & \text{if } z > \lambda \\ z + \lambda & \text{if } z < -\lambda \end{cases}$$

A summary of the optimization routine for SaENET(w) is presented in Algorithm 1. One thing to note is that we never directly compute z_j after each update of β_j . A more computationally efficient approach is to update z_j through the residual $r_{d,i} = \tilde{y}_{d,i} - \mu - x_{d,i}^T \boldsymbol{\beta}$. Specifically,

$$z_j \leftarrow \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} x_{d,ij} r_{d,i} + \sum_{d=1}^D \sum_{i=1}^n o_i w_{d,i} x_{d,ij}^2 \beta_j^{(t)} \quad (8)$$

$$r_{d,i} \leftarrow r'_{d,i} - x_{d,ij} (\beta_j^{(t)} - \beta_j^{(t+1)})$$

where $r'_{d,i}$ indicates the previous estimate of $r_{d,i}$ (we do not use the (t) superscript notation here because $r_{d,i}$ is updated multiple times within the same iteration). Updating the intercept parameter μ is similar, in that we assign $r_{d,i} \leftarrow r'_{d,i} - (\mu^{(t)} - \mu^{(t+1)})$ and update z_0 accordingly. Once a coefficient is shrunk to zero, it will stay at zero for the remaining iterations.

3.2. Optimization of GaLASSO

Although Chen and Wang (2013) initially proposed the idea of grouped objective functions, their optimization procedure relies on a local quadratic approximation argument to rewrite the grouped objective function as the sum of D separate ridge regressions. Because the ridge penalty does not shrink

Algorithm 1: Algorithm for SaENET(w)

Stack all imputed datasets

while *stop criteria is not met* **do**

update the linear predictor

$$\boldsymbol{\eta} \leftarrow \left\{ \left\{ \mu^{(t)} + x_{d,i}^T \boldsymbol{\beta}^{(t)} \right\}_{i=1}^n \right\}_{d=1}^D$$

$$\tilde{\boldsymbol{p}} \leftarrow \left\{ \left\{ \tilde{p}(x_{d,i}) \right\}_{i=1}^n \right\}_{d=1}^D$$

 update working response $\tilde{\boldsymbol{y}}$ according to (5)

$$\text{update residual } \boldsymbol{r} \leftarrow \left\{ \left\{ \tilde{y}_{d,i} - \eta_{d,i} \right\}_{i=1}^n \right\}_{d=1}^D$$

 update the intercept $\mu^{(t+1)}$ according to (6)

update the residual

$$\boldsymbol{r} \leftarrow \left\{ \left\{ r'_{d,i} - (\mu^{(t)} - \mu^{(t+1)}) \right\}_{i=1}^n \right\}_{d=1}^D$$

for j *in the active set* **do** update z_j according to (8) update $\beta_j^{(t+1)}$ according to (7)

update residual

$$\boldsymbol{r} \leftarrow \left\{ \left\{ r'_{d,i} - x_{d,ij}(\beta_j^{(t)} - \beta_j^{(t+1)}) \right\}_{i=1}^n \right\}_{d=1}^D$$

end**end**

regression coefficient estimates all the way to zero, the user is then forced to threshold the resulting regression coefficient estimates in an ad hoc manner. Another limitation of the algorithm proposed in Chen and Wang (2013) is that they only consider continuous outcome variables. In many biomedical applications, including our motivating example, the outcome is a binary indicator of disease status. To address both of these concerns, we developed a MM algorithm for GaLASSO with binary outcome data. The primary computational advantage of the MM-Algorithm is that it allows us to transform the optimization of a nonlinear loss function into an optimization of a linear function, called the majorizing function. Block coordinate descent can then be applied to optimize the linear majorizing function, where the blocks correspond to regression coefficient estimates for each covariate across all imputed datasets, that is, $(\beta_{1,j}, \beta_{2,j}, \dots, \beta_{D,j})$ for the j th covariate.

Without loss of generality, we standardize the data by letting $\sum_{i=1}^n x_{d,ij} = 0$, $n^{-1} \sum_{i=1}^n x_{d,ij}^2 = 1$, for $j = 1, 2, \dots, p$ and $d = 1, 2, \dots, D$. Let $\boldsymbol{\eta} = (\boldsymbol{\eta}_1^T, \boldsymbol{\eta}_2^T, \dots, \boldsymbol{\eta}_D^T)^T$ be the linear predictor based on $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D)$, where $\boldsymbol{\theta}_d = (\mu_d, \boldsymbol{\beta}_d^T)^T$, $\boldsymbol{\eta}_d = (\eta_{d,1}, \eta_{d,2}, \dots, \eta_{d,n})^T$ and $\eta_{d,i} = \mu_d + \beta_{d,1}x_{d,i1} + \dots + \beta_{d,p}x_{d,ip}$. For binary outcome data, the sum of loss functions in (4) can be written as

$$L(\boldsymbol{\eta}) = \sum_{d=1}^D \sum_{i=1}^n \left\{ -y_{d,i} \eta_{d,i} + \log(1 + \exp(\eta_{d,i})) \right\}.$$

Let $\boldsymbol{\eta}^{(t)}$ be the linear predictor at the t -th iteration, $L(\boldsymbol{\eta}^{(t)})$ be the loss function given $\boldsymbol{\eta}^{(t)}$, and $\tilde{L}(\boldsymbol{\eta}|\boldsymbol{\eta}^{(t)})$ be the majorizing function of $L(\boldsymbol{\eta}^{(t)})$ (Breheny and Huang 2015), which can be

expressed as

$$\begin{aligned} \tilde{L}(\boldsymbol{\eta}|\boldsymbol{\eta}^{(t)}) &= L(\boldsymbol{\eta}^{(t)}) + (\boldsymbol{\eta} - \boldsymbol{\eta}^{(t)})^T \nabla L(\boldsymbol{\eta}^{(t)}) \\ &\quad + \frac{\nu}{2} (\boldsymbol{\eta} - \boldsymbol{\eta}^{(t)})^T (\boldsymbol{\eta} - \boldsymbol{\eta}^{(t)}). \end{aligned}$$

Here, $\nu = \max_d \max_i \sup_{\boldsymbol{\eta}} \{\nabla^2 L_{d,i}(\boldsymbol{\eta})\} = 0.25$ (Breheny and Huang 2015). Ignoring constants not involving $\boldsymbol{\eta}$, $\tilde{L}(\boldsymbol{\eta}|\boldsymbol{\eta}^{(t)})$ can be written in terms of $(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D)$:

$$\begin{aligned} \tilde{L}(\boldsymbol{\eta}|\boldsymbol{\eta}^{(t)}) &= \tilde{L}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D) \\ &\propto \frac{\nu}{2n} \sum_{d=1}^D \sum_{i=1}^n (\tilde{y}_{d,i} - \mu_d - x_{d,i}^T \boldsymbol{\beta}_d)^2 \end{aligned}$$

where

$$\tilde{y}_{d,i} = \mu_d^{(t)} + x_{d,i}^T \boldsymbol{\beta}_d^{(t)} + \frac{y_{d,i} - \tilde{p}(x_{d,i})}{\nu} \quad (9)$$

is the working response and

$$\begin{aligned} \tilde{p}(x_{d,i}) &= P(Y_{d,i} = 1 | X_{d,i} = x_{d,i}) \\ &= \left(1 + \exp \left(- \left(\mu_d^{(t)} + x_{d,i}^T \boldsymbol{\beta}_d^{(t)} \right) \right) \right)^{-1}. \end{aligned}$$

Going forward, let $M_Q = \tilde{L}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D) + \lambda P(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_D)$ denote the penalized majorizing function, which we want to optimize

$$\underset{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D}{\operatorname{argmin}} M_Q = \underset{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D}{\operatorname{argmin}} \left\{ \tilde{L}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_D) + \lambda P(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \dots, \boldsymbol{\beta}_D) \right\}.$$

To derive the block coordinate descent updates for the majorizing function optimization, we first need to introduce some new notation to distinguish coefficients within an imputed dataset from those across the imputed datasets. Let $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_D)^T$ be the vector of all intercept parameters across the imputed datasets, let $\boldsymbol{\beta}_{\cdot,j} = (\beta_{1,j}, \beta_{2,j}, \dots, \beta_{D,j})^T$, for $j = 1, 2, \dots, p$ denote the vector of the j th regression coefficients across all imputed datasets, and let $\boldsymbol{\beta}_{d,\cdot} = (\beta_{d,1}, \beta_{d,2}, \dots, \beta_{d,p})^T$ for $d = 1, 2, \dots, D$ indicate the regression coefficients for the d th imputed dataset. If we take the subdifferential of M_Q with respect to $\boldsymbol{\mu}$, we get:

$$\frac{\partial M_Q}{\partial \boldsymbol{\mu}} = -\frac{\nu}{n} \boldsymbol{z}_0 + \nu \boldsymbol{\mu}$$

where

$$\boldsymbol{z}_0 = \begin{bmatrix} \sum_{i=1}^n (\tilde{y}_{1,i} - x_{1,i}^T \boldsymbol{\beta}_{1,\cdot}) \\ \sum_{i=1}^n (\tilde{y}_{2,i} - x_{2,i}^T \boldsymbol{\beta}_{2,\cdot}) \\ \vdots \\ \sum_{i=1}^n (\tilde{y}_{D,i} - x_{D,i}^T \boldsymbol{\beta}_{D,\cdot}) \end{bmatrix}.$$

Since $x_{d,ij}$ have been centered, the intercept can be updated as:

$$\boldsymbol{\mu}^{(t+1)} \leftarrow \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n \tilde{y}_{1,i} \\ \frac{1}{n} \sum_{i=1}^n \tilde{y}_{2,i} \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n \tilde{y}_{D,i} \end{bmatrix}. \quad (10)$$

We now derive the update for the block of coefficients $\beta_{\cdot,j}$, $j = 1, \dots, p$. If $\beta_{\cdot,j} \neq \mathbf{0}$, the subdifferential of M_Q with respect to $\beta_{\cdot,j}$ is

$$\frac{\partial M_Q}{\partial \beta_{\cdot,j}} = -\frac{\nu}{n} \mathbf{z}_j + \nu \beta_{\cdot,j} + \lambda \hat{a}_j \frac{\beta_{\cdot,j}}{\|\beta_{\cdot,j}\|}$$

where

$$\mathbf{z}_j = \begin{bmatrix} \sum_{i=1}^n x_{1,ij} (\tilde{y}_{1,i} - \mu_1 - x_{1,i(-j)}^T \beta_{1,(-j)}) \\ \sum_{i=1}^n x_{2,ij} (\tilde{y}_{2,i} - \mu_2 - x_{2,i(-j)}^T \beta_{2,(-j)}) \\ \vdots \\ \sum_{i=1}^n x_{D,ij} (\tilde{y}_{D,i} - \mu_D - x_{D,i(-j)}^T \beta_{D,(-j)}) \end{bmatrix}.$$

Here, $x_{d,i(-j)}$ is the vector of covariates for the i th observation in the d th dataset after removing $x_{d,ij}$, and $\beta_{d,(-j)}$ is the corresponding vector of regression coefficients. The form of the partial derivative with respect to $\beta_{\cdot,j}$ indicates that the $\beta_{\cdot,j}$ update must lie on the line segment joining the zero vector, $\mathbf{0}$, and \mathbf{z}_j . Therefore, $\beta_{\cdot,j}$ can be updated as

$$\beta_{\cdot,j}^{(t+1)} \leftarrow \frac{1}{\nu} S\left(\frac{\nu}{n} \|\mathbf{z}_j\|, \lambda \hat{a}_j\right) \frac{\mathbf{z}_j}{\|\mathbf{z}_j\|}. \quad (11)$$

As with SaENET(w), we do not directly calculate \mathbf{z}_j but use the residual $r_{d,i}$ to update \mathbf{z}_j . Let $r_{d,i} = \tilde{y}_{d,i} - \mu_d - x_{d,i}^T \beta_{d,\cdot}$. Then

$$\mathbf{z}_j \leftarrow \begin{bmatrix} z_{1,j} \\ z_{2,j} \\ \vdots \\ z_{D,j} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n x_{1,ij} r_{1,i} + n \beta_{1,j}^{(t)} \\ \sum_{i=1}^n x_{2,ij} r_{2,i} + n \beta_{2,j}^{(t)} \\ \vdots \\ \sum_{i=1}^n x_{D,ij} r_{D,i} + n \beta_{D,j}^{(t)} \end{bmatrix} \quad (12)$$

$$r_{d,i} \leftarrow r'_{d,i} - x_{d,ij} (\beta_{d,j}^{(t)} - \beta_{d,j}^{(t+1)}).$$

Updating the intercept parameter μ is similar, in that we assign $r_{d,i} \leftarrow r'_{d,i} - (\mu_d^{(t)} - \mu_d^{(t+1)})$ and update μ_0 accordingly. Once the block of regression coefficients corresponding to the j th covariate is shrunk to zero, it will stay at zero for the remaining iterations. A summary of the optimization routine is described in Algorithm 2.

3.3. Tuning Parameters

The implementation of SaENET(w) and GaLASSO in the *mis-elect* package sets tuning parameter values based on a 5-fold cross-validation routine combined with a “one-standard-error” rule (Hastie et al. 2009, 2015). More specifically, GaLASSO chooses the largest λ , such that the cross-validation error is within one standard error of the minimum cross-validation error. SaENET(w) selects an (α, λ) pair by first identifying all dyads whose cross-validation errors are within one standard error of the minimum cross-validation error, and then picking the (α, λ) pair with the largest L_1 penalty, that is, $\lambda\alpha$. For both GaLASSO and SaENET(w), the default sequence for λ ranges from λ_{\min} to λ_{\max} on the log scale where λ_{\max} is chosen to be the smallest value where all the coefficients are shrunk to zero, and $\lambda_{\min} = 10^{-6} \lambda_{\max}$. Since $\alpha \in [0, 1]$, then for SaENET(w) we can choose a sequence of values, say $(0.2, 0.4, \dots, 0.8, 1)$, to

Algorithm 2: Algorithm for GaLASSO

```

while stop criteria is not met do
  for  $d = 1, 2, \dots, D$  do
    update the linear predictor
     $\eta_d \leftarrow \{\mu_d^{(t)} + x_{d,i}^T \beta_{d,\cdot}^{(t)}\}_{i=1}^n$ 
     $\tilde{p}_d \leftarrow \{\tilde{p}_{d,i}\}_{i=1}^n$ 
    update working response  $\tilde{y}_d \leftarrow \{\tilde{y}_{d,i}\}_{i=1}^n$  according to (9)
    update the residual  $r_d \leftarrow \{\tilde{y}_{d,i} - \eta_{d,i}\}_{i=1}^n$ 
  end
  update  $\mu^{(t+1)}$  according to (10)
  update the residual
   $r \leftarrow \left\{ \{r'_{d,i} - (\mu_d^{(t)} - \mu_d^{(t+1)})\}_{i=1}^n \right\}_{d=1}^D$ 
  for  $j$  in the active set do
    update  $\mathbf{z}_j$  according to (12)
    update  $\beta_{\cdot,j}^{(t+1)}$  according to (11)
    update the residual
     $r \leftarrow \left\{ \{r'_{d,i} - x_{d,ij} (\beta_{d,j}^{(t)} - \beta_{d,j}^{(t+1)})\}_{i=1}^n \right\}_{d=1}^D$ 
  end
end

```

fully explore the tradeoff between L_1 and L_2 shrinkage. When the nonadaptive methods are used, that is, GLASSO, SENET(w), SLASSO(w), we set $\lambda_{\min} = 10^{-3} \lambda_{\max}$. The reason that λ_{\min} is smaller for the adaptive methods is to prevent large adaptive weights from overwhelming the overall shrinkage.

One caveat with cross-validation to select tuning parameters for the stacked objective function approach is that, by stacking the imputed datasets on top of one another, there are now D rows corresponding to the same subject. Therefore, if the D rows corresponding to the same subject are distributed across the validation folds, then we are prone to overfitting. To prevent this issue, we restrict the fold assignment such that all of the rows corresponding to the same subject are assigned to a single validation fold.

4. Data Example

Our motivating example is derived from data in the University of Michigan ALS Patient Biorepository which aims to identify environmental risk factors associated with ALS (Yu et al. 2014; Su et al. 2016; Goutman et al. 2019). ALS is progressive disease primarily involving motor neuron cells in the brain and spinal cord leading to weakness of voluntary muscles and death within 2–4 years due to respiratory failure (Goutman 2017). ALS has a complex etiology driven by the combination of genetic susceptibility and environmental exposures (Al-Chalabi and Hardiman 2013; Paez-Colasante et al. 2015; Goutman et al. 2018). In this particular study we are interested in characterizing the relationship between persistent organic pollutant (POP) exposure and ALS susceptibility. In total, 167 ALS cases and 99 healthy controls were recruited between 2011 and 2014. All participants provided written informed consent and the study was IRB

approved (Su et al. 2016; Goutman et al. 2019). Participants provided their weight and height measurements at the time of study enrollment and five years prior to enrollment and their educational attainment. Plasma samples were collected from each study participant to measure 122 potentially neurotoxic POPs, which can be broadly partitioned into three chemical classes: organochlorine pesticides, polychlorinated biphenyls (PCBs), and polybrominated diphenyl ethers (PBDEs). Of the 122 POPs, a subset of 23 POPs with less than 60% nondetects were used for further analysis. More information regarding data collection and study protocols can be found in (Su et al. 2016; Goutman et al. 2019).

From a statistical perspective, the target model of interest is a penalized logistic regression model, where the outcome is ALS case/control status and the design matrix of predictors and covariates contains the 23 log-transformed POPs, age, sex, body mass index (BMI), rate of BMI change over the five years prior to survey consent, and education. The confounders are selected based on existing ALS literature (Chio et al. 2009). Elastic net regularization is of particular interest here, because many of the POPs have medium pairwise correlations with one another (Figure S1) (Goutman et al. 2019). We only penalize regression coefficients corresponding to POPs to ensure that adjustment covariates are retained in the final model. If we look at the percent missingness (Table S2), we observe that 9 of the 23 variables have more than 30% below the detection limit and 24.4% of subjects have incomplete BMI. Complete-case analysis in this context is infeasible, as every control is missing at least one covariate or has at least one measured POP below its respective detection limit. POP concentrations below their respective detection limits were imputed 10 or 50 times conditional on case/control status following the censored likelihood multiple imputation strategy outlined in Boss et al. (2019). After imputing the exposure nondetects, multiple imputation by chained equations (mice) was used to impute the missing adjustment covariates (Van Buuren and Groothuis-Oudshoorn 2011).

To illustrate the problem with fitting separate penalized regression routines, we first apply LASSO, aLASSO, ENET, and aENET to each imputed dataset and calculate the proportion of imputed datasets in which each variable is selected. In Table 1,

we summarize the results for scenarios with 10 and 50 imputed datasets. The final active set for each method is determined using the ad hoc rules outlined in Wood et al. (2008), namely (i) a variable is considered selected if it is selected in all imputed datasets, (ii) a variable is considered selected if it is selected in at least one imputed dataset, (iii) a variable is considered selected if it is selected in at least half of the imputed datasets. Note that depending on whether we use (i), (ii), or (iii), the number of selected POPs changes, especially when $D = 50$. For example, if ad hoc combining rule (i) is used, then ENET selects PBDE 153, pentachlorobenzene (PeCB), trans-chlordane, cis-nonachlor, and PCB 151. However, ad hoc combining rule (ii) additionally selects PBDE 28, PBDE 99, PCB 110, PCB 174, and PCB 180. As expected, the final active set determined by aLASSO and aENET is much more sparse than their nonadaptive counterparts; if ad hoc combining rule (iii) is used then aLASSO and aENET only select PeCB and cis-nonachlor.

A more subtle point that deserves further comment is that nonuniform POPs selection across imputed datasets makes it difficult to obtain final regression coefficient estimates. For example, consider aLASSO when $D = 50$, which selects PeCB in all 50 imputed datasets. Although PeCB is always selected, the interpretation of the regression coefficient for PeCB is conditional on the other selected exposures. That is, although PeCB is selected for all imputed datasets, the PeCB regression coefficient estimates are not necessarily comparable across imputed datasets because they condition on different sets of selected exposures.

The regression coefficient estimates obtained from the stacked and grouped objective function methods with 50 imputed datasets are in Table 2 and with 10 imputed datasets are in Table S3. Because the grouped methods produce different regression coefficient estimates across imputed datasets, the final estimates are the average of the regression coefficient estimates across imputed datasets. All coefficient estimates are positive, which is consistent with the hypothesis that higher POP exposure positively associates with a higher ALS risk. All of the nonadaptive methods select PBDE 153, PeCB, trans-chlordane, cis-nonachlor, and PCB 151. Similarly, the adaptive methods all select PeCB and cis-nonachlor, however, the equally weighted stacked objective functions additionally select PCB 151 while

Table 1. The proportion of imputed datasets in which each POP is selected.

POPs	LASSO		aLASSO		ENET		aENET	
	D = 10	D = 50	D = 10	D = 50	D = 10	D = 50	D = 10	D = 50
PBDE28	—	—	—	—	0.1	0.1	—	—
PBDE99	—	—	—	—	—	0.04	—	—
PBDE153	1	1	—	0.02	1	1	—	0.02
PeCB	1	1	1	1	1	1	1	1
trans-chlordane	0.9	0.94	0.1	0.08	1	1	0.1	0.08
cis-nonachlor	1	1	1	1	1	1	1	1
PCB 110	—	0.02	—	—	—	0.06	—	—
PCB 151	0.9	0.96	—	0.04	1	1	—	0.06
PCB 174	—	0.02	—	—	0.1	0.08	—	—
PCB 180	—	—	—	—	—	0.02	—	—
Union	5	7	3	5	7	10	3	5
50%-cutoff	5	5	2	2	5	5	2	2
Intersection	3	3	2	2	5	5	2	2

NOTE: Note that only 10 out of the 23 POPs are listed because the other 13 POPs were never selected by LASSO, aLASSO, ENET, or aENET. Bolded entries indicate a selection proportion over 0.50 and dashes denote a selection proportion of zero. D is the number of imputed datasets. The total number of POPs in the final active set, as determined by three ad hoc combining rules, are presented in the rows titled, Union (in at least one dataset), 50%-cutoff (in over 50% datasets), and Intersection (in all datasets).

Table 2. Regression coefficient estimates (log odds ratio corresponding to one IQR increase in each contaminant) for five POPs collected as part of the University of Michigan ALS Patients Biorepository case-control study (167 ALS cases and 99 healthy controls). The last two rows are the averaged number of predictors selected/removed by each method.

POPs	SLASSO	SaLASSO	SENET	SaENET	SLASSO (w)	SaLASSO (w)	SENET (w)	SaENET (w)	GLASSO	GaLASSO
PBDE 153	0.083	—	0.092	—	0.086	—	0.095	—	0.075	—
PeCB	0.331	0.744	0.367	0.715	0.284	0.670	0.322	0.657	0.299	0.968
trans-chlordane	0.074	—	0.109	—	0.073	0.034	0.109	0.036	0.044	—
cis-nonachlor	0.247	0.576	0.303	0.558	0.253	0.615	0.309	0.601	0.198	0.191
PCB151	0.080	0.156	0.128	0.151	0.057	0.189	0.111	0.184	0.038	—
# selected	5	3	5	3	5	4	5	4	5	2
# removed	18	20	18	20	18	19	18	19	18	21

NOTE: Results are based on 50 imputed datasets. Only five of the 23 POPs are displayed because the other 18 POPs were not selected by any method.

Table 3. Data generation details for all simulation settings.

	Correlation structure	Signal structure	No. Signals (%)	Signals
Case 1	(X_1, X_2, X_3) as 0.9 (X_6, X_7, X_8) as 0.5 (X_{11}, X_{12}, X_{13}) as 0.3	Concentrated	5 (25)	$\beta_1 = 2, \beta_4 = 1.5, \beta_7 = 1.5,$ $\beta_{11} = 1, \beta_{14} = 1;$
Case 2	(X_1, X_2, X_3) as 0.9 (X_6, X_7, X_8) as 0.5 (X_{11}, X_{12}, X_{13}) as 0.3	Distributed	5 (25)	$\beta_1 = 2, \beta_2 = 1, \beta_4 = 2, \beta_7 = 1, \beta_{11} = 1$
Case 3	(X_1, \dots, X_6) as 0.9 (X_{11}, \dots, X_{16}) as 0.5 (X_{21}, \dots, X_{26}) as 0.3	Concentrated	10 (10)	$\beta_2 = 2, \beta_7 = 0.8, \beta_9 = 0.8, \beta_{12} = 0.5,$ $\beta_{17} = 1.5, \beta_{27} = 1, \beta_{37} = 0.8, \beta_{47} = 0.4,$ $\beta_{48} = 1, \beta_{49} = 1;$
Case 4	(X_1, \dots, X_6) as 0.9 (X_{11}, \dots, X_{16}) as 0.5 (X_{21}, \dots, X_{26}) as 0.3	Distributed	10 (10)	$\beta_1 = 1.2, \beta_2 = 0.8, \beta_3 = 0.4, \beta_4 = 0.4,$ $\beta_{12} = 1.2, \beta_{13} = 1, \beta_{17} = 1.2, \beta_{27} = 1,$ $\beta_{37} = 1, \beta_{47} = 1;$

NOTE: In the correlation structure column, covariates within the same parentheses have an exchangeable pairwise correlation structure with one another, but are independent from all other covariates. For Case 1 and Case 2, $n = 500, p = 20$, approximately 25% subjects are missing $\{X_1, \dots, X_5\}$, 35% are missing $\{X_6, \dots, X_{13}\}$, 45% are missing $\{X_{14}, \dots, X_{17}\}$, and 55% are missing $\{X_{18}, X_{19}\}$. For Case 3 and Case 4, $n = 1000, p = 100$, about 25% subjects are missing $\{X_1, \dots, X_{30}\}$, 35% are missing $\{X_{31}, \dots, X_{60}\}$, 45% are missing $\{X_{61}, \dots, X_{82}\}$, 55% are missing $\{X_{83}, \dots, X_{95}\}$, and 60% are missing $\{X_{96}, \dots, X_{99}\}$. The No. Signals (%) column refers to the number of covariates that have nonnull coefficients and the corresponding relative percentage.

the stacked objective functions with observational weights additionally select PBDE 153, trans-chlordane and PCB 151. Since all methods select PeCB and cis-nonachlor, we conclude that further studies should be conducted to assess the PeCB and cis-nonachlor neurotoxicity.

5. Simulation

We now evaluate the performance of the stacked methods and the grouped methods mentioned in Section 2, including SLASSO, SaLASSO, SENET, SaENET, SLASSO(w), SaLASSO(w), SENET(w), SaENET(w), GLASSO, and GaLASSO.

5.1. Simulation setting

We simulate 1000 datasets of size n containing outcome Y and p covariates X under 4 different cases. For Cases 1 and 2 we take $n = 500$ and $p = 20$, and for Cases 3 and 4 we take $n = 1000$ and $p = 100$. In all cases, covariates are generated from a multivariate normal distribution with zero mean and unit variance. The correlation structure of the covariates is block-diagonal, in order to mimic the correlation structure of the POPs. A more comprehensive breakdown of the correlation structure is detailed in Table 3.

Given X , we generate a binary Y from $\text{logit}(P(Y = 1|X)) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$. The true value of β is specified according to different simulation cases, where Cases 1 and 3 correspond to concentrated signals and Cases 2 and 4 correspond to distributed signals. Here, signals are concentrated if there is only one nonnull coefficient in a group, and signals are distributed if there is more than one nonnull coefficient in

a group. Regression coefficient magnitudes are set to fix the prevalence of $Y = 1$ at about 50% and maintain the Cox-Snell pseudo- R^2 at approximately 0.5. The four simulation settings are summarized in Table 3.

Missing values are generated under the Missing at Random (MAR) assumption (Little and Rubin 2019). In all cases the outcome and the last covariate X_p are fully observed and the missingness indicator M_j for covariate X_j is generated from the logistic regression model

$$\text{logit}(\Pr(M_j = 1)) = \alpha_0 + \alpha_{1j} X_p + \alpha_{2j} Y$$

where $R_j = 1$ indicates that covariate X_j is missing, and α_0, α_{1j} and α_{2j} are chosen to control the percentage of missingness for X_j . For Case 1 and Case 2, about 25% of subjects are missing $\{X_1, \dots, X_5\}$, 35% are missing $\{X_6, \dots, X_{13}\}$, 45% are missing $\{X_{14}, \dots, X_{17}\}$ and 55% are missing $\{X_{18}, X_{19}\}$. For Case 3 and Case 4, about 25% of subjects are missing $\{X_1, \dots, X_{30}\}$, 35% are missing $\{X_{31}, \dots, X_{60}\}$, 45% are missing $\{X_{61}, \dots, X_{82}\}$, 55% are missing $\{X_{83}, \dots, X_{95}\}$, and 60% are missing $\{X_{96}, \dots, X_{99}\}$. In total, less than 5% of subjects have complete data, and about 13% of subjects have more than 90% data for all covariates. R package *mice* (Van Buuren and Groothuis-Oudshoorn 2011) is used to multiply impute the missing data using predictive mean matching. Each simulated dataset is imputed 10 or 50 times. To obtain independent imputed values, we set the number of iterations to 30 in *mice*. The stacked and grouped methods are then applied to perform variable selection and parameter estimation.

5.2. Simulation Results

We evaluate simulation results in terms of the following four metrics. In the following definitions, T and F are the number

of nonnull and null coefficients in the data generating model, respectively, R is the total number of simulation runs, $\hat{\beta}_j^r$ is the coefficient estimate for β_j in the r th run. Since the estimates for the j th coefficient by GLASSO and GaLASSO are different across imputed datasets, the mean $\frac{1}{D} \sum_{d=1}^D \hat{\beta}_{d,j}^r$ is used to approximate $\hat{\beta}_j^r$.

- Sensitivity (SENS) = $\frac{1}{RT} \sum_{r=1}^R \left(\# \text{ of selected nonnull coefficients in the } r_{\text{th}} \text{ run} \right)$
- Specificity (SPEC) = $\frac{1}{RF} \sum_{r=1}^R \left(\# \text{ of selected null coefficients in the } r_{\text{th}} \text{ run} \right)$
- Mean squared error for nonnull coefficients ($\text{MSE}_{\text{nonnull}}$) = $\frac{1}{R} \sum_{r=1}^R \sum_{j=1}^p (\hat{\beta}_j^r - \beta_j)^2 I(\beta_j \neq 0)$
- Mean squared error for null coefficients (MSE_{null}) = $\frac{1}{R} \sum_{r=1}^R \sum_{j=1}^p (\hat{\beta}_j^r - \beta_j)^2 I(\beta_j = 0)$

Sensitivity and specificity capture the accuracy of variable selection and vary between 0 and 1, with larger values indicating better performance. The MSEs quantify the estimation accuracy with smaller values being preferred.

Figures 3 and 4 and Table S4 present sensitivity, specificity, MSE for nonnull coefficients, and MSE for null coefficients for all four cases with 50 imputed datasets. Using SLASSO as the benchmark, the ratio to SLASSO for each measure is presented in Table S5. Overall, compared to the corresponding nonadaptive methods, the adaptive methods perform better with respect to estimation and selection. Specifically, the adaptive methods have similar sensitivity but considerably higher specificity, and considerably smaller MSE for nonnull coefficients, except for

GaLASSO under Case 4. For null coefficients, the adaptive methods have similar MSE to those of nonadaptive methods under Cases 1 and 2, and noticeably larger MSE under Cases 3 and 4, except for GaLASSO under Cases 3 and 4. Adaptive methods have high sensitivity across all cases and high specificity under Cases 1 and 2. Under Cases 3 and 4, GaLASSO has the highest specificity, followed by SaLASSO(w) and SaENET(w). The stacked methods with adaptive weights have smaller MSE for nonnull coefficients than the grouped method with adaptive weights under all four cases. On the other hand, under Case 4, GaLASSO has the largest MSE for nonnull coefficients and the smallest MSE for null coefficients, which is due to relatively low sensitivity and high specificity compared to other adaptive methods. In addition, the stacked approach performs better than the grouped approach in terms of estimation and selecting important variables. Compared to SLASSO, GLASSO has slightly lower sensitivity in all four cases, but has noticeably lower specificity in Cases 3 and 4. The MSE for nonnull coefficients is larger for GLASSO than SLASSO in all four cases, in addition to the MSE for null coefficients in Cases 3 and 4, due to large variability in the coefficient estimates.

The average runtime for each method is presented in Table 4. Compared to the grouped methods, the stacked methods are faster in all cases. SLASSO is 5.7 times faster than GLASSO in Case 1 when the sample size is small. When the sample size is large, that is, Case 3, SLASSO is 8.9 times faster than GLASSO. The adaptive methods have longer average runtimes than the nonadaptive methods because the adaptive methods require an elastic net initialization step to construct the adaptive weights.

To understand how the performance of these methods is affected by the number of imputed datasets, we provide simulation results for a smaller number of 10 imputed datasets in the supplementary materials and compare sensitivity, specificity, MSE for nonnull coefficients, and MSE for null coefficients (Fig-

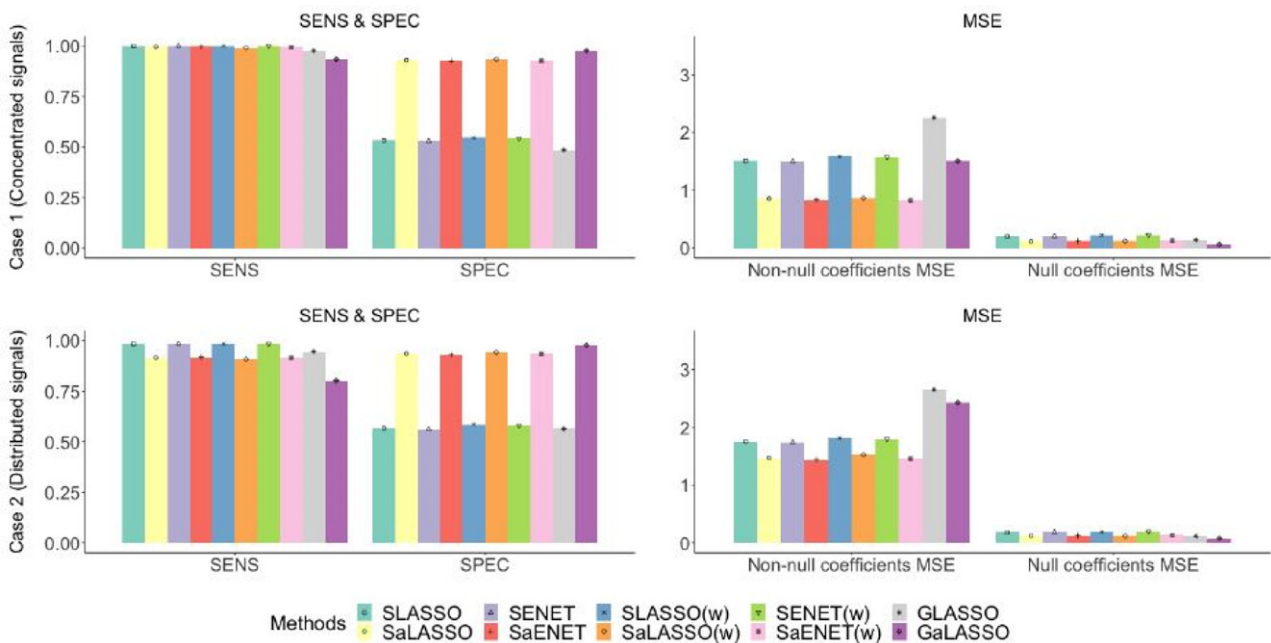


Figure 3. Simulation results for Case 1 (top panel) and Case 2 (bottom panel) where $n = 500$ and $p = 20$ for 50 imputed datasets. Sensitivity (SENS) and specificity (SPEC) are on the left and mean squared error (MSE) for nonnull and null coefficients are on the right.

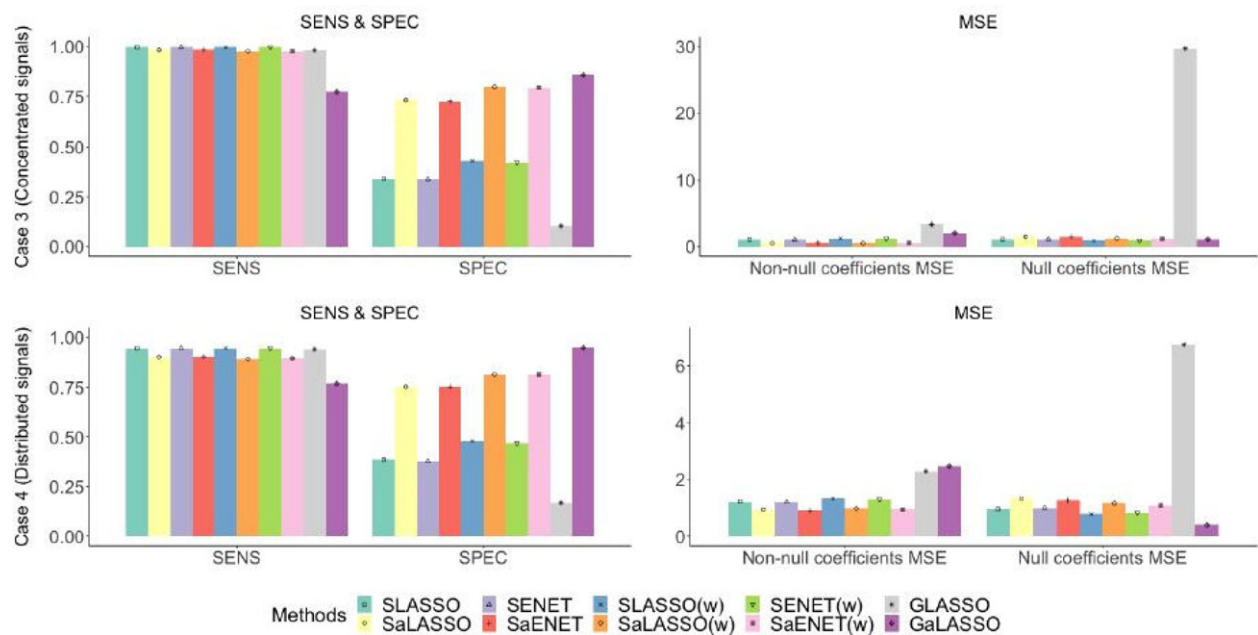


Figure 4. Simulation results for Case 3 (top panel) and Case 4 (bottom panel) where $n = 1000$ and $p = 100$ for 50 imputed datasets. Sensitivity (SENS) and specificity (SPEC) are on the left and mean squared error (MSE) for nonnull and null coefficients are on the right.

Table 4. Average runtime (in minutes) of each method for four simulation cases with 50 imputed datasets (1000 replications).

	SLASSO	SaLASSO	SENET	SaENET	SLASSO (w)	SaLASSO (w)	SENET (w)	SaENET (w)	GLASSO	GaLASSO
Case 1	12.8	56.7	51.3	77.1	12.7	52.4	47.7	69.1	73.1	86.1
Case 2	15.6	65.7	57.8	94.4	16.3	62.3	55.1	85.8	106.0	126.7
Case 3	127.0	585.0	557.8	704.5	129.6	549.2	527.9	638.2	1129.1	1182.8
Case 4	219.9	830.5	769.3	1086.5	194.7	736.4	686.3	931.4	1417.5	1467.0

NOTE: Case 1 and Case 2 have 500 observations and 20 covariates. Case 3 and Case 4 have 1000 observations and 100 covariates.

ures S2 and S3, Tables S6 and S7). Overall, increasing the number of imputed datasets improves the MSE for both nonnull and null coefficients, but has little impact on sensitivity and specificity, especially when p is large. The improvement coming from an increased number of imputed datasets is more significant for the grouped methods than the stacked methods. For example, in Case 3, the MSE for nonnull and null coefficients for SLASSO with 10 imputed datasets are 1.07 and 1.22, respectively, and with 50 imputed datasets are 1.04 and 1.08, respectively. While for GLASSO, the MSE for nonnull and null coefficients with 10 imputed datasets are 4.09 and 42.27, respectively, and with 50 imputed datasets are 3.30 and 29.70, respectively.

6. Discussion

In this article, we elucidated the difference between stacked and grouped pooled objective functions, which are both designed to achieve uniform variable selection across multiple imputed datasets. The stacked pooled objective function assumes that the underlying true signals are the same across imputed datasets, including the signal magnitude, while the grouped pooled objective function assumes uniform signal selection but allows for different active signal magnitudes across imputed datasets. We extended existing methods to handle binary outcomes, developed a MM algorithm combined with block coordinate descent updates to optimize grouped pooled objective functions for LASSO and aLASSO regularization, and derived cyclic

coordinate descent algorithm for the stacked pooled objective functions with ENET and aENET regularization. Algorithms for implementing the stacked and grouped approaches outlined in Section 2, for both continuous and binary outcomes, are available in the *miselect* R package on CRAN.

From a practical perspective, there are several reasons that one might prefer stacked over grouped objective functions. Based on our simulations, the overall MSE for estimates generated by optimizing stacked objective functions was either smaller than or equal to the estimates generated by optimizing the grouped objective function, provided that adaptive weighting is used. We also observed that the total runtime for optimizing stacked pooled objective functions is noticeably lower compared to the grouped objective function optimization routine. An additional practical consideration is that the stacked objective functions are much easier to extend beyond ENET penalization. For example, if one wanted to use a hierarchical interaction detection penalty, such as *hierNet* (Bien et al. 2013), one would only need to download the *hierNet* package from CRAN, and use the existing *hierNet* implementation on the stacked imputed datasets. Conversely, grouped methods would necessitate the development of additional algorithms to optimize an objective function with both a group LASSO penalty and a hierarchical interaction detection penalty. Finally, although we did not observe a substantial difference between equal weights and the observational weights proposed by Wan et al. (2015), there are still conceptual concerns with having the

observational weights depend on the fraction of missingness. As we mentioned earlier, upweighting observations with more data artificially moves the analysis in the direction of complete-case analysis, which is known to be biased under MAR and MNAR missing data mechanisms.

The positive ALS-POPs associations identified in the data example add to a growing body of literature on environmental risk factors for ALS (McGuire et al. 1997; Sutedja et al. 2008; Kamel et al. 2012). A major advantage of the data collected in this study is that POP concentrations were measured in plasma samples, rather than through surveys. Since ALS results from the complex interplay of multiple risks combined with neurotoxic environmental exposures (the gene-time-environment hypothesis), we contend that additional work is needed to more fully understand gene and pesticide exposure interaction (Al-Chalabi and Hardiman 2013; Paez-Colasante et al. 2015; Goutman et al. 2018).

Supplementary Materials

A summary table describing existing methods for handling missing data and variable selection in the literature. A table describing the proportion of missing values in each variable in the data example. A figure of pairwise correlation for the 23 exposures considered in the data example. A version of Table 2 (in the main text) with 10 imputed datasets. Tables and figures for additional simulation results with 10 or 50 imputed datasets. Example codes for using stacked and grouped methods in R package miselect.

Acknowledgments

The authors thank Alexander Rix, BS for software development support, Crystal Pacut, Jayna Duell, RN, Blake Swihart, and Daniel Burger for study support.

Disclosure Statement

Dr. Stephen A. Goutman: scientific advisory to Biogen and ITF Pharma. DSMB service for Watermark Research Partners; additional research support from ALS Association.

Funding

This work was supported by the NSF DMS under Grant 1712933; NIH under Grant P30 CA 046591; NIEHS under Grant K23ES027221; National ALS Registry/CDC/ATSDR CDCP-DHHS-US (CDC/ATSDR under Grant 200-2013-56856); NIH/NINDS under Grant R01NS127188; NIH/NIEHS under Grant R01ES030049; NeuroNetwork for Emerging Therapies at Michigan Medicine, University of Michigan; Robert and Katherine Jacobs Environmental Health Initiative; Sinai Medical Staff Foundation and Scott L. Pranger.

ORCID

Lauren J. Beesley  <https://orcid.org/0000-0002-3788-5944>

References

Al-Chalabi, A., and Hardiman, O. (2013), "The Epidemiology of ALS: A Conspiracy of Genes, Environment and Time," *Nature Reviews Neurology*, 9, 617–628. [1069,1074]

- Bien, J., Taylor, J., and Tibshirani, R. (2013), "A Lasso for Hierarchical Interactions," *The Annals of Statistics*, 41, 1111–1141. [1073]
- Boss, J., Mukherjee, B., Ferguson, K. K., Aker, A., Alshawabkeh, A. N., Cordero, J. F., Meeker, J. D., and Kim, S. (2019), "Estimating Outcome-Exposure Associations When Exposure Biomarker Detection Limits Vary Across Batches," *Epidemiology*, 30, 746–755. [1070]
- Breheny, P., and Huang, J. (2015), "Group Descent Algorithms for Non-convex Penalized Linear and Logistic Regression Models with Grouped Predictors," *Statistics and Computing*, 25, 173–187. [1068]
- Chen, Q., and Wang, S. (2013), "Variable Selection for Multiply-Imputed Data with Application to Dioxin Exposure Study," *Statistics in Medicine*, 32, 3646–3659. [1064,1066,1067,1068]
- Chio, A., Logroscino, G., Hardiman, O., Swingler, R., Mitchell, D., Beghi, E., Traynor, B. G., Consortium, E. (2009), "Prognostic Factors in ALS: A Critical Review," *Amyotrophic Lateral Sclerosis*, 10, 310–323. [1070]
- Friedman, J., Hastie, T., and Tibshirani, R. (2010), "Regularization Paths for Generalized Linear Models via Coordinate Descent," *Journal of Statistical Software*, 33, 1–22. [1066]
- Garcia, R. I., Ibrahim, J. G., and Zhu, H. (2010), "Variable Selection in the Cox Regression Model with Covariates Missing at Random," *Biometrics*, 66, 97–104. [1063]
- Geronimi, J., and Saporta, G. (2017), "Variable Selection for Multiply-Imputed Data with Penalized Generalized Estimating Equations," *Computational Statistics & Data Analysis*, 110, 103–114. [1066]
- Ghosh, D., Zhu, Y., and Coffman, D. L. (2015), "Penalized Regression Procedures for Variable Selection in the Potential Outcomes Framework," *Statistics in Medicine*, 34, 1645–1658. [1063]
- Goutman, S. A. (2017), "Diagnosis and Clinical Management of Amyotrophic Lateral Sclerosis and Other Motor Neuron Disorders," *CONTINUUM: Lifelong Learning in Neurology*, 23, 1332–1359. [1069]
- Goutman, S. A., J. Boss, A. Patterson, B. Mukherjee, S. Batterman, and E. L. Feldman (2019), "High Plasma Concentrations of Organic Pollutants Negatively Impact Survival in Amyotrophic Lateral Sclerosis," *Journal of Neurology, Neurosurgery & Psychiatry*, 90, 907–912. [1069,1070]
- Goutman, S. A., Chen, K. S., Paez-Colasante, X., and Feldman, E. L. (2018), "Emerging Understanding of the Genotype-Phenotype Relationship in Amyotrophic Lateral Sclerosis," in *Handbook of Clinical Neurology* (Vol. 148), eds. D. H. Geschwind, H. L. Paulson, and C. Klein, pp. 603–623, Amsterdam: Elsevier. [1069,1074]
- Hastie, T., and Qian, J. (2014), "Glmnet vignette." Available at http://www.web.stanford.edu/hastie/Papers/Glmnet_Vignette.pdf. [1064]
- Hastie, T., Tibshirani, R., and Friedman, J. (2009), *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, New York: Springer. [1069]
- Hastie, T., Tibshirani, R., and Wainwright, M. (2015), *Statistical Learning with Sparsity: The Lasso and Generalizations*, Boca Raton, FL: Chapman and Hall/CRC. [1069]
- Heymans, M. W., Van Buuren, S., Knol, D. L., Van Mechelen, W., and De Vet, H. C. (2007), "Variable Selection Under Multiple Imputation Using the Bootstrap in a Prognostic Study," *BMC Medical Research Methodology*, 7, 1–10. [1063]
- Johnson, B. A., Lin, D., and Zeng, D. (2008), "Penalized Estimating Functions and Variable Selection in Semiparametric Regression Models," *Journal of the American Statistical Association*, 103, 672–680. [1063]
- Kamel, F., Umbach, D. M., Bedlack, R. S., Richards, M., Watson, M., Alavanja, M. C. R., Blair, A., Hoppin, J. A., Schmidt, S., and Sandler, D. P. (2012), "Pesticide Exposure and Amyotrophic Lateral Sclerosis," *Neurotoxicology*, 33, 457–462. [1074]
- Lachenbruch, P. A. (2011), "Variable selection When Missing Values are Present: A Case Study," *Statistical Methods in Medical Research*, 20, 429–444. [1063]
- Little, R. J., and Rubin, D. B. (2019), *Statistical Analysis with Missing Data* (Vol. 793), Hoboken, NJ: Wiley. [1071]
- Liu, Y., Wang, Y., Feng, Y., and Wall, M. M. (2016), "Variable Selection and Prediction with Incomplete High-Dimensional Data," *Annals of Applied Statistics*, 10, 418–450. [1063]
- Long, Q., and Johnson, B. A. (2015), "Variable Selection in the Presence of Missing Data: Resampling and Imputation," *Biostatistics*, 16, 596–610. [1063]

- McGuire, V., Longstreth Jr., W. T., Nelson, L. M., Koepsell, T. D., Checkoway, H., Morgan, M. S., and van Belle, G. (1997), "Occupational Exposures and Amyotrophic Lateral Sclerosis. A Population-Based Case-Control Study," *American Journal of Epidemiology*, 145, 1076–1088. [1074]
- Paez-Colasante, X., Figueroa-Romero, C., Sakowski, S. A., Goutman, S. A., and Feldman, E. L. (2015), "Amyotrophic Lateral Sclerosis: Mechanisms and Therapeutics in the Epigenomic Era," *Nature Reviews Neurology*, 11, 266. [1069,1074]
- Robins, J. M., and Wang, N. (2000), "Inference for Imputation Estimators," *Biometrika*, 87, 113–124. [1065]
- Schafer, J. L. (2001), "Analyzing the NHANES III Multiply Imputed Data Set: Methods and Examples," Prepared for the National Center for Health, Statistics, 1–61. [1063]
- Su, F.-C., Goutman, S. A., Chernyak, S., Mukherjee, B., Callaghan, B. C., Batterman, S., and Feldman, E. L. (2016), "Association of Environmental Toxins with Amyotrophic Lateral Sclerosis," *JAMA Neurology*, 73, 803–811. [1069,1070]
- Sutedja, N. A., Veldink, J. H., Fischer, K., Kromhout, H., Heederik, D., Huisman, M. H., Wokke, J. H. J., and van den Berg, L. H. (2008), "Exposure to Chemicals and Metals and Risk of Amyotrophic Lateral Sclerosis: A Systematic Review," *Amyotrophic Lateral Sclerosis*, 10, 302–309. [1074]
- Tibshirani, R. (1996), "Regression Shrinkage and Selection via the Lasso," *Journal of the Royal Statistical Society, Series B*, 58, 267–288. [1064]
- Van Buuren, S., and Groothuis-Oudshoorn, K. (2011), "mice: Multivariate Imputation by Chained Equations in r," *Journal of Statistical Software*, 45, 1–67. [1070,1071]
- Wan, Y., Datta, S., Conklin, D. J., and Kong, M. (2015), "Variable Selection Models Based on Multiple Imputation with an Application for Predicting Median Effective Dose and Maximum Effect," *Journal of Statistical Computation and Simulation*, 85, 1902–1916. [1064,1065,1073]
- Wang, N., and Robins, J. M. (1998), "Large-Sample Theory for Parametric Multiple Imputation Procedures," *Biometrika*, 85, 935–948. [1065]
- Wood, A. M., White, I. R., and Royston, P. (2008), "How Should Variable Selection be Performed with Multiply Imputed Data?" *Statistics in Medicine*, 27, 3227–3246. [1063,1070]
- Yang, X., Belin, T. R., and Boscardin, W. J. (2005), "Imputation and Variable Selection in Linear Regression Models with Missing Covariates," *Biometrics*, 61, 498–506. [1063]
- Yang, Y., and Zou, H. (2015), "A Fast Unified Algorithm for Solving Group-Lasso Penalize Learning Problems," *Statistics and Computing*, 25, 1129–1141. [1064]
- Yang, Y., Zou, H., Yang, M. Y., and Matrix, D. (2017), "Package 'gcdnet'" [1064]
- Yu, Y., Su, F.-C., Callaghan, B. C., Goutman, S. A., Batterman, S. A., and Feldman, E. L. (2014), "Environmental Risk Factors and Amyotrophic Lateral Sclerosis (ALS): A Case-Control Study of ALS in Michigan," *PloS One*, 9, e101186. [1069]
- Yuan, M., and Lin, Y. (2006), "Model Selection and Estimation in Regression with Grouped Variables," *Journal of the Royal Statistical Society, Series B*, 68, 49–67. [1066]
- Zou, H. (2006), "The Adaptive Lasso and Its Oracle Properties," *Journal of the American Statistical Association*, 101, 1418–1429. [1064,1065]
- Zou, H., and Hastie, T. (2005), "Regularization and Variable Selection via the Elastic Net," *Journal of the Royal Statistical Society, Series B*, 67, 301–320. [1064]
- Zou, H., and Zhang, H. H. (2009), "On the Adaptive Elastic-Net with a Diverging Number of Parameters," *The Annals of Statistics*, 37, 1733–1751. [1064,1065,1066]

Supplementary materials for *Variable selection with multiply-imputed datasets: choosing between stacked and grouped methods*

Jiacong Du¹, Jonathan Boss¹, Peisong Han¹, Lauren J Beesley¹, Michael Kleinsasser¹
Stephen A Goutman², Stuart Batterman³, Eva L Feldman², Bhramar Mukherjee¹

¹Department of Biostatistics, University of Michigan, Ann Arbor, MI

²Department of Neurology, University of Michigan, Ann Arbor, MI

³Department of Environmental Health Science, University of Michigan, Ann Arbor, MI

Corresponding author Jiacong Du jiacong@umich.edu

January 3, 2022

Table S1: Summary of methods for handling missing data and variable selection in the literature.

Class	Method	Key reference First author, year	Method description	Limitations
Sequential handling missing data by multiple imputation and variable selection	Ad-hoc thresholding	Wood et al. (2008)	First apply the existing variable selection method, e.g., stepwise selection, on each imputed dataset. A variable is considered selected if it is selected in at least πD imputed datasets, $0 < \pi \leq 1$.	No clear guideline on how to choose the thresholding parameter π .
	Bootstrapped imputation and variable selection	Long and Johnson (2015)	This method imputes missing values in the bootstrapped samples from the original data and applies penalized regression on each bootstrapped imputed dataset. Proportions of each covariate selected over all bootstrap imputed datasets are provided and the final active set is determined by a thresholding.	1. Coefficient estimates can only be obtained after the final active set is determined. 2. Still requires ad-hoc thresholding.
	Elastic net regularization on the stacked imputed datasets	Wan et al. (2015)	Simultaneously performs variable selection and regression coefficient estimation by maximizing the data likelihood after stacking all imputed datasets subject to elastic net penalization. Each data point is weighted by an observational weight which quantifies the proportion of observed covariates out of the total number of covariates.	1. Only continuous outcomes are considered in the paper. 2. Using the proposed observational weights can introduce bias if the missing data mechanism is dependent on the outcome.
	Group LASSO regularization on the pooled likelihood for imputed datasets	Chen and Wang (2013)	Group LASSO penalized optimization over all imputed datasets, where the groups are defined by regression coefficients corresponding to the same variable across all imputed datasets. Variable selection and point estimation are obtained simultaneously.	1. Existing optimization algorithm uses local quadratic approximation, which is reliant on a subjective hard thresholding rule to select variables. 2. The existing implementation is limited to continuous outcomes and a LASSO penalty function.
	Bayesian variable selection by data augmentation strategy	Yang et al. (2005)	The Simultaneously Impute and Select method jointly models the missing data mechanism, model parameters, and selection indicators in a Bayesian framework, and samples joint posterior draws using Gibbs sampler.	Can be computationally intensive for a large amount of missingness.
Simultaneously handling missing data and variable selection	Penalized maximum likelihood estimators combined with EM algorithm	Garcia et al. (2010)	Variable selection and coefficient estimates are simultaneously obtained by maximizing the penalized observed data likelihood. SCAD and ALASSO penalty functions have been studied. Additional information criteria involving penalized observed data likelihood have been developed for choosing tuning parameters.	Computational challenges for integrating the augmented likelihood with no closed forms over the missing data distribution to find the observed data likelihood.
	Penalized estimating equations combined with IPW/AIPW	Johnson et al. (2008)	Missing data is handled by IPW for each complete-observed data point, and variable selection is achieved by solving the penalized estimating equations.	Assumes that the missing data pattern is monotone.

Table S2: Missing data proportion for 30 variables in the ALS data. The data in total contains 266 observations with 167 cases and 99 controls.

Variables	Proportion	Variables	Proportion
PCB 174	0.508	PCB 153	0.237
PCB 110	0.429	PCB 202	0.199
cis-chlordane	0.398	cis-nonachlor	0.147
trans-nonachlor	0.380	beta-HCH	0.132
PCB 175	0.368	PBDE 99	0.086
trans-chlordane	0.365	p,p'-DDE	0.083
PCB 118	0.361	PeCB	0.060
PCB 180	0.350	PBDE 100	0.056
PBDE 28	0.305	Education1	0.038
PCB 138	0.278	Education2	0.038
PBDE 154	0.267	PCB 151	0.034
PBDE 153	0.256	PBDE 47	0.015
BMI	0.244	Age	0.011
BMI_slope	0.244	Sex	0.000
PBDE 85	0.241	ALS	0.000

Pearson correlation plot of 23 environmental pollutants

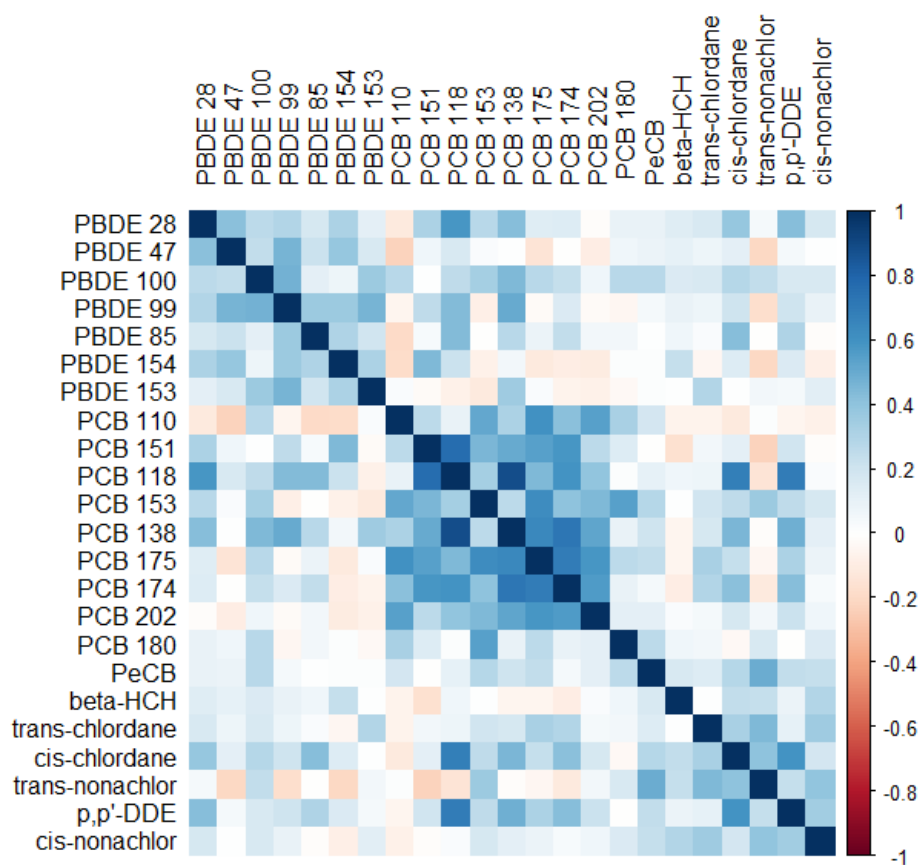


Figure S1: Pairwise Pearson correlation matrix for 23 POPs collected through the University of Michigan ALS Patient Biorepository. The dataset contains 266 observations (167 cases and 99 controls).

Table S3: Regression coefficient estimates for five POPs collected as part of the University of Michigan ALS Patients Biorepository case-control study (167 ALS cases and 99 healthy controls). Results are based on 10 imputed datasets. Only five of the 23 POPs are displayed because the other 18 POPs were not selected by any method.

POPs	SLASSO	SaLASSO	SENET	SaENET	SLASSO (w)	SaLASSO (w)	SENET (w)	SaENET (w)	GLASSO	GaLASSO
PBDE 153	0.087	-	0.095	-	0.090	0.075	0.099	0.067	0.081	-
PeCB	0.316	0.751	0.354	0.743	0.268	0.645	0.307	0.655	0.295	0.969
trans-chlordane	0.067	0.001	0.104	0.003	0.068	0.081	0.105	0.068	0.047	-
cis-nonachlor	0.220	0.578	0.278	0.572	0.225	0.524	0.284	0.530	0.185	0.194
PCB151	0.060	0.173	0.109	0.171	0.035	-	0.090	-	0.031	-
# selected	5	4	5	4	5	4	5	4	5	2
# removed	18	19	18	19	18	19	18	19	18	21

Table S4: Sensitivity (SENS), specificity (SPEC), MSE for non-null coefficients ($MSE_{non-null}$), and MSE for null coefficients (MSE_{null}) for simulation Cases 1-4 with 50 imputed datasets. Cases 1 and 2 have 500 observations with 20 covariates, while Cases 3 and 4 have 1000 observations with 100 covariates. Cases 1 and 3 correspond to concentrated signals and Cases 2 and 4 correspond to distributed signals.

	SLASSO	SaLASSO	SENET	SaENET	SLASSO (w)	SaLASSO (w)	SENET (w)	SaENET (w)	GLASSO	GaLASSO
Case 1										
<i>SENS</i>	1.00	1.00	1.00	1.00	1.00	0.99	1.00	0.99	0.98	0.94
<i>SPEC</i>	0.53	0.93	0.53	0.92	0.55	0.93	0.54	0.93	0.49	0.98
$MSE_{non-null}$	1.51	0.86	1.50	0.84	1.58	0.86	1.57	0.83	2.26	1.51
MSE_{null}	0.20	0.12	0.20	0.12	0.22	0.12	0.22	0.13	0.14	0.06
Case 2										
<i>SENS</i>	0.98	0.92	0.98	0.92	0.98	0.91	0.98	0.92	0.95	0.80
<i>SPEC</i>	0.57	0.94	0.56	0.93	0.59	0.94	0.58	0.94	0.56	0.98
$MSE_{non-null}$	1.75	1.47	1.74	1.44	1.81	1.53	1.79	1.46	2.66	2.43
MSE_{null}	0.18	0.12	0.18	0.13	0.19	0.12	0.20	0.13	0.12	0.08
Case 3										
<i>SENS</i>	1.00	0.98	1.00	0.98	1.00	0.98	1.00	0.98	0.98	0.77
<i>SPEC</i>	0.34	0.73	0.34	0.72	0.43	0.80	0.42	0.79	0.10	0.86
$MSE_{non-null}$	1.04	0.52	1.03	0.55	1.21	0.51	1.19	0.55	3.30	2.02
MSE_{null}	1.08	1.47	1.11	1.44	0.87	1.23	0.91	1.18	29.70	1.09
Case 4										
<i>SENS</i>	0.95	0.90	0.95	0.90	0.95	0.89	0.95	0.90	0.94	0.77
<i>SPEC</i>	0.39	0.75	0.38	0.75	0.48	0.81	0.47	0.81	0.17	0.95
$MSE_{non-null}$	1.21	0.93	1.19	0.90	1.31	0.97	1.29	0.93	2.28	2.46
MSE_{null}	0.95	1.31	0.99	1.25	0.78	1.16	0.82	1.08	6.75	0.39

Table S5: Sensitivity (SENS), specificity (SPEC), MSE for non-null coefficients ($MSE_{non-null}$), and MSE for null coefficients (MSE_{null}) for simulation Cases 1-4 with 50 imputed datasets. Cases 1 and 2 have 500 observations with 20 covariates, while Cases 3 and 4 have 1000 observations with 100 covariates. Cases 1 and 3 correspond to concentrated signals and Cases 2 and 4 correspond to distributed signals. For each measure, SLASSO is used as the benchmark and the ratio to SLASSO is presented.

	SLASSO	SaLASSO	SENET	SaENET	SLASSO (w)	SaLASSO (w)	SENET (w)	SaENET (w)	GLASSO	GaLASSO
Case 1										
<i>SENS</i>	1.00	1.00	1.00	1.00	1.00	0.99	1.00	0.99	0.98	0.94
<i>SPEC</i>	1.00	1.74	0.99	1.74	1.02	1.75	1.02	1.74	0.91	1.83
<i>MSE_{non-null}</i>	1.00	0.57	0.99	0.56	1.05	0.57	1.04	0.55	1.50	1.00
<i>MSE_{null}</i>	1.00	0.58	1.01	0.61	1.10	0.60	1.11	0.66	0.70	0.32
Case 2										
<i>SENS</i>	1.00	0.93	1.00	0.93	1.00	0.92	1.00	0.93	0.96	0.82
<i>SPEC</i>	1.00	1.65	0.99	1.64	1.03	1.66	1.02	1.65	0.99	1.72
<i>MSE_{non-null}</i>	1.00	0.84	0.99	0.82	1.04	0.87	1.02	0.83	1.52	1.39
<i>MSE_{null}</i>	1.00	0.69	1.02	0.71	1.06	0.67	1.08	0.74	0.66	0.45
Case 3										
<i>SENS</i>	1.00	0.99	1.00	0.99	1.00	0.98	1.00	0.98	0.98	0.77
<i>SPEC</i>	1.00	2.14	0.99	2.12	1.26	2.33	1.24	2.32	0.30	2.51
<i>MSE_{non-null}</i>	1.00	0.50	0.99	0.53	1.16	0.48	1.14	0.52	3.16	1.94
<i>MSE_{null}</i>	1.00	1.36	1.02	1.33	0.81	1.13	0.84	1.09	27.42	1.01
Case 4										
<i>SENS</i>	1.00	0.95	1.00	0.96	1.00	0.94	1.00	0.95	1.00	0.81
<i>SPEC</i>	1.00	1.95	0.98	1.95	1.24	2.11	1.21	2.11	0.43	2.46
<i>MSE_{non-null}</i>	1.00	0.77	0.99	0.74	1.08	0.80	1.06	0.77	1.89	2.03
<i>MSE_{null}</i>	1.00	1.38	1.04	1.32	0.82	1.22	0.87	1.13	7.11	0.42

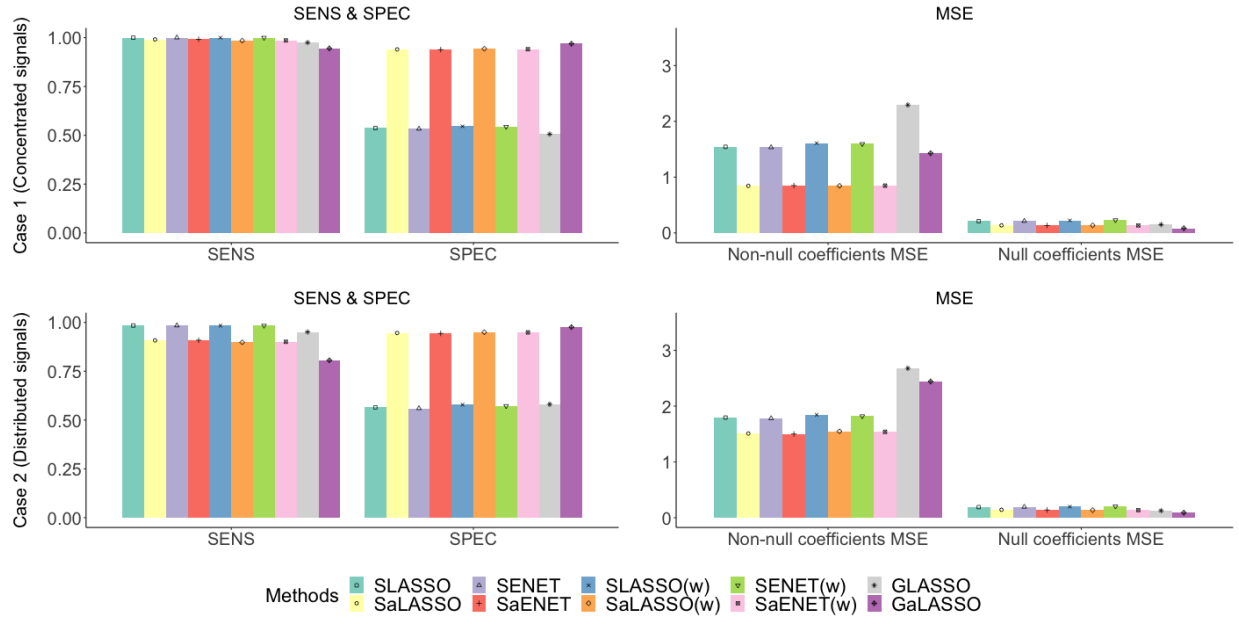


Figure S2: Simulation results for Case 1 (top panel) and Case 2 (bottom panel) where $n=500$ and $p=20$ for 10 imputed datasets. Sensitivity (SENS) and specificity (SPEC) are on the left and MSE for non-null and null coefficients are on the right.

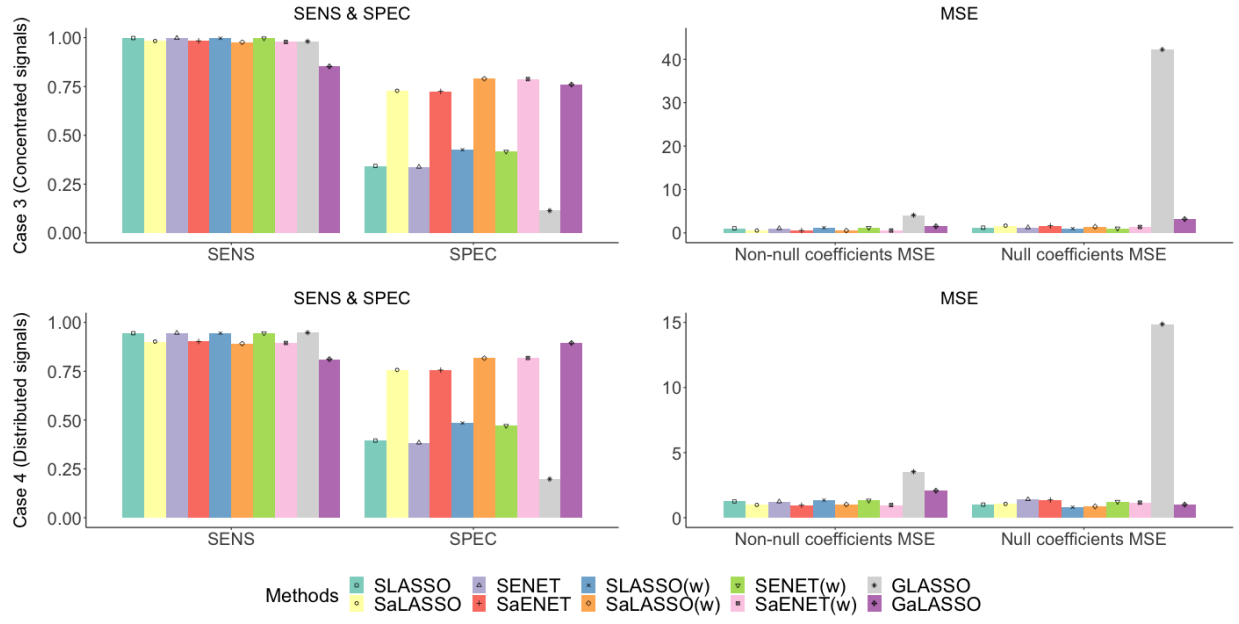


Figure S3: Simulation results for Case 3 (top panel) and Case 4 (bottom panel) where $n=1000$ and $p=100$ for 10 imputed datasets. Sensitivity (SENS) and specificity (SPEC) are on the left and MSE for non-null and null coefficients are on the right.

Table S6: Sensitivity (SENS), specificity (SPEC), MSE for non-null coefficients ($MSE_{non-null}$), and MSE for null coefficients (MSE_{null}) for simulation Cases 1-4 with 10 imputed datasets. Cases 1 and 2 have 500 observations with 20 covariates, while Cases 3 and 4 have 1000 observations with 100 covariates. Cases 1 and 3 correspond to concentrated signals and Cases 2 and 4 correspond to distributed signals.

	SLASSO	SaLASSO	SENET	SaENET	SLASSO (w)	SaLASSO (w)	SENET (w)	SaENET (w)	GLASSO	GaLASSO
Case 1										
<i>SENS</i>	1.00	0.99	1.00	0.99	1.00	0.98	1.00	0.99	0.98	0.94
<i>SPEC</i>	0.54	0.94	0.53	0.94	0.55	0.94	0.54	0.94	0.51	0.97
$MSE_{non-null}$	1.54	0.84	1.53	0.85	1.61	0.84	1.60	0.85	2.29	1.43
MSE_{null}	0.21	0.14	0.21	0.13	0.22	0.14	0.23	0.13	0.15	0.08
Case 2										
<i>SENS</i>	0.98	0.91	0.98	0.91	0.98	0.90	0.98	0.90	0.95	0.81
<i>SPEC</i>	0.56	0.95	0.56	0.94	0.58	0.95	0.57	0.95	0.58	0.98
$MSE_{non-null}$	1.79	1.51	1.78	1.50	1.84	1.55	1.82	1.54	2.68	2.44
MSE_{null}	0.19	0.14	0.19	0.14	0.20	0.14	0.20	0.14	0.13	0.09
Case 3										
<i>SENS</i>	1.00	0.98	1.00	0.98	1.00	0.98	1.00	0.98	0.98	0.85
<i>SPEC</i>	0.34	0.73	0.34	0.72	0.43	0.79	0.42	0.79	0.11	0.76
$MSE_{non-null}$	1.07	0.55	1.05	0.56	1.22	0.54	1.20	0.55	4.09	1.54
MSE_{null}	1.22	1.68	1.26	1.64	0.99	1.41	1.05	1.37	42.27	3.19
Case 4										
<i>SENS</i>	0.94	0.90	0.95	0.90	0.94	0.89	0.94	0.90	0.95	0.81
<i>SPEC</i>	0.40	0.76	0.38	0.76	0.48	0.82	0.47	0.82	0.20	0.90
$MSE_{non-null}$	1.26	0.98	1.24	0.96	1.35	1.02	1.32	0.97	3.53	2.08
MSE_{null}	1.00	1.42	1.06	1.36	0.80	1.23	0.87	1.16	14.87	1.01

Table S7: Sensitivity (SENS), specificity (SPEC), MSE for non-null coefficients ($MSE_{non-null}$), and MSE for null coefficients (MSE_{null}) for simulation Cases 1-4 with 10 imputed datasets. Cases 1 and 2 have 500 observations with 20 covariates, while Cases 3 and 4 have 1000 observations with 100 covariates. Cases 1 and 3 correspond to concentrated signals and Cases 2 and 4 correspond to distributed signals. For each measure, SLASSO is used as the benchmark and the ratio to SLASSO is presented.

	SLASSO	SaLASSO	SENET	SaENET	SLASSO (w)	SaLASSO (w)	SENET (w)	SaENET (w)	GLASSO	GaLASSO
Case 1										
<i>SENS</i>	1.00	0.99	1.00	0.99	1.00	0.98	1.00	0.99	0.98	0.94
<i>SPEC</i>	1.00	1.75	0.99	1.75	1.02	1.76	1.01	1.75	0.94	1.80
<i>MSE_{non-null}</i>	1.00	0.55	0.99	0.55	1.04	0.55	1.03	0.55	1.49	0.93
<i>MSE_{null}</i>	1.00	0.64	1.01	0.64	1.07	0.65	1.10	0.64	0.71	0.40
Case 2										
<i>SENS</i>	1.00	0.92	1.00	0.92	1.00	0.91	1.00	0.92	0.97	0.82
<i>SPEC</i>	1.00	1.67	0.99	1.67	1.02	1.68	1.01	1.68	1.03	1.73
<i>MSE_{non-null}</i>	1.00	0.84	0.99	0.84	1.03	0.86	1.01	0.86	1.49	1.36
<i>MSE_{null}</i>	1.00	0.75	1.02	0.73	1.05	0.71	1.07	0.72	0.66	0.48
Case 3										
<i>SENS</i>	1.00	0.98	1.00	0.99	1.00	0.98	1.00	0.98	0.98	0.85
<i>SPEC</i>	1.00	2.12	0.99	2.11	1.24	2.30	1.22	2.30	0.33	2.22
<i>MSE_{nonnull}</i>	1.00	0.52	0.99	0.52	1.14	0.50	1.12	0.52	3.83	1.44
<i>MSE_{null}</i>	1.00	1.38	1.03	1.34	0.81	1.16	0.86	1.12	34.62	2.61
Case 4										
<i>SENS</i>	1.00	0.95	1.00	0.96	1.00	0.94	1.00	0.95	1.00	0.86
<i>SPEC</i>	1.00	1.92	0.97	1.91	1.23	2.07	1.19	2.07	0.50	2.27
<i>MSE_{nonnull}</i>	1.00	0.78	0.98	0.76	1.08	0.81	1.05	0.77	2.81	1.66
<i>MSE_{null}</i>	1.00	1.41	1.05	1.36	0.80	1.22	0.86	1.15	14.80	1.00

References

- Chen, Q. and S. Wang (2013). Variable selection for multiply-imputed data with application to dioxin exposure study. *Statistics in Medicine* 32(21), 3646–3659.
- Garcia, R. I., J. G. Ibrahim, and H. Zhu (2010). Variable selection in the cox regression model with covariates missing at random. *Biometrics* 66(1), 97–104.
- Johnson, B. A., D. Lin, and D. Zeng (2008). Penalized estimating functions and variable selection in semiparametric regression models. *Journal of the American Statistical Association* 103(482), 672–680.
- Long, Q. and B. A. Johnson (2015). Variable selection in the presence of missing data: resampling and imputation. *Biostatistics* 16(3), 596–610.
- Wan, Y., S. Datta, D. J. Conklin, and M. Kong (2015). Variable selection models based on multiple imputation with an application for predicting median effective dose and maximum effect. *Journal of Statistical Computation and Simulation* 85(9), 1902–1916.
- Wood, A. M., I. R. White, and P. Royston (2008). How should variable selection be performed with multiply imputed data? *Statistics in Medicine* 27(17), 3227–3246.
- Yang, X., T. R. Belin, and W. J. Boscardin (2005). Imputation and variable selection in linear regression models with missing covariates. *Biometrics* 61(2), 498–506.