

# A Sparse Additive Model for High-Dimensional Interactions with an Exposure Variable

Sahir R Bhatnagar<sup>a,b,\*</sup>, Tianyuan Lu<sup>c,d</sup>, Amanda Lovato<sup>e</sup>, David L Olds<sup>f</sup>, Michael S Kobor<sup>g</sup>, Michael J Meaney<sup>h</sup>, Kieran O'Donnell<sup>i</sup>, Yi Yang<sup>j</sup> and Celia MT Greenwood<sup>a,c,e</sup>

<sup>a</sup>Department of Epidemiology, Biostatistics and Occupational Health, McGill University, Montréal, Canada

<sup>b</sup>Department of Diagnostic Radiology, McGill University, Montréal, Canada

<sup>c</sup>Quantitative Life Sciences, McGill University

<sup>d</sup>Lady Davis Institute, Jewish General Hospital, Montréal, QC

<sup>e</sup>Statistics Canada, Ottawa, ON

<sup>f</sup>Department of Pediatrics, University of Colorado School of Medicine, Denver

<sup>g</sup>Department of Medical Genetics, University of British Columbia, BC

<sup>h</sup>Singapore Institute for Clinical Sciences, Singapore; McGill University

<sup>i</sup>Department of Psychiatry, McGill University

<sup>j</sup>Department of Mathematics and Statistics, McGill University

<sup>k</sup>Departments of Oncology and Human Genetics, McGill University

---

## ARTICLE INFO

### Keywords:

Gene-environment interaction

Strong heredity property

Blockwise coordinate descent

High-dimensional data

Variable selection

---

## ABSTRACT

A conceptual paradigm for onset of a new disease is often considered to be the result of changes in entire biological networks whose states are affected by a complex interaction of genetic and environmental factors. However, when modelling a relevant phenotype as a function of high dimensional measurements, power to estimate interactions is low, the number of possible interactions could be enormous and their effects may be non-linear. In this work, we introduce a method called `sail` for detecting non-linear interactions with a key environmental or exposure variable in high-dimensional settings which respects the strong or weak heredity constraints. We prove that asymptotically, our method possesses the oracle property, i.e., it performs as well as if the true model were known in advance. We develop a computationally efficient fitting algorithm with automatic tuning parameter selection, which scales to high-dimensional datasets. Through an extensive simulation study, we show that `sail` outperforms existing penalized regression methods in terms of prediction accuracy and support recovery when there are non-linear interactions with an exposure variable. We apply `sail` to detect non-linear interactions between genes and a prenatal psychosocial intervention program on cognitive performance in children at 4 years of age. Results show that individuals who are genetically predisposed to lower educational attainment are those who stand to benefit the most from the intervention. Our algorithms are implemented in an R package available on CRAN (<https://cran.r-project.org/package=sail>).

---

## 1. Introduction

Computational approaches to variable selection have become increasingly important with the advent of high-throughput technologies in genomics and brain imaging studies, where the data has become massive, yet where it is believed that the number of truly important variables is small relative to the total number of variables. Although many approaches have been developed for main effects, there is an enduring interest in powerful methods for estimating interactions, since interactions may reflect important modulation of a genomic system by an external factor and vice versa (Bhatnagar, Yang, Khundrakpam, Evans, Blanchette, Bouchard and Greenwood, 2018).

Interactions may occur in numerous types and of varying complexities. In this paper, we consider one specific type of interaction model, where one exposure variable  $E$  is involved in possibly non-linear interactions with a high-dimensional set of measures  $\mathbf{X}$  leading to effects on a response variable,  $Y$ . We propose a multivariable penalization procedure for detecting non-linear interactions between  $\mathbf{X}$  and  $E$ . Our method is motivated by the Nurse Family

---

\*Corresponding author

 sahir.bhatnagar@mcgill.ca (S.R. Bhatnagar)

 PurvisHall, 1020PineAve., W., MontrealQC, H3G1A2 (S.R. Bhatnagar)

ORCID(s): 0000-0001-8956-2509 (S.R. Bhatnagar); 0000-0002-5664-5698 (T. Lu); 0000-0002-2427-5696 (C.M. Greenwood)

Partnership (NFP); a program of prenatal and infancy home visiting by nurses for low-income mothers and their children (Olds, Henderson Jr, Cole, Eckenrode, Kitzman, Luckey, Pettitt, Sidora, Morris and Powers, 1998). In this intervention, NFP nurses guided pregnant women and parents of young children to improve the outcomes of pregnancy, their children's health and development, and their economic self-sufficiency, with the goal of reducing disparities over the life-course. Early intervention in young children has been shown to positively impact intellectual abilities (Campbell and Ramey, 1994), and more recent studies have shown that cognitive performance is also strongly influenced by genetic factors (Rietveld, Medland, Derringer, Yang, Esko, Martin, Westra, Shakhsabzov, Abdellaoui, Agrawal et al., 2013). Given the important role of both environment and genetics, we are interested in finding interactions between these two components on cognitive function in children.

### 1.1. A sparse additive interaction model

Let  $Y \in \mathbb{R}$  be a continuous outcome variable,  $E \in \mathbb{R}$  a binary or continuous environment/exposure vector of known importance, and  $X \in \mathbb{R}^p$  a vector of additional predictors, possibly high-dimensional. Assume that we have  $n$  observations of each quantity denoted by  $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$ ,  $X_E = (E_1, \dots, E_n) \in \mathbb{R}^n$ , and  $X = (X_1^\top, \dots, X_p^\top) \in \mathbb{R}^{n \times p}$ . Furthermore let  $f_j : \mathbb{R} \rightarrow \mathbb{R}$  be a smoothing method for variable  $X_j$  by a projection on to a set of basis functions:

$$f_j(X_j) = \sum_{\ell=1}^{m_j} \psi_{j\ell}(X_j) \beta_{j\ell}. \quad (1)$$

Here, the  $\{\psi_{j\ell}\}_{\ell=1}^{m_j}$  are a family of basis functions in  $X_j$  (Hastie, Tibshirani and Wainwright, 2015). Let  $\Psi_j$  be the  $n \times m_j$  matrix of evaluations of the  $\psi_{j\ell}$  and  $\theta_j = (\beta_{j1}, \dots, \beta_{jm_j}) \in \mathbb{R}^{m_j}$  for  $j = 1, \dots, p$  ( $\theta_j$  is a  $m_j$ -dimensional column vector of basis coefficients for the  $j$ th main effect). In this article we consider an additive interaction regression model of the form

$$Y = \beta_0 \cdot \mathbf{1}_n + \sum_{j=1}^p \Psi_j \theta_j + \beta_E X_E + \sum_{j=1}^p (X_E \circ \Psi_j) \tau_j + \epsilon, \quad (2)$$

where  $\beta_0 \in \mathbb{R}$  is the intercept,  $\beta_E \in \mathbb{R}$  is the coefficient for the environment variable,  $\tau_j = (\tau_{j1}, \dots, \tau_{jm_j}) \in \mathbb{R}^{m_j}$  are the basis coefficients for the  $j$ th interaction term,  $(X_E \circ \Psi_j)$  is the  $n \times m_j$  matrix formed by the component-wise multiplication of the column vector  $X_E$  by each column of  $\Psi_j$ , and  $\epsilon \in \mathbb{R}^n$  is a vector of i.i.d errors with mean zero and finite variance. Here we assume that  $p$  is large relative to  $n$ , and particularly that  $\sum_{j=1}^p m_j/n$  is large. Due to the large number of parameters to estimate with respect to the number of observations, one commonly-used approach in the penalization literature is to shrink the regression coefficients by placing a constraint on the values of  $(\beta_E, \theta_j, \tau_j)$ . Certain constraints have the added benefit of producing a sparse model in the sense that many of the coefficients will be set exactly to 0 (Bühlmann and Van De Geer, 2011). Such a reduced predictor set can lead to a more interpretable model with smaller prediction variance, albeit at the cost of having biased parameter estimates (Fan, Han and Liu, 2014). In light of these goals, consider the following penalized objective function:

$$Q(\Phi) = -L(\Phi) + \lambda(1-\alpha) \left( w_E \|\beta_E\| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} \|\tau_j\|_2, \quad (3)$$

where  $\Phi = (\beta_0, \beta_E, \theta_1, \dots, \theta_p, \tau_1, \dots, \tau_p)$ ,  $L(\Phi)$  is the log-likelihood function of the observations  $V_i = (Y_i, \Psi_i, X_{iE})$  for  $i = 1, \dots, n$ ,  $\|\theta_j\|_2 = \sqrt{\sum_{k=1}^{m_j} \beta_{jk}^2}$ ,  $\|\tau_j\|_2 = \sqrt{\sum_{k=1}^{m_j} \tau_{jk}^2}$ ,  $\lambda > 0$  and  $\alpha \in (0, 1)$  are adjustable tuning parameters,  $w_E, w_j, w_{jE}$  are non-negative penalty factors for  $j = 1, \dots, p$  which serve as a way of allowing parameters to be penalized differently (see Algorithm 2 for more details on how to estimate these weights). The first term in the penalty penalizes the main effects while the second term penalizes the interactions. The parameter  $\alpha$  controls the relative weight on the two penalties. Note that we do not penalize the intercept.

An issue with (3) is that since no constraint is placed on the structure of the model, it is possible that an estimated interaction term is non-zero while the corresponding main effects are zero. While there may be certain situations where this is plausible, statisticians have generally argued that interactions should only be included if the corresponding main effects are also in the model (McCullagh and Nelder, 1989). This is known as the strong heredity principle (Chipman, 1996). Indeed, large main effects are more likely to lead to detectable interactions (Cox, 1984). In the next section we discuss how a simple reparametrization of the model (3) can lead to this desirable property.

**Table 1**

Summary of reparametrization and penalty terms for strong and weak heredity *sail* model. Note that the penalty terms are identical for both model types, i.e., the reparametrization only affects the likelihood term of the objective function.

Model	Reparametrization	Penalty
Strong heredity	$\tau_j = \gamma_{jE}\beta_E\theta_j$	$\lambda(1-\alpha) \left( w_E  \beta_E  + \sum_{j=1}^p w_j \ \theta_j\ _2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE}  \gamma_{jE} $
Weak heredity	$\tau_j = \gamma_{jE}(\beta_E \cdot \mathbf{1}_{m_j} + \theta_j)$	$\lambda(1-\alpha) \left( w_E  \beta_E  + \sum_{j=1}^p w_j \ \theta_j\ _2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE}  \gamma_{jE} $

## 1.2. Strong and weak heredity

The strong heredity principle states that an interaction term can only have a non-zero estimate if its corresponding main effects are estimated to be non-zero, whereas the weak heredity principle allows for a non-zero interaction estimate as long as one of the corresponding main effects is estimated to be non-zero (Chipman, 1996). In the context of penalized regression methods, these principles can be formulated as structured sparsity (Bach, Jenatton, Mairal, Obozinski et al., 2012) problems. Several authors have proposed to modify the type of penalty in order to achieve the heredity principle (Radchenko and James, 2010; Bien, Taylor, Tibshirani et al., 2013; Lim and Hastie, 2015; Haris, Witten and Simon, 2016). We take an alternative approach. Following Choi et al. (Choi, Li and Zhu, 2010), we introduce a new set of parameters  $\gamma = (\gamma_{1E}, \dots, \gamma_{pE}) \in \mathbb{R}^p$  and reparametrize the coefficients for the interaction terms  $\tau_j$  in (2) as a function of  $\gamma_{jE}$  and the main effect parameters  $\theta_j$  and  $\beta_E$ . This reparametrization for both strong and weak heredity is summarized in Table 1.

To perform variable selection in this new parametrization, we penalize  $\gamma = (\gamma_{1E}, \dots, \gamma_{pE})$  instead of penalizing  $\tau$  as in (3), leading to the following penalized objective function:

$$Q(\Phi) = -L(\Phi) + \lambda(1-\alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_{jE}|. \quad (4)$$

An estimate of the regression parameters is given by  $\hat{\Phi} = \arg \min_{\Phi} Q(\Phi)$ . This penalty allows for the possibility of excluding the interaction term from the model even if the corresponding main effects are non-zero. Furthermore, smaller values for  $\alpha$  would lead to more interactions being included in the final model while values approaching 1 would favor main effects. Similar to the elastic net (Zou and Zhang, 2009), we fix  $\alpha$  and obtain a solution path over a sequence of  $\lambda$  values.

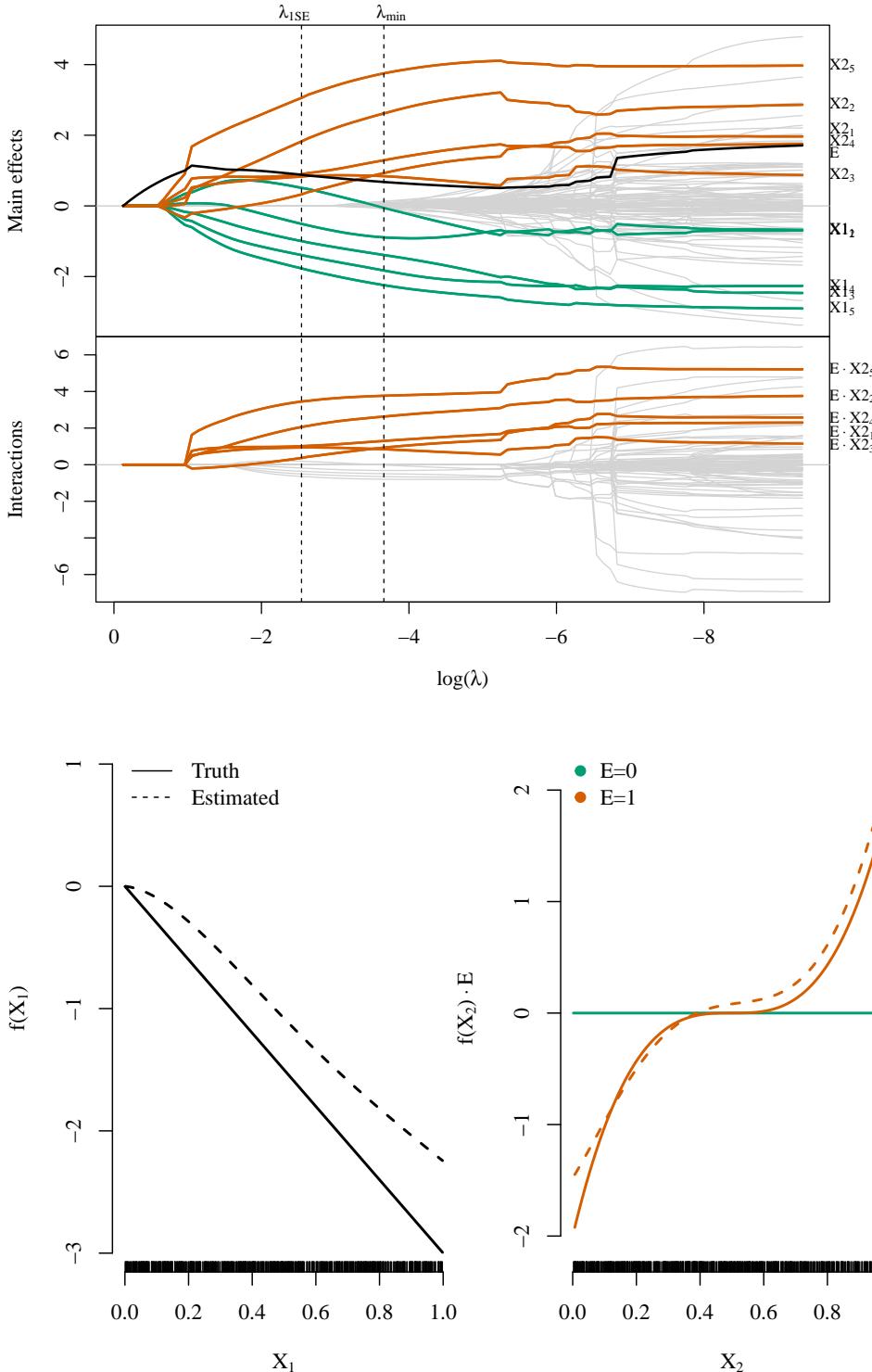
## 1.3. Toy example

We present here a toy example to better illustrate the methods proposed in this paper. With a sample size of  $n = 100$ , we sample  $p = 20$  covariates  $X_1, \dots, X_p$  independently from a  $N(0, 1)$  distribution truncated to the interval  $[0, 1]$ . Data were generated from a model which follows the strong heredity principle, but where only one covariate,  $X_2$ , is involved in an interaction with a binary exposure variable ( $E$ ):

$$Y = f_1(X_1) + f_2(X_2) + 1.75E + 1.5E \cdot f_2(X_2) + \epsilon.$$

For illustration, function  $f_1(\cdot)$  is assumed to be linear, whereas function  $f_2(\cdot)$  is non-linear:  $f_1(x) = -3x$ ,  $f_2(x) = 2(2x - 1)^3$ . The error term  $\epsilon$  is generated from a normal distribution with variance chosen such that the signal-to-noise ratio (SNR) is 2. We generated a single simulated dataset and used the strong heredity *sail* method (described below) with B-splines (df=5) to estimate the functional forms. 10-fold cross-validation (CV) was used to choose the optimal value of penalization. We used  $\alpha = 0.5$  and default values for all other arguments. We plot the solution path for both main effects and interactions in Figure 1 (top panel), coloring lines to correspond to the selected model. We see that our method is able to correctly identify the true model. We can also visually see the effect of the penalty and strong heredity principle working in tandem, i.e., the interaction term  $E \cdot f_2(X_2)$  (orange lines in the bottom panel) can only be non-zero if the main effects  $E$  and  $f_2(X_2)$  (black and orange lines respectively in the top panel) are non-zero, while non-zero main effects does not imply a non-zero interaction.

In Figure 1 (bottom panel), we plot the true and estimated component functions  $\hat{f}_1(X_1)$  and  $E \cdot \hat{f}_2(X_2)$ , and their estimates from this analysis with *sail*. We are able to capture the shape of the correct functional form. Lack-of-fit for  $f_1(X_1)$  can be partially explained by acknowledging that *sail* is trying to fit a spline to a linear function. Nevertheless, this example demonstrates that *sail* can still identify trends reasonably well.



**Figure 1: Top:** Toy example solution path for main effects (top) and interactions (bottom).  $\{X_{11}, X_{12}, X_{13}, X_{14}, X_{15}\}$  and  $\{X_{21}, X_{22}, X_{23}, X_{24}, X_{25}\}$  are the five basis coefficients for  $X_1$  and  $X_2$ , respectively.  $\lambda_{ISE}$  is the largest value of penalization for which the CV error is within one standard error of the minimizing value  $\lambda_{min}$ . **Bottom:** Estimated smooth functions for  $X_1$  and the  $X_2 \cdot E$  interaction by the sail method based on  $\lambda_{min}$ .

## 1.4. Related work

Methods for variable selection of interactions can be broken down into two categories: linear and non-linear interaction effects. Many of the linear effect methods consider all pairwise interactions in  $\mathbf{X}$  (Zhao, Rocha and Yu, 2009; Choi et al., 2010; Bien et al., 2013; She and Jiang, 2014) which can be computationally prohibitive when  $p$  is large. More recent proposals for selection of interactions allow the user to restrict the search space to interaction candidates (Lim and Hastie, 2015; Haris et al., 2016). This is useful when the researcher wants to impose prior information on the model. Two-stage procedures, where interaction candidates are considered from an original screen of main effects, have shown good performance when  $p$  is large (Hao, Feng and Zhang, 2018; Shah, 2016) in the linear setting. There are many fewer methods available for estimating non-linear interactions. For example, Radchenko and James (2010) (Radchenko and James, 2010) proposed a model of the form  $Y = \beta_0 + \sum_{j=1}^p f_j(X_j) + \sum_{j>k} f_{jk}(X_j, X_k) + \varepsilon$ , where  $f(\cdot)$  are smooth component functions. This method is more computationally expensive than `sail` since it considers all pairwise interactions between the basis functions, and its effectiveness in simulations or real-data applications is unknown as there is no software implementation.

The main contributions of this paper are five-fold. First, we develop a model for non-linear interactions with a key exposure variable, following either the weak or strong heredity principle, that is computationally efficient and scales to the high-dimensional setting ( $n \ll p$ ). Second, through simulation studies, we show improved performance in terms of prediction accuracy and support recovery over existing methods that only consider linear interactions or additive main effects. Third, we show that our method possesses the oracle property (Fan and Li, 2001), i.e., it performs as well as if the true model were known in advance. Fourth, we demonstrate the performance of our method in two applications: 1) gene-environment interactions in a prenatal psychosocial intervention program Olds et al. (1998) and 2) a study aimed at identifying which clinical variables influence mortality rates amongst seriously ill hospitalized patients (Connors, Dawson, Desbiens, Fulkerson, Goldman, Knaus, Lynn, Oye, Bergner, Damiano et al., 1995). Fifth, we implement our algorithms in the `sail` R package on CRAN (<https://cran.r-project.org/package=sail>), along with extensive documentation. In particular, our implementation also allows for linear interaction models, user-defined basis expansions, a cross-validation procedure for selecting the optimal tuning parameter, and differential shrinkage parameters to apply the adaptive lasso idea (Zou, 2006).

The rest of the paper is organized as follows. Section 2 describes our optimization procedure and some details about the algorithm used to fit the `sail` model for the least squares case. Theoretical results are given in Section 3. In Section 4, through simulation studies we compare the performance of our proposed approach and demonstrate the scenarios where it can be advantageous to use `sail` over existing methods. Section 5 contains two real data examples and Section 6 discusses some limitations and future directions.

## 2. Computation

In this section we describe a blockwise coordinate descent algorithm for fitting the least-squares version of the `sail` model in (4). We fix the value for  $\alpha$  and minimize the objective function over a decreasing sequence of  $\lambda$  values ( $\lambda_{\max} > \dots > \lambda_{\min}$ ). We use the subgradient equations to determine the maximal value  $\lambda_{\max}$  such that all estimates are zero. Due to the heredity principle, this reduces to finding the largest  $\lambda$  such that all main effects  $(\beta_E, \theta_1, \dots, \theta_p)$  are zero. Following Friedman et al. (Friedman, Hastie and Tibshirani, 2010), we construct a  $\lambda$ -sequence of 100 values decreasing from  $\lambda_{\max}$  to  $0.001\lambda_{\max}$  on the log scale, and use the warm start strategy where the solution for  $\lambda_\ell$  is used as a starting value for  $\lambda_{\ell+1}$ .

### 2.1. Blockwise coordinate descent for least-squares loss

We assume that  $Y$ ,  $\Psi_j$ ,  $X_E$  and  $X_E \circ \Psi_j$  have been centered by their sample means  $\bar{Y}$ ,  $\bar{\Psi}_j$ ,  $\bar{X}_E$ , and  $\bar{X}_E \circ \Psi_j$ , respectively. Here,  $\bar{\Psi}_j \in \mathbb{R}^{m_j}$  and  $\bar{X}_E \circ \Psi_j \in \mathbb{R}^{m_j}$  represent the column means of  $\Psi_j$  and  $X_E \circ \Psi_j$ , respectively. Since the intercept ( $\beta_0$ ) is not penalized and all variables have been centered, we can omit it from the loss function and compute it once the algorithm has converged for all other parameters. The strong heredity `sail` model with least-squares loss has the form:

$$\hat{Y} = \sum_{j=1}^p \Psi_j \theta_j + \beta_E X_E + \sum_{j=1}^p \gamma_{jE} \beta_E (X_E \circ \Psi_j) \theta_j , \quad (5)$$

and the objective function is given by

$$Q(\Phi) = \frac{1}{2n} \|Y - \hat{Y}\|_2^2 + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_{jE}|. \quad (6)$$

Solving (6) in a blockwise manner allows us to leverage computationally fast algorithms for  $\ell_1$  and  $\ell_2$  norm penalized regression. Indeed, by careful construction of pseudo responses and pseudo design matrices, existing efficient algorithms can be used to estimate the parameters. The objective function simplifies to a modified lasso problem when holding all  $\theta_j$  fixed, and a modified group lasso problem when holding  $\beta_E$  and all  $\gamma_{jE}$  fixed. The main computations are provided in Algorithm 1. A more detailed version of the derivations are given in Supplemental Section B.1.

**Algorithm 1** Blockwise Coordinate Descent for Least-Squares sail with Strong Heredity

---

```

1: function sail( $X, Y, X_E, \text{basis}, \lambda, \alpha, w_j, w_E, w_{jE}, \epsilon$ ) ▷ Algorithm for solving (6)
2:    $\Psi_j \leftarrow \text{basis}(X_j), \tilde{\Psi}_j \leftarrow X_E \circ \Psi_j$  for  $j = 1, \dots, p$ 
3:   Center all variables by their sample means
4:   Initialize:  $\beta_E^{(0)} = \theta_j^{(0)} = \gamma_j^{(0)} \leftarrow 0$  for  $j = 1, \dots, p$ .
5:   Set iteration counter  $k \leftarrow 0$ 
6:    $R^* \leftarrow Y - \beta_E^{(k)} X_E - \sum_j (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)}$ 
7:   repeat
8:     • To update  $\gamma = (\gamma_1, \dots, \gamma_p)$ 
9:        $\tilde{X}_j \leftarrow \beta_E^{(k)} \tilde{\Psi}_j \theta_j^{(k)}$  for  $j = 1, \dots, p$ 
10:       $R \leftarrow R^* + \sum_{j=1}^p \gamma_j^{(k)} \tilde{X}_j$ 
11:

$$\gamma^{(k)(new)} \leftarrow \arg \min_{\gamma} \frac{1}{2n} \left\| R - \sum_j \gamma_j \tilde{X}_j \right\|_2^2 + \lambda \alpha \sum_j w_{jE} |\gamma_j|$$

12:       $\Delta = \sum_j (\gamma_j^{(k)} - \gamma_j^{(k)(new)}) \tilde{X}_j$ 
13:       $R^* \leftarrow R^* + \Delta$ 
14:     • To update  $\theta = (\theta_1, \dots, \theta_p)$ 
15:        $\tilde{X}_j \leftarrow \Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j$  for  $j = 1, \dots, p$ 
16:       for  $j = 1, \dots, p$  do
17:          $R \leftarrow R^* + \tilde{X}_j \theta_j^{(k)}$ 
18:

$$\theta_j^{(k)(new)} \leftarrow \arg \min_{\theta_j} \frac{1}{2n} \|R - \tilde{X}_j \theta_j\|_2^2 + \lambda(1 - \alpha) w_j \|\theta_j\|_2$$

19:       $\Delta = \tilde{X}_j (\theta_j^{(k)} - \theta_j^{(k)(new)})$ 
20:       $R^* \leftarrow R^* + \Delta$ 
21:     • To update  $\beta_E$ 
22:        $\tilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \tilde{\Psi}_j \theta_j^{(k)}$ 
23:        $R \leftarrow R^* + \beta_E^{(k)} \tilde{X}_E$ 
24:

$$\beta_E^{(k)(new)} \leftarrow \frac{1}{\tilde{X}_E^\top \tilde{X}_E} S \left( \frac{1}{n \cdot w_E} \tilde{X}_E^\top R, \lambda(1 - \alpha) \right)$$

 $\triangleright S(x, t) = \text{sign}(x)(|x| - t)_+$ 
25:       $\Delta = (\beta_E^{(k)} - \beta_E^{(k)(new)}) \tilde{X}_E$ 
26:       $R^* \leftarrow R^* + \Delta$ 
27:       $k \leftarrow k + 1$ 
28:
29:      until convergence criterion is satisfied:  $|Q(\Phi^{(k-1)}) - Q(\Phi^{(k)})| / Q(\Phi^{(k-1)}) < \epsilon$ 
30:      Compute the intercept  $\beta_0$ 
31:       $\beta_0 \leftarrow \bar{Y} - \sum_{j=1}^p \bar{\Psi}_j \hat{\theta}_j - \hat{\beta}_E \bar{X}_E - \sum_{j=1}^p \hat{\gamma}_j \hat{\beta}_E (\bar{X}_E \circ \bar{\Psi}_j) \hat{\theta}_j$ 

```

---

## 2.2. Details on Update for $\theta$

Here we discuss a computational speedup in the updates for the  $\theta$  parameter. The partial residual ( $R_s$ ) used for updating  $\theta_s$  ( $s \in 1, \dots, p$ ) at the  $k$ th iteration is given by

$$R_s = Y - \tilde{Y}_{(-s)}^{(k)}, \tag{7}$$

where  $\tilde{Y}_{(-s)}^{(k)}$  is the fitted value at the  $k$ th iteration excluding the contribution from  $\Psi_s$ :

$$\tilde{Y}_{(-s)}^{(k)} = \beta_E^{(k)} X_E + \sum_{\ell \neq s} \Psi_\ell \theta_\ell^{(k)} + \sum_{\ell \neq s} \gamma_\ell^{(k)} \beta_E^{(k)} \tilde{\Psi}_\ell \theta_\ell^{(k)}. \quad (8)$$

Using (8), (7) can be re-written as

$$\begin{aligned} R_s &= Y - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \\ &= R^* + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)}, \end{aligned} \quad (9)$$

where

$$R^* = Y - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)}. \quad (10)$$

Denote  $\theta_s^{(k)(new)}$  the solution for predictor  $s$  at the  $k$ th iteration, given by:

$$\theta_s^{(k)(new)} = \arg \min_{\theta_j} \frac{1}{2n} \|R_s - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_j\|_2^2 + \lambda(1-\alpha) w_s \|\theta_j\|_2. \quad (11)$$

Now we want to update the parameters for the next predictor  $\theta_{s+1}$  ( $s+1 \in 1, \dots, p$ ) at the  $k$ th iteration. The partial residual used to update  $\theta_{s+1}$  is given by

$$R_{s+1} = R^* + (\Psi_{s+1} + \gamma_{s+1}^{(k)} \beta_E^{(k)} \tilde{\Psi}_{s+1}) \theta_{s+1}^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) (\theta_s^{(k)} - \theta_s^{(k)(new)}), \quad (12)$$

where  $R^*$  is given by (10),  $\theta_s^{(k)}$  is the parameter value prior to the update, and  $\theta_s^{(k)(new)}$  is the updated value given by (11). Taking the difference between (9) and (12) gives

$$\begin{aligned} \Delta &= R_t - R_s \\ &= (\Psi_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\Psi}_t) \theta_t^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) (\theta_s^{(k)} - \theta_s^{(k)(new)}) - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \\ &= (\Psi_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\Psi}_t) \theta_t^{(k)} - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)(new)}. \end{aligned} \quad (13)$$

Therefore  $R_t = R_s + \Delta$ , and the partial residual for updating the next predictor can be computed by updating the previous partial residual by  $\Delta$ , given by (13). This formulation can lead to computational speedups especially when  $\Delta = 0$ , meaning the partial residual does not need to be re-calculated.

### 2.3. Weak Heredity

Our method can be easily adapted to enforce the weak heredity property. That is, an interaction term can only be present if at least one of its corresponding main effects is non-zero. To do so, we reparametrize the coefficients for the interaction terms in (2) as  $\tau_j = \gamma_{jE}(\beta_E \cdot \mathbf{1}_{m_j} + \theta_j)$ , where  $\mathbf{1}_{m_j}$  is a vector of ones with dimension  $m_j$  (i.e. the length of  $\theta_j$ ). We defer the algorithm details for fitting the sail model with weak heredity in Supplemental Section B.4, as it is very similar to Algorithm 1 for the strong heredity sail model.

### 2.4. Adaptive sail

The weights for the environment variable, main effects and interactions are given by  $w_E$ ,  $w_j$  and  $w_{jE}$  respectively. These weights serve as a means of allowing a different penalty to be applied to each variable. In particular, any variable with a weight of zero is not penalized at all. This feature is usually selected for one of two reasons:

1. Prior knowledge about the importance of certain variables is known. Larger weights will penalize the variable more, while smaller weights will penalize the variable less
2. Allows users to apply the adaptive sail, similar to the adaptive lasso (Zou, 2006)

We describe the adaptive sail in Algorithm 2. This is a general procedure that can be applied to the weak and strong heredity settings. We provide this capability in the sail package using the `penalty.factor` argument.

**Algorithm 2** Adaptive sail algorithm

- 
1. For a decreasing sequence  $\lambda = \lambda_{max}, \dots, \lambda_{min}$  and fixed  $\alpha$  run the sail algorithm
  2. Use cross-validation or a data splitting procedure to determine the optimal value for the tuning parameter:  $\lambda^{[opt]} \in \{\lambda_{max}, \dots, \lambda_{min}\}$
  3. Let  $\widehat{\beta}_E^{[opt]}, \widehat{\theta}_j^{[opt]}$  and  $\widehat{\tau}_j^{[opt]}$  for  $j = 1, \dots, p$  be the coefficient estimates corresponding to the model at  $\lambda^{[opt]}$
  4. Set the weights to be  

$$w_E = \left( \left| \widehat{\beta}_E^{[opt]} \right| + 1/n \right)^{-1}, w_j = \left( \|\widehat{\theta}_j^{[opt]}\|_2 + 1/n \right)^{-1}, w_{jE} = \left( \|\widehat{\tau}_j^{[opt]}\|_2 + 1/n \right)^{-1} \text{ for } j = 1, \dots, p$$
  5. Run the sail algorithm with the weights defined in step 4), and use cross-validation or a data splitting procedure to choose the optimal value of  $\lambda$
- 

**2.5. Flexible design matrix**

The definition of the basis expansion functions in (1) is very flexible, in the sense that our algorithms are independent of this choice. As a result, the user can apply any basis expansion they desire. In the extreme case, one could apply the identity map, i.e.,  $f_j(X_j) = X_j$  which leads to a linear interaction model (referred to as **linear sail**). When little information is known a priori about the relationship between the predictors and the response, by default, we choose to apply the same basis expansion to all columns of  $X$ . This is a reasonable approach when all the variables are continuous. However, there are often situations when the data contains a combination of categorical and continuous variables. In these cases it may be sub-optimal to apply a basis expansion to the categorical variables. Owing to the flexible nature of our algorithm, we can handle this scenario in our implementation by allowing a user-defined design matrix. The only extra information needed is the group membership of each column in the design matrix. We illustrate such an example in a vignette of the sail R package.

**3. Theory**

In this section we study the asymptotic behaviour of the sail estimator  $\widehat{\Phi}$ , defined as the minimizer of (4), as well as the model selection properties. We show that sail possesses the oracle property when the sample size approaches infinity and the number of predictors is fixed. That is, under certain regularity conditions, it performs as well as if the true model were known in advance and has the optimal estimation rate (Zou, 2006). The regularity conditions and proofs are given in Supplemental Section 1.

Let  $\Phi^* = (\beta_E^*, \theta_1^{*\top}, \dots, \theta_p^{*\top}, \gamma_{1E}^*, \dots, \gamma_{pE}^*)^\top$  denote the unknown vector of true coefficients in (4). To simplify the notation, we use the representation  $\Phi^* = (\phi_1^{*\top}, \phi_2^{*\top}, \dots, \phi_{p+1}^{*\top}, \phi_{p+2}^{*\top}, \dots, \phi_{2p+1}^{*\top})^\top$ , where  $\phi_1^* = \beta_E^*$ ,  $\phi_2^* = \theta_1^*, \dots, \phi_{p+1}^* = \theta_p^*$ , and  $\phi_{p+2}^* = \gamma_{1E}^*, \dots, \phi_{2p+1}^* = \gamma_{pE}^*$ . Denote by  $\mathcal{A} = \{m : \phi_m^* \neq \mathbf{0}\}$  the unknown sparsity pattern of  $\Phi^*$ , and  $\widehat{\mathcal{A}} = \{m : \widehat{\phi}_m \neq \mathbf{0}\}$  the estimated sail model selector. We can rewrite the penalty terms in (4), and consider the sail estimates  $\widehat{\Phi}_n$  given b

$$\widehat{\Phi}_n = \arg \min_{\Phi} Q_n(\Phi) = -L_n(\Phi) + n\lambda_m \sum_{m=1}^{2p+1} \|\phi_m\|_2, \quad (14)$$

where  $\lambda_1 = \lambda(1 - \alpha)w_E$ ,  $\lambda_m = \lambda(1 - \alpha)w_m$  for  $m = 2, \dots, p + 1$ , and  $\lambda_m = \lambda\alpha w_{mE}$  for  $m = p + 2, \dots, 2p + 1$ . Define

$$\mathcal{A}_1 = \{m : \phi_m^* \neq \mathbf{0} (1 \leq m \leq p + 1)\}, \quad \mathcal{A}_2 = \{m : \phi_m^* \neq \mathbf{0} (p + 2 \leq m \leq 2p + 1)\}, \quad \mathcal{A} = \mathcal{A}_1 \cup \mathcal{A}_2$$

that is,  $\mathcal{A}_1$  contains the indices for main effects whose true coefficients are non-zero, and  $\mathcal{A}_2$  contains the indices for interaction terms whose true coefficients are non-zero. Let

$$a_n = \max \{ \lambda_m, \lambda_{m'} : m \in \mathcal{A}_1, m' \in \mathcal{A}_2 \}$$

$$b_n = \min \left\{ \lambda_m, \lambda_{m'} : m \in \mathcal{A}_1^c, m' \in \mathcal{A}_2^c \text{ s.t. } \phi_{m'}^* = \gamma_{jE}^* = 0 \text{ but } \beta_E^* \neq 0 \text{ and } \theta_j^* \neq \mathbf{0} \quad (1 \leq j \leq p) \right\}$$

Note that our asymptotic results are stated for the main effects and interaction terms only, even though our formulation includes an unpenalized intercept. Consistency results immediately follow for  $\beta_0$  since we assume the data has been centered, leading to a closed form solution for the intercept in the least-squares setting.

**Lemma 1.** [Existence of a local minimizer] If  $a_n = o(\frac{1}{\sqrt{n}})$  as  $n \rightarrow \infty$ , i.e.  $\sqrt{n}a_n \rightarrow 0$ , then  $\|\hat{\Phi}_n - \Phi^*\|_2 = O_p(\frac{1}{\sqrt{n}})$

Lemma (1) states that if the tuning parameters corresponding to the non-zero coefficients converge to 0 at a speed faster than  $\frac{1}{\sqrt{n}}$ , then there exists a local minimizer of  $Q_n(\Phi)$  which is  $\sqrt{n}$ -consistent (Wang, Li and Tsai, 2007; Choi et al., 2010).

**Theorem 1** (Model selection consistency). *If  $\sqrt{n}a_n \rightarrow 0$  and  $\sqrt{nb_n} \rightarrow \infty$ , then*

$$P\left(\hat{\Phi}_{\mathcal{A}_1^c} = \mathbf{0}\right) \rightarrow 1 \quad \text{and} \quad P\left(\hat{\Phi}_{\mathcal{A}_2^c} = \mathbf{0}\right) \rightarrow 1 \quad (15)$$

Theorem (1) shows that `sail` can consistently remove the main effects and interaction terms which are not associated with the response with high probability. Together with Lemma (1), we see that the asymptotic behaviour of the penalty terms for the zero and non-zero predictors must be different to satisfy the model selection consistency property (15) (Nardi, Rinaldo et al., 2008). Specifically, when the tuning parameters for the non-zero coefficients converge to 0 faster than  $1/\sqrt{n}$  (i.e.  $\sqrt{n}a_n \rightarrow 0$ ) and those for zero coefficients are large enough (i.e.  $\sqrt{nb_n} \rightarrow \infty$ ), the Lemma (1) and Theorem (1) imply that the  $\sqrt{n}$ -consistent estimator  $\hat{\Phi}_n$  satisfies  $P\left(\hat{\Phi}_{\mathcal{A}_2^c} = \mathbf{0}\right) \rightarrow 1$ .

Next, we obtain the asymptotic distribution of the `sail` estimator.

**Theorem 2** (Asymptotic normality). *Denote  $\mathcal{A} = \mathcal{A}_1 \cup \mathcal{A}_2$ . Assume that  $\sqrt{n}a_n \rightarrow 0$  and  $\sqrt{nb_n} \rightarrow \infty$ . Under the regularity conditions, the subvector  $\hat{\Phi}_{\mathcal{A}}$  of the local minimizer  $\hat{\Phi}_n$  given in Lemma (1) satisfies*

$$\sqrt{n}\left(\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*\right) \xrightarrow{d} N\left(\mathbf{0}, \mathbf{I}^{-1}(\Phi_{\mathcal{A}}^*)\right), \quad (16)$$

where  $\mathbf{I}(\Phi_{\mathcal{A}}^*)$  is the Fisher information matrix for  $\Phi_{\mathcal{A}}$  at  $\Phi_{\mathcal{A}} = \Phi_{\mathcal{A}}^*$ , assuming  $\mathcal{A}_c$  is known in advance.

Together, Theorems (1) and (2) establish that if the tuning parameters satisfy the conditions  $\sqrt{n}a_n \rightarrow 0$  and  $\sqrt{nb_n} \rightarrow \infty$ , then as the sample size grows large, `sail` has the oracle property (Fan and Li, 2001). In order for the conditions on the tuning parameters to be satisfied, we follow the strategies outlined for the adaptive Lasso (Zou, 2006), the adaptive group Lasso (Nardi et al., 2008) and the adaptive elastic-net (Zou and Zhang, 2009). That is, we define the adaptive weights as  $w_m = \|\hat{\phi}_m^{\text{init}} + 1/n\|_2^{-\xi}$  for  $m = 1, \dots, 2p+1$ , where  $\xi$  is a positive constant and  $\hat{\phi}_m^{\text{init}}$  is an initial  $\sqrt{n}$ -consistent estimate of  $\phi_m^*$ . Here, the  $1/n$  is to avoid division by zero.

## 4. Simulation Study

In this section, we use simulated data to understand the performance of `sail` in different scenarios.

### 4.1. Comparator Methods

Since there are no other packages that directly address our chosen problem, we selected comparator methods based on the following criteria: 1) penalized regression methods that can handle high-dimensional data ( $n < p$ ), 2) allowing at least one of linear effects, non-linear effects or interaction effects, and 3) having a software implementation in R. The selected methods can be grouped into three categories:

1. Linear main effects: `lasso` (Tibshirani, 1996), `adaptive lasso` (Zou, 2006)
2. Linear interactions: `lassoBT` (Shah, 2016), `GLinternet` (Lim and Hastie, 2015)
3. Non-linear main effects: `HierBasis` (Haris, Shojaie and Simon, 2019), `SPAM` (Ravikumar, Lafferty, Liu and Wasserman, 2009), `gamsel` (Chouldechova and Hastie, 2015)

For GLinternet we specified the `interactionCandidates` argument so as to only consider interactions between the environment and all other  $X$  variables. For all other methods we supplied  $(X, X_E)$  as the data matrix, 100 for the number of tuning parameters to fit, and used the default values otherwise (R code for each method available at [https://github.com/sahirbhatnagar/sail/blob/master/my\\_sims/method\\_functions.R](https://github.com/sahirbhatnagar/sail/blob/master/my_sims/method_functions.R)). lassoBT considers all pairwise interactions as there is no way for the user to restrict the search space. SPAM applies the same basis expansion to every column of the data matrix; we chose 5 basis spline functions. HierBasis and gamsel selects whether a term in an additive model is non-zero, linear, or a non-linear spline up to a specified max degrees of freedom per variable.

We compare the above listed methods with our main proposal method `sail`, as well as with `adaptive sail` (Algorithm 2) and `sail weak` which has the weak heredity property. For each function  $f_j$ , we use a B-spline basis matrix with `degree=5` implemented in the `bs` function in R (R Core Team, 2017). We center the environment variable and the basis functions before running the `sail` method.

## 4.2. Simulation Design

To make the comparisons with other methods as fair as possible, we followed a simulation framework that has been previously used for variable selection methods in additive models (Lin, Zhang et al., 2006; Huang, Horowitz and Wei, 2010). We extend this framework to include interaction effects as well. The covariates are simulated as follows. First, we generate  $x_1, \dots, x_{1000}$  independently from a standard normal distribution truncated to the interval  $[0,1]$  for  $i = 1, \dots, n$ . The first four variables are non-zero (i.e. active in the response), while the rest of the variables are zero (i.e. are noise variables). The exposure variable ( $X_E$ ) is generated from a standard normal distribution truncated to the interval  $[-1,1]$ . The outcome  $Y$  is then generated following one of the models and assumptions described below. We evaluate the performance of our method on three of its defining characteristics: 1) the strong heredity property, 2) non-linearity of predictor effects and 3) interactions. Simulation scenarios are designed specifically to test the performance of these characteristics.

### 1. Heredity simulation

Scenario (a) Truth obeys strong heredity. In this situation, the true model for  $Y$  contains main effect terms for all covariates involved in interactions.

$$Y = \sum_{j=1}^4 f_j(X_j) + \beta_E \cdot X_E + X_E \cdot f_3(X_3) + X_E \cdot f_4(X_4) + \epsilon$$

Scenario (b) Truth obeys weak heredity. Here, in addition to the interaction, the  $E$  variable has its own main effect but the covariates  $X_3$  and  $X_4$  do not.

$$Y = f_1(X_1) + f_2(X_2) + \beta_E \cdot X_E + X_E \cdot f_3(X_3) + X_E \cdot f_4(X_4) + \epsilon$$

Scenario (c) Truth only has interactions. In this simulation, the covariates involved in interactions do not have main effects as well.

$$Y = X_E \cdot f_3(X_3) + X_E \cdot f_4(X_4) + \epsilon$$

### 2. Non-linearity simulation scenario

Truth is linear. `sail` is designed to model non-linearity; here we assess its performance if the true model is completely linear.

$$Y = 5X_1 + 3(X_2 + 1) + 4X_3 + 6(X_4 - 2) + \beta_E \cdot X_E + X_E \cdot 4X_3 + X_E \cdot 6(X_4 - 2) + \epsilon$$

### 3. Interactions simulation scenario

Truth only has main effects. `sail` is designed to capture interactions; here we assess its performance when there are none in the true model.

$$Y = \sum_{j=1}^4 f_j(X_j) + \beta_E \cdot X_E + \epsilon$$

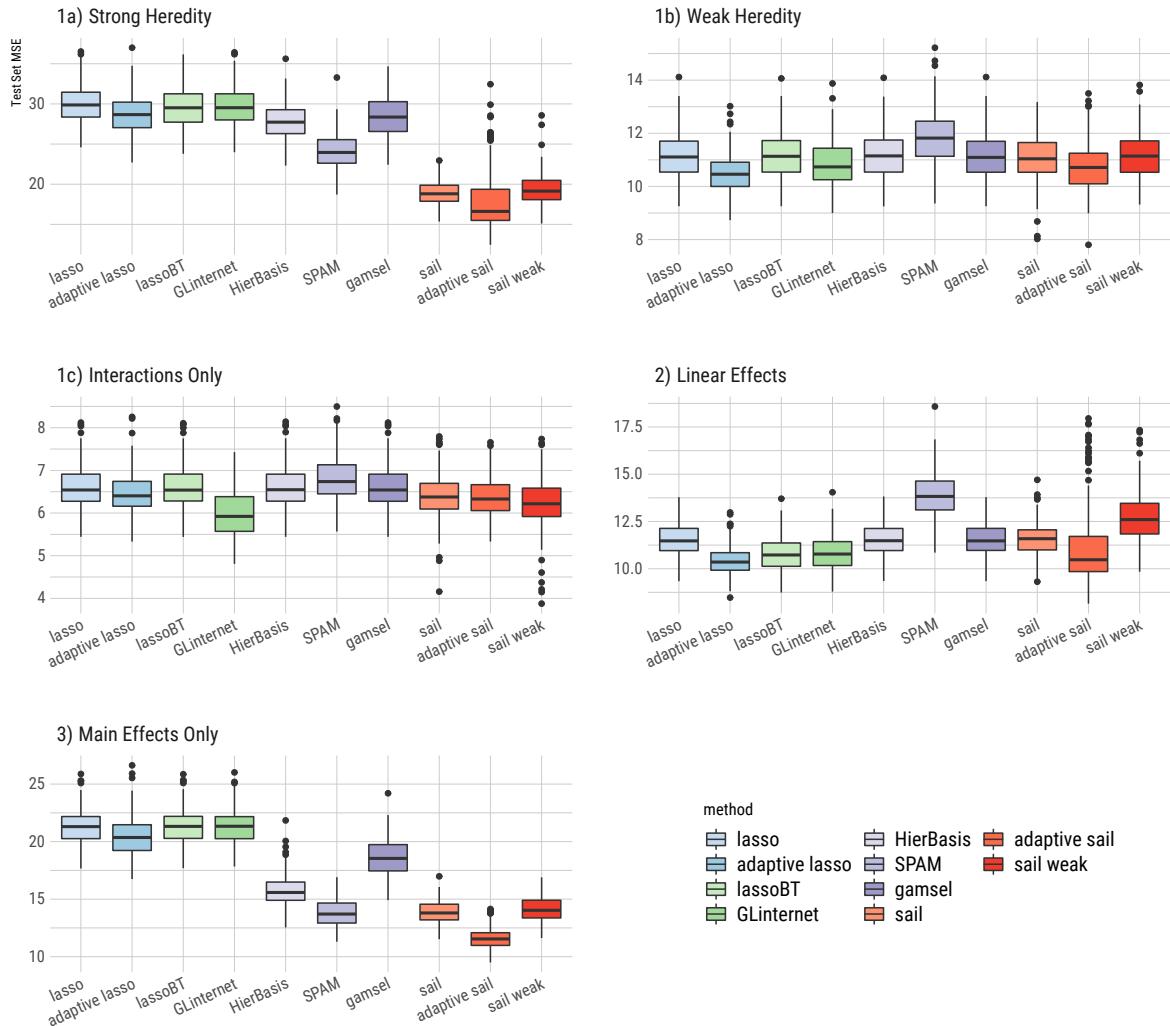
The true component functions are the same as in (Lin et al., 2006; Huang et al., 2010) and are given by  $f_1(t) = 5t$ ,  $f_2(t) = 3(2t - 1)^2$ ,  $f_3(t) = 4 \sin(2\pi t)/(2 - \sin(2\pi t))$ ,  $f_4(t) = 6(0.1 \sin(2\pi t) + 0.2 \cos(2\pi t) + 0.3 \sin(2\pi t)^2 + 0.4 \cos(2\pi t)^3 + 0.5 \sin(2\pi t)^3)$ . We set  $\beta_E = 2$  and draw  $\epsilon$  from a normal distribution with variance chosen such that the signal-to-noise ratio is 2. Using this setup, we generated 200 replications consisting of a training set of  $n = 200$ , a validation set of  $n = 200$  and a test set of  $n = 800$ . The training set was used to fit the model and the validation set was used to select the optimal tuning parameter corresponding to the minimum prediction mean squared error (MSE). Variable selection results including true positive rate, false positive rate and number of active variables (the number of variables with a non-zero coefficient estimate) were assessed on the training set, and MSE was assessed on the test set.

### 4.3. Results

The prediction accuracy and variable selection results for each of the five simulation scenarios are shown in Figure 2 and Table 2, respectively. We see that `sail`, `adaptive sail` and `sail weak` have the best performance in terms of both MSE and yielding correct sparse models when the truth follows a strong heredity (scenario 1a), as we would expect, since this is exactly the scenario that our method is trying to target. Our method is also competitive when only main effects are present (scenario 3) and performs just as well as methods that only consider linear and non-linear main effects (`HierBasis`, `SPAM`), owing to the penalization applied to the interaction parameter. Due to the heredity property being violated in scenario 1c), no method can identify the correct model with the exception of `GLinternet`. When only linear effects and interactions are present (scenario 2), we see that `adaptive sail` has similar MSE compared to the other linear interaction methods (`lassoBT` and `GLinternet`) with a better TPR and FPR. It is important to note that the variable selection performance of `sail` is highly dependent on being able to correctly select the exposure variable ( $X_E$ ). In Supplemental Section C, we show the selection rates of  $X_E$ . We see that `sail` is able to consistently identify the exposure variable across all simulation scenarios and replications. Overall, our simulation study suggests that `sail` outperforms existing methods when the true model contains non-linear interactions, and is competitive even when the truth only has either linear or additive main effects.

## Test Set MSE

Based on 200 simulations



**Figure 2:** Boxplots of the test set mean squared error from 200 replications for each of the five simulation scenarios.

Table 2: Mean (standard deviation) of the number of selected variables ( $|\hat{\mathcal{J}}|$ ), true positive rate (TPR) and false positive rate (FPR) as a percentage from 200 replications for each of the five scenarios.  $|\mathcal{J}|$  is the number of truly associated variables.

	Linear		Linear		Non-linear		Non-linear	
	Main Effects	adaptive lasso	lassoBT	GLinternet	HierBasis	SPAM	gamsel	sail
<b>1a) Strong heredity (<math> \mathcal{J}  = 7</math>)</b>								
$ \hat{\mathcal{J}} $	28 (15)	8 (4)	35 (18)	40 (20)	133 (48)	42 (19)	46 (21)	37 (15)
TPR	53.9 (8.4)	49.3 (10.1)	61.7 (11.5)	66.4 (14.0)	65.2 (8.1)	60.9 (8.5)	56.9 (7.7)	89.5 (8.2)
FPR	1.2 (0.7)	0.2 (0.2)	1.5 (0.9)	1.8 (1.0)	6.5 (2.4)	1.9 (0.9)	2.1 (1.1)	1.5 (0.7)
<b>1b) Weak heredity (<math> \mathcal{J}  = 5</math>)</b>								
$ \hat{\mathcal{J}} $	19 (12)	4 (2)	20 (13)	38 (23)	24 (23)	28 (16)	21 (15)	24 (19)
TPR	40.7 (3.6)	40.1 (1.4)	40.8 (3.8)	64.1 (14.9)	42.2 (6.3)	53.9 (9.4)	42.7 (6.8)	52.4 (11.4)
FPR	0.9 (0.6)	0.1 (0.1)	0.9 (0.7)	1.7 (1.1)	1.1 (1.1)	1.2 (0.8)	1.0 (0.7)	1.0 (0.9)
<b>1c) Interactions Only (<math> \mathcal{J}  = 2</math>)</b>								
$ \hat{\mathcal{J}} $	12 (12)	3 (2)	14 (13)	38 (21)	12 (13)	13 (12)	12 (12)	10 (18)
TPR	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	81.4 (27.0)	0.0 (0.0)	0.0 (0.0)	0.0 (0.0)	1.0 (6.9)
FPR	0.6 (0.6)	0.6 (6.9)	0.7 (0.7)	1.8 (1.0)	0.6 (0.7)	0.7 (0.6)	0.6 (0.6)	0.5 (0.9)
<b>2) Linear Effects (<math> \mathcal{J}  = 7</math>)</b>								
$ \hat{\mathcal{J}} $	37 (17)	8 (3)	48 (19)	51 (23)	37 (19)	42 (19)	37 (16)	34 (18)
TPR	70.4 (3.7)	67.2 (6.7)	72.3 (6.3)	93.4 (8.5)	70.3 (3.8)	65.0 (8.1)	70.4 (3.7)	93.9 (9.9)
FPR	1.6 (0.8)	0.2 (0.2)	2.2 (1.0)	2.2 (1.2)	1.6 (0.9)	1.9 (0.9)	1.6 (0.8)	1.4 (0.9)
<b>3) Main Effects Only (<math> \mathcal{J}  = 5</math>)</b>								
$ \hat{\mathcal{J}} $	29 (14)	7 (4)	31 (15)	34 (18)	154 (17)	46 (21)	56 (20)	44 (19)
TPR	75.9 (10.9)	66.5 (15.3)	76.0 (10.9)	77.0 (9.5)	97.5 (6.6)	93.1 (10.7)	81.3 (9.5)	91.5 (10.3)
FPR	1.3 (0.7)	0.2 (0.2)	1.3 (0.8)	1.5 (0.9)	7.5 (0.9)	2.1 (1.0)	2.6 (1.0)	2.0 (0.9)
								0.2 (0.1)
								0.9 (0.1)

We also plotted the true and predicted curves for scenario 1a) in Supplemental Section C, to visually inspect whether our method could correctly capture the shape of the association between the predictors and the response for both main and interaction effects. In general, we see the non-linear effects are clearly being captured by sail.

## 5. Real data applications

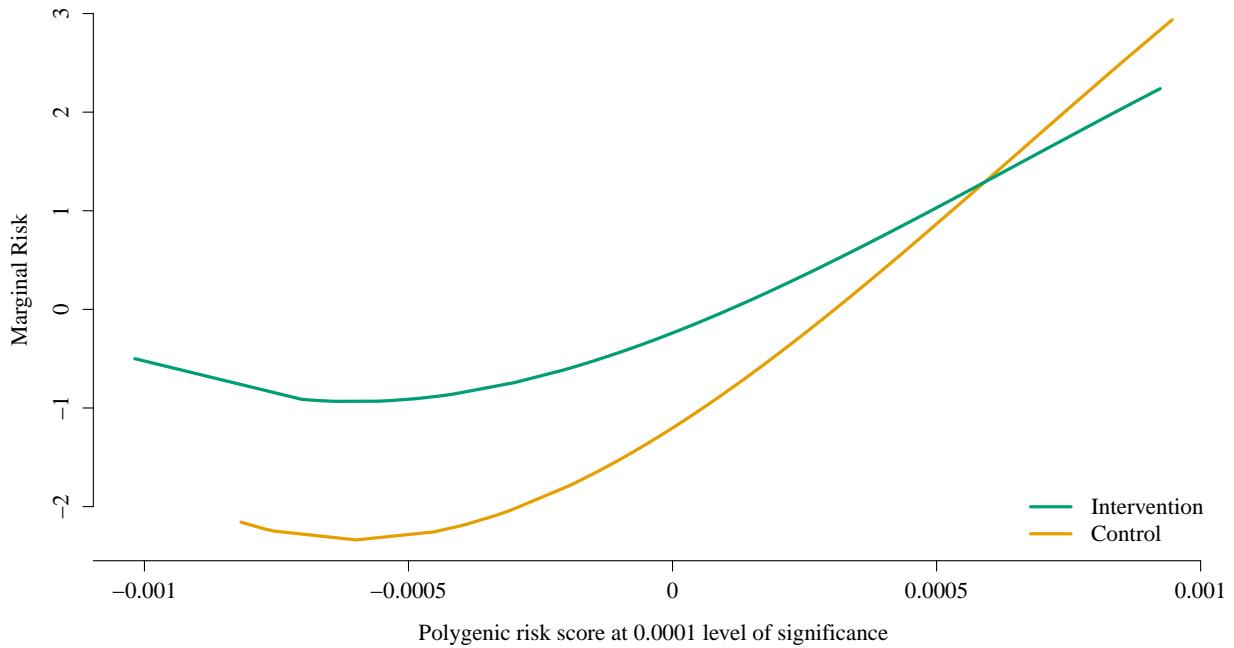
### 5.1. Gene-environment interactions in the Nurse Family Partnership program

It is well known that environmental exposures can have an important impact on academic achievement. Indeed, early intervention in young children has been shown to positively impact intellectual abilities (Campbell and Ramey, 1994). More recent studies have shown that cognitive performance, a trait that measures the ability to learn, reason and solve problems, is also strongly influenced by genetic factors. Genome-wide association studies (GWAS) suggest that 20% of the variance in educational attainment (years of education) may be accounted for by common genetic variation (Rietveld et al., 2013; Okbay, Beauchamp, Fontana, Lee, Pers, Rietveld, Turley, Chen, Emilsson, Meddins et al., 2016). Unsurprisingly, there is significant overlap in the SNPs that predict educational attainment and measures of cognitive function. An interesting query that arises is how the environment interacts with these genetics variants to predict measures of cognitive function. To address this question, we analyzed data from the Nurse Family Partnership (NFP), a psychosocial intervention program that begins in pregnancy and targets maternal health, parenting and mother-infant interactions (Olds et al., 1998). The Stanford Binet IQ scores at 4 years of age were collected for 189 subjects (including 19 imputed using mice (Buuren and Groothuis-Oudshoorn, 2010)) born to women randomly assigned to control ( $n = 100$ ) or nurse-visited intervention groups ( $n = 89$ ). For each subject, we calculated a polygenic risk score (PRS) for educational attainment at different p-value thresholds using weights from the GWAS conducted in Okbay et al. (Okbay et al., 2016). In this context, individuals with a higher PRS have a propensity for higher educational attainment. The goal of this analysis was to determine if there was an interaction between genetic predisposition to educational attainment ( $X$ ) and maternal participation in the NFP program ( $E$ ) on child IQ at 4 years of age ( $Y$ ). We applied the weak heredity sail with cubic B-splines and  $\alpha = 0.1$  to encourage interactions, and selected the optimal tuning parameter using 10-fold cross-validation. Our method identified an interaction between the intervention and PRS which included genetic variants at the 0.0001 level of significance. This interaction is shown in Figure 3. We see that the intervention has a much larger effect on IQ for lower PRS compared to a higher PRS. In other words, perinatal home visitation by nurses can impact IQ scores in children who are genetically predisposed to lower educational attainment. Similar results were obtained for the other imputed datasets (Supplemental Section D). We also compared sail with two other interaction selection methods, lassoBT and GLinternet with default settings, on 200 bootstrap samples of the data. The average and standard deviation of the MSE and size of the active set ( $|\hat{J}|$ ) across the 200 bootstrap samples are given in Table 3. We see that sail tends to select sparser models while maintaining similar prediction performance compared to lassoBT. The GLinternet statistics are omitted here since the algorithm did not converge for many of the 200 simulations.

### 5.2. Study to Understand Prognoses Preferences Outcomes and Risks of Treatment

The Study to Understand Prognoses Preferences Outcomes and Risks of Treatment (SUPPORT) aimed at identifying which clinical variables influence medium-term (half-year) mortality rate amongst seriously ill hospitalized patients and improving clinical decision making (Connors et al., 1995). With a relatively large sample size of 9,105 and detailed documentation of clinical variables, the SUPPORT dataset allows detection of potential interactions using the strategy implemented in sail. We applied sail to test for non-linear interactions between acute renal failure or multiple organ system failure (ARF/MOSF), an important predictor for survival rate, and 13 other variables that were deemed clinically relevant. These variables included the number of comorbidities (excluding ARF/MOSF), age, sex, as well as multiple physiological and blood biochemical indices. The response was whether a patient survived after six months since hospitalization.

A total of 8,873 samples had complete data on all variables of interest. We randomly divided these samples into equal sized training/validation/test splits and ran lassoBT, GLinternet, and the weak heredity sail with cubic B-splines and  $\alpha = 0.1$  (as was done in the Nurse Family Partnership program case study). A binomial distribution family was specified for GLinternet, whereas lassoBT had the same default settings as the simulation study since it did not support a specialized implementation for binary outcomes. We again ran each method on the training data, determined the optimal tuning parameter on the validation data based on the area under the receiver operating characteristic curve (AUC), and assessed AUC on the test data. We repeated this process 200 times and report the results in Table 3.



**Figure 3:** Estimated interaction effect identified by the weak heredity `sail` using cubic B-splines and  $\alpha = 0.1$  for the Nurse Family Partnership data. The selected model, chosen via 10-fold cross-validation, contained three variables: the main effects for the intervention and the PRS for educational attainment using genetic variants significant at the 0.0001 level, as well as their interaction.

**Table 3**

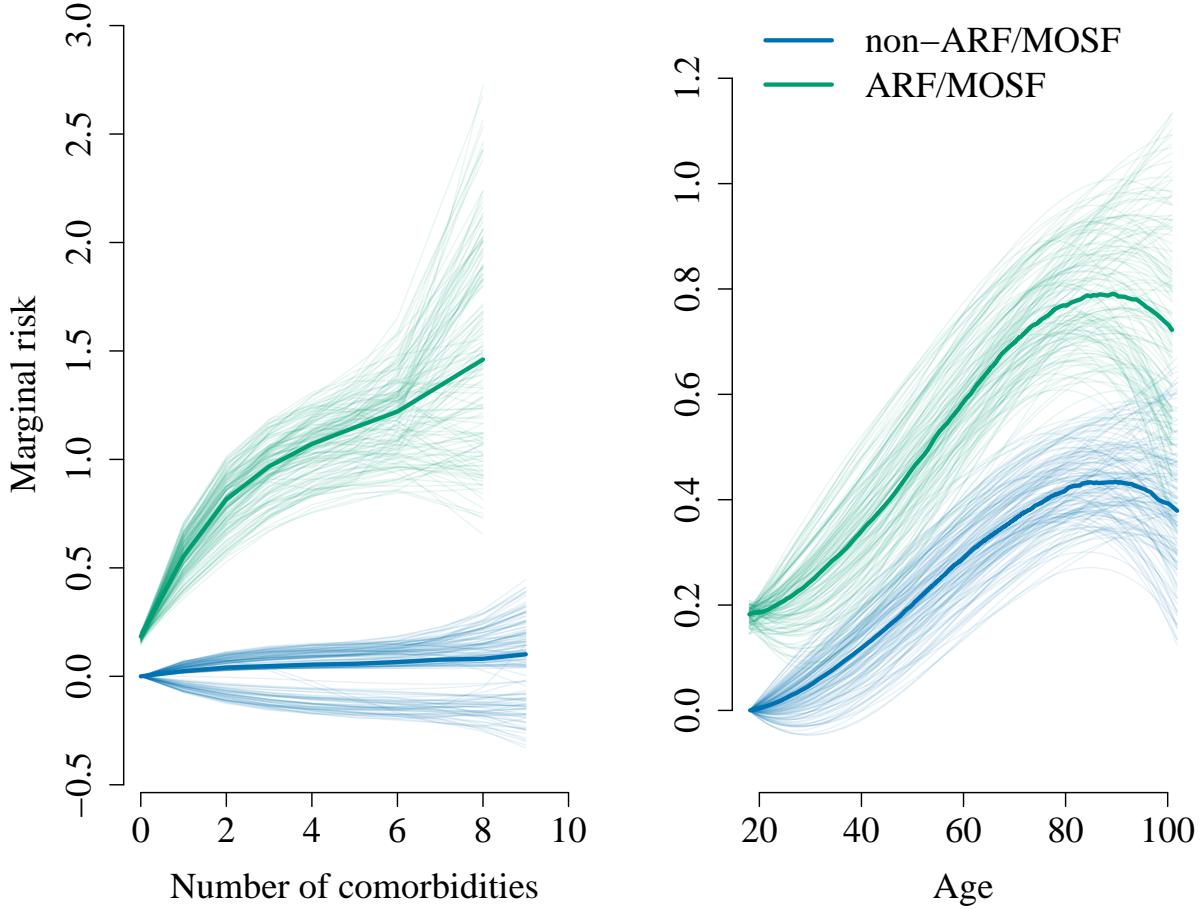
Comparison of analytic methods for selecting interactions using the Nurse Family Partnership program and the SUPPORT datasets. Averages (standard deviations in parentheses) are based on 200 bootstrap samples.  $|\hat{J}|$  is the number of variables selected by the method. GLinternet results not reported for NFP data since the algorithm did not converge in many of the bootstrap samples.

Method	Nurse Family Partnership		SUPPORT	
	Mean Squared Error	$ \hat{J} $	AUC	$ \hat{H} $
<code>sail</code>	3.5 (0.6)	4 (3)	0.66 (0.01)	25 (3)
<code>lassoBT</code>	3.53 (0.477)	11 (6)	0.65 (0.009)	49 (14)
<code>GLinternet</code>	–	–	0.65 (0.009)	58 (7)

We found that `sail` achieved similar prediction accuracy to `lassoBT` and `GLinternet`. However, the predictive performance of `lassoBT` and `GLinternet` relied on models which included many more variables. In Figure 4, we visualize the two strongest interaction effects associated with the number of comorbidities and age, respectively. For those having undergone ARF/MOSF, an increased number of comorbidities decreases their chance of survival, while there seems to be no such relationship for non-ARF/MOSF patients. The interaction between ARF/MOSF and age shows the risk incurred by ARF/MOSF is most distinguishing among patients between the ages of 70 and 80.

## 6. Discussion

In this article we have introduced the sparse additive interaction learning model `sail` for detecting non-linear interactions with a key environmental or exposure variable in high-dimensional settings. Using a simple reparametrization, we are able to achieve either the weak or strong heredity property without using a complex penalty function. We



**Figure 4:** Illustration of estimated interaction effects identified by sail for the SUPPORT data. Median prediction curves in dark colors based on 200 train/validate/test splits represent the estimated marginal interaction effects. Coefficients estimated in each of the 200 train/validate/test splits were used to generate prediction curves representing a 90% confidence interval colored in corresponding light colors.

developed a blockwise coordinate descent algorithm to solve the sail objective function for the least-squares loss. We further studied the asymptotic properties of our method and showed that under certain conditions, it possesses the oracle property. All our algorithms have been implemented in a computationally efficient, well-documented and freely available R package on CRAN. Furthermore, our method is flexible enough to handle any type of basis expansion including the identity map, which allows for linear interactions. Our implementation allows the user to selectively apply the basis expansions to the predictors, allowing for example, a combination of continuous and categorical predictors. An extensive simulation study shows that sail, adaptive sail and sail weak outperform existing penalized regression methods in terms of prediction accuracy, sensitivity and specificity when there are non-linear main effects only, as well as interactions with an exposure variable. We then demonstrated the utility of our method to identify non-linear interactions in both biological and epidemiological data. In the NFP program, we showed that individuals who are genetically predisposed to lower educational attainment are those who stand to benefit the most from the intervention. Analysis of the SUPPORT data revealed that those having undergone ARF/MOSF, an increased number of comorbidities decreased their chances of survival, while there seemed to be no such relationship for non-ARF/MOSF patients. In a bootstrap analysis of both datasets, we observed that sail tended to select sparser models while maintaining similar prediction performance compared to other interaction selection methods.

Our method however does have its limitations. sail can currently only handle  $X_E \cdot f(X)$  or  $f(X_E) \cdot X$  and

does not allow for  $f(X, X_E)$ , i.e., only one of the variables in the interaction can have a non-linear effect and we do not consider the tensor product. The reparametrization leads to a non-convex optimization problem which makes convergence rates difficult to assess, though we did not experience any major convergence issues in our simulations and real data analysis. The memory footprint can also be an issue depending on the degree of the basis expansion and the number of variables. Furthermore, the functional form of the covariate effects is treated as known in our method. Being able to automatically select for example, linear vs. nonlinear components, is currently an active area of research in main effects models (Haris et al., 2019). To our knowledge, our proposal is the first to allow for non-linear interactions with a key exposure variable following the weak or strong heredity property in high-dimensional settings. We also provide a first software implementation for these models.

## Acknowledgments

SRB and CMTG were supported by the Ludmer Centre for Neuroinformatics and Mental Health and the Canadian Institutes for Health Research PJT 148620. SRB acknowledges the support of the Natural Sciences and Engineering Research Council of Canada (NSERC), RGPIN-2020-05133. This research was enabled in part by support provided by Calcul Québec ([www.calculquebec.ca](http://www.calculquebec.ca)) and Compute Canada ([www.computecanada.ca](http://www.computecanada.ca)). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

## References

- Bach, F., Jenatton, R., Mairal, J., Obozinski, G., et al., 2012. Structured sparsity through convex optimization. *Statistical Science* 27, 450–468.
- Bhatnagar, S.R., Yang, Y., Khundrakpam, B., Evans, A.C., Blanchette, M., Bouchard, L., Greenwood, C.M., 2018. An analytic approach for interpretable predictive models in high-dimensional data in the presence of interactions with exposures. *Genetic epidemiology* 42, 233–249.
- Bien, J., Taylor, J., Tibshirani, R., et al., 2013. A lasso for hierarchical interactions. *The Annals of Statistics* 41, 1111–1141.
- Bühlmann, P., Van De Geer, S., 2011. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media.
- Buuren, S.v., Groothuis-Oudshoorn, K., 2010. mice: Multivariate imputation by chained equations in r. *Journal of statistical software* , 1–68.
- Campbell, F.A., Ramey, C.T., 1994. Effects of early intervention on intellectual and academic achievement: a follow-up study of children from low-income families. *Child development* 65, 684–698.
- Chipman, H., 1996. Bayesian variable selection with related predictors. *Canadian Journal of Statistics* 24, 17–36.
- Choi, N.H., Li, W., Zhu, J., 2010. Variable selection with the strong heredity constraint and its oracle property. *Journal of the American Statistical Association* 105, 354–364.
- Chouldechova, A., Hastie, T., 2015. Generalized additive model selection. arXiv preprint arXiv:1506.03850 .
- Connors, A.F., Dawson, N.V., Desbiens, N.A., Fulkerson, W.J., Goldman, L., Knaus, W.A., Lynn, J., Oye, R.K., Bergner, M., Damiano, A., et al., 1995. A controlled trial to improve care for seriously ill hospitalized patients: The study to understand prognoses and preferences for outcomes and risks of treatments (support). *Jama* 274, 1591–1598.
- Conway, J.R., Lex, A., Gehlenborg, N., 2017. UpSetR: an R package for the visualization of intersecting sets and their properties. *Bioinformatics* 33, 2938–2940. URL: <https://doi.org/10.1093/bioinformatics/btx364>, doi:10.1093/bioinformatics/btx364, arXiv:<https://academic.oup.com/bioinformatics/article-pdf/33/18/2938/25164302/btx364.pdf>.
- Cox, D.R., 1984. Interaction. *International Statistical Review/Revue Internationale de Statistique* , 1–24.
- Fan, J., Han, F., Liu, H., 2014. Challenges of big data analysis. *National science review* 1, 293–314.
- Fan, J., Li, R., 2001. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association* 96, 1348–1360.
- Friedman, J., Hastie, T., Tibshirani, R., 2010. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software* 33, 1.
- Hao, N., Feng, Y., Zhang, H.H., 2018. Model selection for high-dimensional quadratic regression via regularization. *Journal of the American Statistical Association* 113, 615–625.
- Haris, A., Shojaie, A., Simon, N., 2019. Nonparametric regression with adaptive truncation via a convex hierarchical penalty. *Biometrika* 106, 87–107.
- Haris, A., Witten, D., Simon, N., 2016. Convex modeling of interactions with strong heredity. *Journal of Computational and Graphical Statistics* 25, 981–1004.
- Hastie, T., Tibshirani, R., Wainwright, M., 2015. *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press.
- Huang, J., Horowitz, J.L., Wei, F., 2010. Variable selection in nonparametric additive models. *Annals of statistics* 38, 2282–2313.
- Lim, M., Hastie, T., 2015. Learning interactions via hierarchical group-lasso regularization. *Journal of Computational and Graphical Statistics* 24, 627–654.
- Lin, Y., Zhang, H.H., et al., 2006. Component selection and smoothing in multivariate nonparametric regression. *The Annals of Statistics* 34, 2272–2297.
- McCullagh, P., Nelder, J.A., 1989. *Generalized linear models*. volume 37. CRC press.
- Nardi, Y., Rinaldo, A., et al., 2008. On the asymptotic properties of the group lasso estimator for linear models. *Electronic Journal of Statistics* 2, 605–633.

## A Sparse Additive Model for High-Dimensional Interactions with an Exposure Variable

- Okbay, A., Beauchamp, J.P., Fontana, M.A., Lee, J.J., Pers, T.H., Rietveld, C.A., Turley, P., Chen, G.B., Emilsson, V., Meddens, S.F.W., et al., 2016. Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* 533, 539.
- Olds, D., Henderson Jr, C.R., Cole, R., Eckenrode, J., Kitzman, H., Luckey, D., Pettitt, L., Sidora, K., Morris, P., Powers, J., 1998. Long-term effects of nurse home visitation on children's criminal and antisocial behavior: 15-year follow-up of a randomized controlled trial. *Jama* 280, 1238–1244.
- R Core Team, 2017. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria. URL: <https://www.R-project.org/>.
- Radchenko, P., James, G.M., 2010. Variable selection using adaptive nonlinear interaction structures in high dimensions. *Journal of the American Statistical Association* 105, 1541–1553.
- Ravikumar, P., Lafferty, J., Liu, H., Wasserman, L., 2009. Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 71, 1009–1030.
- Rietveld, C.A., Medland, S.E., Derringer, J., Yang, J., Esko, T., Martin, N.W., Westra, H.J., Shakhsbazov, K., Abdellaoui, A., Agrawal, A., et al., 2013. Gwas of 126,559 individuals identifies genetic variants associated with educational attainment. *science* 340, 1467–1471.
- Shah, R.D., 2016. Modelling interactions in high-dimensional data with backtracking. *Journal of Machine Learning Research* 17, 1–31.
- She, Y., Jiang, H., 2014. Group regularized estimation under structural hierarchy. *arXiv preprint arXiv:1411.4691* .
- Tibshirani, R., 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 267–288.
- Wang, H., Li, G., Tsai, C.L., 2007. Regression coefficient and autoregressive order shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 69, 63–78.
- Yang, Y., Zou, H., 2015. A fast unified algorithm for solving group-lasso penalize learning problems. *Statistics and Computing* 25, 1129–1141.
- Zhao, P., Rocha, G., Yu, B., 2009. The composite absolute penalties family for grouped and hierarchical variable selection. *The Annals of Statistics* , 3468–3497.
- Zou, H., 2006. The adaptive lasso and its oracle properties. *Journal of the American statistical association* 101, 1418–1429.
- Zou, H., Zhang, H.H., 2009. On the adaptive elastic-net with a diverging number of parameters. *Annals of statistics* 37, 1733–1751.

$\beta_E^*$	$\theta_1^{*\top}$	$\theta_2^{*\top}$	...	$\theta_p^{*\top}$	$\gamma_{1E}^*$	$\gamma_{2E}^*$	...	$\gamma_{pE}^*$
$\phi_1^{*\top}$	$\phi_2^{*\top}$	$\phi_3^{*\top}$	...	$\phi_{p+1}^{*\top}$	$\phi_{p+2}^{*\top}$	$\phi_{p+3}^{*\top}$	...	$\phi_{2p+1}^{*\top}$
$\lambda(1-\alpha)w_E$	$\lambda(1-\alpha)w_2$	$\lambda(1-\alpha)w_3$	...	$\lambda(1-\alpha)w_{p+1}$	$\lambda\alpha w_{p+2,E}$	$\lambda\alpha w_{p+3,E}$	...	$\lambda\alpha w_{2p+1,E}$
$\lambda_1$	$\lambda_2$	$\lambda_3$	...	$\lambda_{p+1}$	$\lambda_{p+2}$	$\lambda_{p+3}$	...	$\lambda_{2p+1}$

**Table 4**

Correspondence between parameters used to simplify the notation in the proofs. The first row shows the actual parameters used in the loss function. The second row shows the corresponding parameters in the simplified notation. The third row shows the actual tuning parameters used in the penalty function. The fourth row shows the corresponding tuning parameters in the simplified notation. This correspondence greatly simplifies the notation used in the proofs.

## A. Proofs

As shown in the main text, we simplified the notation to make the proofs easier to follow. We summarize the original notation and the corresponding simplified notation in Table 4. This notation then allows us to write down the sail estimates as

$$\hat{\Phi}_n = \arg \min_{\Phi} Q_n(\Phi) = -L_n(\Phi) + n\lambda_m \sum_{m=1}^{2p+1} \|\phi_m\|_2, \quad (17)$$

### A.1. Regularity Conditions

(C1) The observation  $\{\mathbf{V}_i : i = 1, \dots, n\}$  are independent and identically distributed with a probability density  $f(\mathbf{V}, \Phi)$ , which has a common support. We assume the density  $f$  satisfies the following equations:

$$E_{\Phi} \left[ \nabla_{\phi_j} \log f(\mathbf{V}, \Phi) \right] = \mathbf{0} \quad \text{for } j = 1, \dots, 2p+1.$$

and

$$\begin{aligned} \mathbf{I}_{j_1 k_1 j_2 k_2}(\Phi) &= E_{\Phi} \left[ \frac{\partial}{\partial \phi_{j_1 k_1}} \log f(\mathbf{V}, \Phi) \cdot \frac{\partial}{\partial \phi_{j_2 k_2}} \log f(\mathbf{V}, \Phi) \right] \\ &= E_{\Phi} \left[ -\frac{\partial^2}{\partial \phi_{j_1 k_1} \partial \phi_{j_2 k_2}} \log f(\mathbf{V}, \Phi) \right], \end{aligned}$$

for any  $j_1, j_2 = 1, \dots, 2p+1$ , and  $k_1 = 1, \dots, p_{j_1}$ ,  $k_2 = 1, \dots, p_{j_2}$ , where  $j_1, j_2$  are the index of group,  $k_1, k_2$  be the index of elements within the corresponding group,  $p_{j_1}, p_{j_2}$  are the group size of  $j_1, j_2$  respectively.

(C2) The Fisher information matrix

$$\mathbf{I}(\Phi) = E \left[ \left( \frac{\partial}{\partial \Phi} \log f(\mathbf{V}, \Phi) \right) \left( \frac{\partial}{\partial \Phi} \log f(\mathbf{V}, \Phi) \right)^{\top} \right]$$

is finite and positive definite at  $\Phi = \Phi^*$ .

(C3) There exists an open set  $\omega$  of  $\Omega$  that contains the true parameter point  $\Phi^*$  such that for almost all  $\mathbf{V}$  the density  $f(\mathbf{V}, \Phi)$  admits all third derivatives  $\frac{\partial^3 f(\mathbf{V}, \Phi)}{\partial \phi_{j_1 k_1} \partial \phi_{j_2 k_2} \partial \phi_{j_3 k_3}}$  for all  $\Phi$  in  $\omega$  and any  $j_1, j_2, j_3 = 1, \dots, 2p+1$ , and  $k_1 = 1, \dots, p_{j_1}$ ,  $k_2 = 1, \dots, p_{j_2}$  and  $k_3 = 1, \dots, p_{j_3}$ . Furthermore, there exist functions  $M_{j_1 k_1 j_2 k_2 j_3 k_3}$  such that

$$\left| \frac{\partial^3}{\partial \phi_{j_1 k_1} \partial \phi_{j_2 k_2} \partial \phi_{j_3 k_3}} \log f(\mathbf{V}, \Phi) \right| \leq M_{j_1 k_1 j_2 k_2 j_3 k_3}(\mathbf{V}) \quad \text{for all } \Phi \in \omega,$$

and  $m_{j_1 k_1 j_2 k_2 j_3 k_3} = E_{\Phi^*}[M_{j_1 k_1 j_2 k_2 j_3 k_3}(\mathbf{V})] < \infty$ .

## A.2. Lemma 1 proof

Let  $\eta_n = \frac{1}{\sqrt{n}} + a_n$  and  $\{\Phi^* + \eta_n \delta : \|\delta\|_2 \leq C\}$  be the ball around  $\Phi^*$  for  $\delta \in \mathbb{R}^d$ , where  $d$  is the dimension of the design matrix and  $C$  is some constant. Under the regularity assumptions, we show that there exists a local minimizer  $\hat{\Phi}_n$  of  $Q_n(\Phi)$  such that  $\|\hat{\Phi}_n - \Phi^*\|_2 = O_p(\frac{1}{\sqrt{n}})$ . For this proof, we adopt the approaches outlined in (Fan and Li, 2001; Choi et al., 2010; Nardi et al., 2008; Wang et al., 2007) and extend it to our situation. Let  $\eta_n = \frac{1}{\sqrt{n}} + a_n$  and  $\{\Phi^* + \eta_n \delta : \|\delta\|_2 \leq C\}$  be the ball around  $\Phi^*$  for  $\delta = (\mathbf{u}_1^\top, \mathbf{u}_2^\top, \dots, \mathbf{u}_{p+1}^\top, \mathbf{u}_{p+2}^\top, \dots, \mathbf{u}_{2p+1}^\top)^\top \in \mathbb{R}^d$ , where  $d$  is the dimension of the design matrix and  $C$  is some constant. The objective function is given by

$$Q_n(\Phi) = -L_n(\Phi) + n\lambda_m \sum_{m=1}^{2p+1} \|\phi_m\|_2,$$

Define

$$D_n(\delta) \equiv Q_n(\Phi^* + \eta_n \delta) - Q_n(\Phi^*).$$

Then for  $\delta$  that satisfies  $\|\delta\|_2 = C$ , we have

$$\begin{aligned} D_n(\delta) &= -L_n(\Phi^* + \eta_n \delta) + L_n(\Phi^*) + n \sum_{m=1}^{2p+1} \lambda_m (\|\theta_m^* + \eta_n \mathbf{u}_m\|_2 - \|\theta_m^*\|_2) \\ &\stackrel{(a)}{\geq} -L_n(\Phi^* + \eta_n \delta) + L_n(\Phi^*) + n \sum_{m \in \mathcal{A}_1} \lambda_m (\|\theta_m^* + \eta_n \mathbf{u}_m\|_2 - \|\theta_m^*\|_2) \\ &\quad + n \sum_{m \in \mathcal{A}_2} \lambda_m (\|\theta_m^* + \eta_n \mathbf{u}_m\|_2 - \|\theta_m^*\|_2) \\ &\stackrel{(b)}{\geq} -L_n(\Phi^* + \eta_n \delta) + L_n(\Phi^*) - n\eta_n \sum_{m \in \mathcal{A}_1} \lambda_m \|\mathbf{u}_m\|_2 - n\eta_n \sum_{m \in \mathcal{A}_2} \lambda_m \|\mathbf{u}_m\|_2 \\ &\stackrel{(c)}{\geq} -L_n(\Phi^* + \eta_n \delta) + L_n(\Phi^*) - n\eta_n^2 \sum_{m \in \mathcal{A}_1} \|\mathbf{u}_m\|_2 - n\eta_n^2 \sum_{m \in \mathcal{A}_2} \|\mathbf{u}_m\|_2 \\ &\geq -L_n(\Phi^* + \eta_n \delta) + L_n(\Phi^*) - n\eta_n^2 (|\mathcal{A}_1| + |\mathcal{A}_2|)C \\ &\stackrel{(d)}{=} -[\nabla L_n(\Phi^*)]^\top (\eta_n \delta) - \frac{1}{2} (\eta_n \delta)^\top [\nabla^2 L_n(\Phi^*)] (\eta_n \delta) (1 + o(1)) \\ &\quad - n\eta_n^2 (|\mathcal{A}_1| + |\mathcal{A}_2|)C \end{aligned} \tag{18}$$

Inequality (a) is by the fact that  $\sum_{m \notin \mathcal{A}_1} \|\phi_m^*\|_2 = 0$  and  $\sum_{m \notin \mathcal{A}_2} \|\phi_m^*\|_2 = 0$ . Inequality (b) is due to the reverse triangle inequality  $\|a\|_2 - \|b\|_2 \geq -\|a - b\|_2$ . Inequality (c) is by  $\lambda_m \leq a_n \leq \eta_n$  for  $m \in \mathcal{A}_1$  and  $m \in \mathcal{A}_2$ . Equality (d) is by the standard argument on the Taylor expansion of the loss function:

$$\begin{aligned} L_n(\Phi^* + \eta_n \delta) &= L_n(\Phi^* + \eta_n \cdot \mathbf{0}) + \eta_n \nabla L_n(\Phi^* + \eta_n \cdot \mathbf{0})^\top (\delta - \mathbf{0}) \\ &\quad + \frac{1}{2} (\delta - \mathbf{0})^\top \nabla^2 L_n(\Phi^* + \eta_n \cdot \mathbf{0}) (\delta - \mathbf{0}) \{1 + o(1)\} \\ &= L_n(\Phi^*) + \eta_n \nabla L_n(\Phi^*)^\top \delta + \frac{1}{2} \delta^\top \nabla^2 L_n(\Phi^*) \delta \eta_n^2 \{1 + o(1)\} \end{aligned}$$

We split (18) into three parts:

$$\begin{aligned} D_1 &= -[\nabla L_n(\Phi^*)]^\top (\eta_n \delta) \\ D_2 &= -\frac{1}{2} (\eta_n \delta)^\top [\nabla^2 L_n(\Phi^*)] (\eta_n \delta) (1 + o(1)) \\ D_3 &= -n\eta_n^2 (|\mathcal{A}_1| + |\mathcal{A}_2|)C \end{aligned}$$

Then

$$\begin{aligned}
D_1 &= -\eta_n [\nabla L_n(\Phi^*)]^\top \delta \\
&= -\sqrt{n}\eta_n \left( \frac{1}{\sqrt{n}} \nabla L_n(\Phi^*) \right)^\top \delta \\
&= -\sqrt{n}\eta_n \left( \sqrt{n} \frac{1}{n} \sum_{i=1}^n \nabla \log f(V_i, \Phi) \Big|_{\Phi=\Phi^*} \right)^\top \delta \\
&= -\sqrt{n}\eta_n \left( \sqrt{n} \left[ \frac{1}{n} \sum_{i=1}^n \nabla \log f(V_i, \Phi) \Big|_{\Phi=\Phi^*} - \mathbf{0} \right] \right)^\top \delta \\
&= -\sqrt{n}\eta_n \left( \sqrt{n} \left[ \frac{1}{n} \sum_{i=1}^n \nabla \log f(V_i, \Phi) \Big|_{\Phi=\Phi^*} - E_{\Phi^*} \nabla L(\Phi^*) \right] \right)^\top \delta \\
&= -\sqrt{n}\eta_n O_P(1) \delta \\
&= -O_P(n\eta_n^2) \delta
\end{aligned} \tag{19}$$

The last equation is by  $a_n = o(\frac{1}{\sqrt{n}})$  and

$$\begin{aligned}
O_P(n\eta_n^2) &= O_P(n(n^{-1/2} + a_n)^2) = O_P(1 + 2n^{1/2}a_n + na_n^2) \\
&= O_P(1 + n^{1/2}a_n + (n^{1/2}a_n)^2) = O_P(1 + n^{1/2}a_n + o(1)) \\
&= O_p(n^{1/2}(n^{-1/2} + a_n)) = O_p(n^{1/2}\eta_n)
\end{aligned}$$

$$\begin{aligned}
D_2 &= \frac{1}{2}n\eta_n^2 \left\{ \delta^\top \left[ -\frac{1}{n} \nabla^2 L_n(\Phi^*) \right] \delta \right\} (1 + o_p(1)) \\
&= \frac{1}{2}n\eta_n^2 \left\{ \delta^\top [I(\Phi^*)] \delta \right\} (1 + o_p(1)) \text{ by the weak law of large numbers.} \\
&= O_p(n\eta_n^2 \|\delta\|_2^2)
\end{aligned} \tag{20}$$

Combining (19) and (20) with (18) gives:

$$\begin{aligned}
D_n(\delta) &\geq D_1 + D_2 + D_3 \\
&= -O_P(n\eta_n^2) \delta + O_p(n\eta_n^2 \|\delta\|_2^2) - n\eta_n^2(|\mathcal{A}_1| + |\mathcal{A}_2|)C
\end{aligned}$$

We can see that the first term  $D_1$  is linear in  $\delta$  and the second term  $D_2$  is quadratic in  $\delta$ . We can conclude that for a large enough constant  $C = \|\delta\|_2$ ,  $D_2$  dominates  $D_1$  and  $D_3$ . Note that this is a positive term since  $I(\Phi)$  is positive definite at  $\Phi = \Phi^*$  by regularity condition (C2). Therefore, for each  $\varepsilon > 0$ , there exists a large enough constant  $C$  such that, for large enough  $n$

$$P \left\{ \inf_{\|\delta\|_2=C} D_n(\delta) > 0 \right\} \geq 1 - \varepsilon$$

This implies with probability at least  $1 - \varepsilon$  that the empirical likelihood  $Q_n$  has a local minimizer in the ball  $\{\Phi^* + \eta_n \delta : \|\delta\|_2 \leq C\}$  (since  $Q_n$  is bounded and  $\{\Phi^* + \alpha_n \delta : \|\delta\|_2 \leq C\}$  is closed). In other words, there exists a local solution  $\hat{\Phi}_n$  such that  $\|\hat{\Phi}_n - \Phi^*\| \leq \eta_n \|\delta\|_2 \leq \eta_n C = O_P(\eta_n) = O_P(\frac{1}{\sqrt{n}} + a_n) = O_p(\frac{1}{\sqrt{n}})$ , since  $a_n = o(\frac{1}{\sqrt{n}})$ . Hence,

$$\|\hat{\Phi}_n - \Phi^*\|_2 = O_P\left(\frac{1}{\sqrt{n}}\right). \square$$

### A.3. Theorem 1 proof

We first consider consistency for the main effects  $P(\hat{\Phi}_{\mathcal{A}_1^c} = \mathbf{0}) \rightarrow 1$ . Following (Fan and Li, 2001; Choi et al., 2010), it is sufficient to show that for all  $m \in \mathcal{A}_1^c$ ,  $P(\hat{\phi}_m = \mathbf{0}) \rightarrow 1$ , which implies that  $P(\hat{\Phi}_{\mathcal{A}_1^c} = \mathbf{0}) \rightarrow 1$ , i.e., the  $\sqrt{n}$ -consistent estimate  $\hat{\Phi}$  has oracle property  $\hat{\phi}_m = \mathbf{0}$  if  $\phi_m^* = \mathbf{0}$ . Denote

$$\hat{\phi}_m = (\hat{\phi}_{m1}, \dots, \hat{\phi}_{mp_m}),$$

where  $p_m$  is the group size of  $\hat{\phi}_m$ . Let  $\hat{\phi}_{mk}$  be the  $k$ -th entry of  $\hat{\phi}_m$ . Note that if  $\hat{\phi}_m \neq \mathbf{0}$ , then  $\hat{\phi}_{mk} \neq 0$  for  $k = 1, \dots, p_m$ , then penalty function  $\|\hat{\phi}_m\|_2$  becomes differentiable. Therefore  $\phi_{mk}$  for  $k = 1, \dots, p_m$  must satisfy the following normal equation

$$\begin{aligned} \frac{\partial Q_n(\hat{\Phi}_n)}{\partial \phi_{mk}} &= -\frac{\partial L_n(\hat{\Phi}_n)}{\partial \phi_{mk}} + n\lambda_m \frac{\hat{\phi}_{mk}}{\|\hat{\phi}_m\|_2} \\ &= -\frac{\partial L_n(\Phi^*)}{\partial \phi_{mk}} - \sum_{j_1=1}^{2p+1} \sum_{k_1=1}^{p_{j_1}} \frac{\partial^2 L_n(\Phi^*)}{\partial \phi_{mk} \partial \phi_{j_1 k_1}} (\hat{\phi}_{j_1 k_1} - \phi_{j_1 k_1}^*) \\ &\quad - \frac{1}{2} \sum_{j_1=1}^{2p+1} \sum_{k_1=1}^{p_{j_1}} \sum_{j_2=1}^{2p+1} \sum_{k_2=1}^{p_{j_2}} \frac{\partial^3 L_n(\tilde{\Phi})}{\partial \phi_{mk} \partial \phi_{j_1 k_1} \partial \phi_{j_2 k_2}} (\hat{\phi}_{j_1 k_1} - \phi_{j_1 k_1}^*) (\hat{\phi}_{j_2 k_2} - \phi_{j_2 k_2}^*) \\ &\quad + n\lambda_m \frac{\hat{\phi}_{mk}}{\|\hat{\phi}_m\|_2} \triangleq I_1 + I_2 + I_3 + I_4 = 0 \end{aligned}$$

where  $\tilde{\Phi}$  lies between  $\hat{\Phi}_n$  and  $\Phi^*$ . By the regularity conditions and Lemma (1) that  $\|\hat{\Phi}_n - \Phi^*\|_2 = O_P\left(\frac{1}{\sqrt{n}}\right)$ , the first term is of the order  $O_p(\sqrt{n})$

$$I_1 = -\frac{\partial L_n(\hat{\Phi}_n)}{\partial \phi_{mk}} = -\sqrt{n} \sqrt{n} \frac{1}{n} \frac{\partial L_n(\hat{\Phi}_n)}{\partial \phi_{mk}} = \sqrt{n} O_p(1) = O_p(\sqrt{n}).$$

Then the second is of the order  $O_P\left(\frac{1}{\sqrt{n}}\right)$  and the third term is of the order  $O_P\left(\frac{1}{n}\right)$ . Hence

$$\frac{\partial Q_n(\hat{\Phi}_n)}{\partial \Phi_m} = \sqrt{n} \left\{ O_p(1) + \sqrt{n} \lambda_m \frac{\hat{\phi}_{mk}}{\|\hat{\phi}_m\|_2} \right\}. \quad (21)$$

As  $\sqrt{n}\lambda_m \geq \sqrt{nb_n} \rightarrow \infty$  for  $m \in \mathcal{A}_1^c$  from the assumption, therefore we know that  $I_4$  dominates  $I_1, I_2$  and  $I_3$  in (21) with probability tending to one. This means that (21) cannot be true as long as the sample size is sufficiently large. As a result, we can conclude that with probability tending to one, the estimate  $\hat{\phi}_m = (\hat{\phi}_{m1}, \dots, \hat{\phi}_{mp_m})$  must be in a position where  $\hat{\phi}_m$  is not differentiable. Hence  $\hat{\phi}_m = \mathbf{0}$  for all  $m \in \mathcal{A}_1^c$ . Hence  $P(\hat{\Phi}_{\mathcal{A}_1^c} = \mathbf{0}) \rightarrow 1$ . This completes the proof.

Next, we prove that for the interactions  $P(\hat{\Phi}_{\mathcal{A}_2^c} = \mathbf{0}) \rightarrow 1$ . For  $m \in \mathcal{A}_2^c$  s.t.  $\phi_m^* = \gamma_{jE}^* = 0$  but  $\beta_E \neq 0$  and  $\theta_j^* \neq \mathbf{0}$  ( $1 \leq j \leq p$ ), we can prove  $P(\hat{\Phi}_{\mathcal{A}_2^c} = \mathbf{0}) \rightarrow 1$  by a similar reasoning, which further implies that  $P(\hat{\gamma}_{jE} = 0) \rightarrow 0$ . For  $m \in \mathcal{A}_2^c$  such that  $\phi_m^* = \gamma_{jE}^* = 0$  and either  $\beta_E = 0$  or  $\theta_j^* = \mathbf{0}$  ( $1 \leq j \leq p$ ): without loss of generality, assume that  $\theta_j^* = \mathbf{0}$ . Notice that  $\hat{\theta}_j = \mathbf{0}$  implies  $\hat{\gamma}_{jE} = 0$ , since if  $\hat{\gamma}_{jE} \neq 0$ , the value of the loss function does not change but the value of the penalty function will increase. Because we already prove  $P(\hat{\Phi}_{\mathcal{A}_1^c} = \mathbf{0}) \rightarrow 1$ , therefore we get  $P(\hat{\Phi}_{\mathcal{A}_2^c} = \mathbf{0}) \rightarrow 1$  as well for this case.  $\square$

#### A.4. Theorem 2 proof

By Lemma 1 and Theorem 1, there exists a  $\hat{\Phi}_{\mathcal{A}}$  that is a  $\sqrt{n}$ -consistent local minimizer of  $Q(\Phi_{\mathcal{A}})$ , therefore  $\|\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*\|_2 = O_P\left(\frac{1}{\sqrt{n}}\right)$  and  $P\left(\hat{\Phi}_{\mathcal{A}^c} = \mathbf{0}\right) \rightarrow 1$ . Thus satisfies (with probability tending to 1):

$$\left. \frac{\partial Q_n(\Phi_{\mathcal{A}})}{\partial \Phi_m} \right|_{\Phi=\begin{pmatrix} \hat{\Phi}_{\mathcal{A}} \\ 0 \end{pmatrix}} = 0, \quad \forall m \in \mathcal{A}, \quad (22)$$

that is

$$\left. \frac{\partial Q_n(\Phi_{\mathcal{A}})}{\partial \Phi_m} \right|_{\Phi_{\mathcal{A}}=\hat{\Phi}_{\mathcal{A}}} = 0, \quad \forall m \in \mathcal{A}, \quad (23)$$

where

$$\begin{aligned} Q_n(\Phi_{\mathcal{A}}) &= -L_n(\Phi_{\mathcal{A}}) + n \underbrace{\sum_{m \in \mathcal{A}_1} \lambda_m \|\phi_m\|_2 + n \sum_{m \in \mathcal{A}_2} \lambda_m \|\phi_m\|_2}_{\triangleq nP(\Phi_{\mathcal{A}})} \\ &= -L_n(\Phi_{\mathcal{A}}) + nP(\Phi_{\mathcal{A}}). \end{aligned} \quad (24)$$

From (23) and (24) we have

$$\nabla_{\mathcal{A}} Q_n(\hat{\Phi}_{\mathcal{A}}) = -\nabla_{\mathcal{A}} L_n(\hat{\Phi}_{\mathcal{A}}) + n \nabla_{\mathcal{A}} P(\hat{\Phi}_{\mathcal{A}}) = \mathbf{0}, \quad (25)$$

with probability tending to 1.

Denote  $\Sigma = \text{diag}\{o_p(1), \dots, o_p(1)\}$ . We then expand  $-\nabla_{\mathcal{A}} L_n(\Phi_{\mathcal{A}})$  at  $\Phi_{\mathcal{A}} = \Phi_{\mathcal{A}}^*$  in (25):

$$\begin{aligned} -\nabla_{\mathcal{A}} L_n(\hat{\Phi}_{\mathcal{A}}) &= -\nabla_{\mathcal{A}} L_n(\Phi_{\mathcal{A}}^*) - [\nabla_{\mathcal{A}}^2 L_n(\Phi_{\mathcal{A}}^*) + \Sigma] (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \\ &= \sqrt{n} \left[ -\frac{1}{\sqrt{n}} \nabla_{\mathcal{A}} L_n(\Phi_{\mathcal{A}}^*) + \left( -\frac{1}{n} \nabla_{\mathcal{A}}^2 L_n(\Phi_{\mathcal{A}}^*) - \Sigma \right) \sqrt{n} (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \right] \\ &= \sqrt{n} \left[ -\frac{1}{\sqrt{n}} \nabla_{\mathcal{A}} L_n(\Phi_{\mathcal{A}}^*) + (\mathbf{I}(\Phi_{\mathcal{A}}^*) - \Sigma) \sqrt{n} (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \right]. \end{aligned}$$

The third line follows by

$$\frac{1}{n} \nabla_{\mathcal{A}}^2 L_n(\Phi_{\mathcal{A}}^*) = E\{\nabla_{\mathcal{A}}^2 L(\Phi_{\mathcal{A}}^*)\} + \Sigma = -\mathbf{I}(\Phi_{\mathcal{A}}^*) + \Sigma.$$

Denote

$$\mathbf{b} = (\lambda_m \text{sgn}(\beta_m^*), \lambda_m \frac{\theta_m^*}{\|\theta_m^*\|_2}, \lambda_m \text{sgn}(\gamma_{mE}^*))^\top, \quad m \in \mathcal{A},$$

We also expand  $n \nabla_{\mathcal{A}} P(\Phi_{\mathcal{A}})$  at  $\Phi_{\mathcal{A}} = \Phi_{\mathcal{A}}^*$  in (25):

$$n \nabla_{\mathcal{A}} P(\hat{\Phi}_{\mathcal{A}}) = n \left[ \mathbf{b} + \Sigma (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \right].$$

And due to the fact that  $\sqrt{n} \lambda_m \leq \sqrt{n} a_n \rightarrow 0$  for  $m \in \mathcal{A}$  and  $\frac{\theta_m^*}{\|\theta_m^*\|_2} \leq 1$  for any  $1 \leq k \leq p_m$ , we know that  $\sqrt{n} \mathbf{b} = (o_p(1), \dots, o_p(1))^\top$ . Thus,

$$\nabla_{\mathcal{A}} Q_n(\hat{\Phi}_{\mathcal{A}}) = \sqrt{n} \left[ -\frac{1}{\sqrt{n}} \nabla_{\mathcal{A}} L_n(\Phi_{\mathcal{A}}^*) + (\mathbf{I}(\Phi_{\mathcal{A}}^*) + \Sigma) \sqrt{n} (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \right]$$

$$\begin{aligned}
& + \sqrt{n} \left[ \sqrt{n} \mathbf{b} + \Sigma \sqrt{n} (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \right] \\
& = \sqrt{n} \left[ -\frac{1}{\sqrt{n}} \nabla_{\mathcal{A}} L_n(\Phi_{\mathcal{A}}^*) + \sqrt{n} \mathbf{b} + (\mathbf{I}(\Phi_{\mathcal{A}}^*) + \Sigma) \sqrt{n} (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \right] \\
& = \mathbf{0}.
\end{aligned}$$

$$(\mathbf{I}(\Phi_{\mathcal{A}}^*) + \Sigma) \sqrt{n} (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) = \sqrt{n} \frac{1}{n} \sum_{i=1}^n \nabla_{\mathcal{A}} \log f(V_i, \Phi_{\mathcal{A}}^*) + o_p(1).$$

Therefore, by the central limit theorem, we know that

$$\sqrt{n} \left[ \frac{1}{n} \sum_{i=1}^n \nabla_{\mathcal{A}} \log f(V_i, \Phi_{\mathcal{A}}^*) \right] \rightarrow N(\mathbf{0}, \mathbf{I}(\Phi_{\mathcal{A}}^*)).$$

Hence,

$$\sqrt{n} (\hat{\Phi}_{\mathcal{A}} - \Phi_{\mathcal{A}}^*) \xrightarrow{d} N(\mathbf{0}, \mathbf{I}^{-1}(\Phi_{\mathcal{A}}^*)).$$

□

## B. Algorithm Details

In this section we provide more specific details about the algorithms used to solve the `sail` objective function. We assume that  $Y$ ,  $\Psi_j$ ,  $X_E$  and  $X_E \circ \Psi_j$  have been centered by their sample means  $\bar{Y}$ ,  $\bar{\Psi}_j$ ,  $\bar{X}_E$ , and  $\bar{X}_E \circ \Psi_j$ , respectively. Here,  $\bar{\Psi}_j \in \mathbb{R}^{m_j}$  and  $\bar{X}_E \circ \Psi_j \in \mathbb{R}^{m_j}$  represent the column means of  $\Psi_j$  and  $X_E \circ \Psi_j$ , respectively. Since the intercept ( $\beta_0$ ) is not penalized and all variables have been centered, we can omit it from the loss function and compute it once the algorithm has converged for all other parameters. The strong heredity `sail` model with least-squares loss has the form

$$\hat{Y} = \sum_{j=1}^p \Psi_j \theta_j + \beta_E X_E + \sum_{j=1}^p \gamma_j \beta_E (X_E \circ \Psi_j) \theta_j \quad (26)$$

and the objective function is given by

$$Q(\Phi) = \frac{1}{2n} \|Y - \hat{Y}\|_2^2 + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda \alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (27)$$

Solving (27) in a blockwise manner allows us to leverage computationally fast algorithms for  $\ell_1$  and  $\ell_2$  norm penalized regression. Denote the  $n$ -dimensional residual column vector  $R = Y - \hat{Y}$ . The subgradient equations are given by

$$\frac{\partial Q}{\partial \beta_E} = -\frac{1}{n} \left( X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^T R + \lambda(1 - \alpha) w_E s_1 = 0 \quad (28)$$

$$\frac{\partial Q}{\partial \theta_j} = -\frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^T R + \lambda(1 - \alpha) w_j s_2 = \mathbf{0} \quad (29)$$

$$\frac{\partial Q}{\partial \gamma_j} = -\frac{1}{n} (\beta_E (X_E \circ \Psi_j) \theta_j)^T R + \lambda \alpha w_{jE} s_3 = 0 \quad (30)$$

where  $s_1$  is in the subgradient of the  $\ell_1$  norm:

$$s_1 \in \begin{cases} \text{sign}(\beta_E) & \text{if } \beta_E \neq 0 \\ [-1, 1] & \text{if } \beta_E = 0, \end{cases}$$

$s_2$  is in the subgradient of the  $\ell_2$  norm:

$$s_2 \in \begin{cases} \frac{\theta_j}{\|\theta_j\|_2} & \text{if } \theta_j \neq \mathbf{0} \\ u \in \mathbb{R}^{m_j} : \|u\|_2 \leq 1 & \text{if } \theta_j = \mathbf{0}, \end{cases}$$

and  $s_3$  is in the subgradient of the  $\ell_1$  norm:

$$s_3 \in \begin{cases} \text{sign}(\gamma_j) & \text{if } \gamma_j \neq 0 \\ [-1, 1] & \text{if } \gamma_j = 0. \end{cases}$$

Define the partial residuals, without the  $j$ th predictor for  $j = 1, \dots, p$ , as

$$R_{(-j)} = Y - \sum_{\ell \neq j} \Psi_\ell \theta_\ell - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell \beta_E (X_E \circ \Psi_\ell) \theta_\ell$$

the partial residual without  $X_E$  as

$$R_{(-E)} = Y - \sum_{j=1}^p \Psi_j \theta_j$$

and the partial residual without the  $j$ th interaction for  $j = 1, \dots, p$ , as

$$R_{(-jE)} = Y - \sum_{j=1}^p \Psi_j \theta_j - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell \beta_E (X_E \circ \Psi_\ell) \theta_\ell$$

From the subgradient equations (28)–(30) we see that

$$\hat{\beta}_E = \frac{S\left(\frac{1}{n \cdot w_E} \left(X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \hat{\theta}_j\right)^\top R_{(-E)}, \lambda(1-\alpha)\right)}{\left(X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \hat{\theta}_j\right)^\top \left(X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \hat{\theta}_j\right)} \quad (31)$$

$$\lambda(1-\alpha) w_j \frac{\theta_j}{\|\theta_j\|_2} = \frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-j)} \quad (32)$$

$$\hat{\gamma}_j = \frac{S\left(\frac{1}{n \cdot w_{jE}} (\beta_E (X_E \circ \Psi_j) \theta_j)^\top R_{(-jE)}, \lambda\alpha\right)}{(\beta_E (X_E \circ \Psi_j) \theta_j)^\top (\beta_E (X_E \circ \Psi_j) \theta_j)} \quad (33)$$

where  $S(x, t) = \text{sign}(x)(|x| - t)$  is the soft-thresholding operator. Given these estimates, the intercept can be computed using the following equation:

$$\hat{\beta}_0 = \bar{Y} - \sum_{j=1}^p \bar{\Psi}_j \hat{\theta}_j - \hat{\beta}_E \bar{X}_E - \sum_{j=1}^p \hat{\gamma}_j \hat{\beta}_E (\bar{X}_E \circ \bar{\Psi}_j) \hat{\theta}_j. \quad (34)$$

We see from (31) that there is a closed form solution for  $\beta_E$ . From (33), each  $\gamma_j$  also has a closed form solution and can be solved efficiently for  $j = 1, \dots, p$  using a coordinate descent procedure (Friedman et al., 2010). Since there is no closed form solution for  $\beta_j$ , we use a quadratic majorization technique (Yang and Zou, 2015) to solve (32). Furthermore, we update each  $\theta_j$  in a coordinate wise fashion and leverage this to implement further computational speedups which are detailed in Supplemental Section B.2. From these estimates, we compute the interaction effects using the reparametrizations presented in Table 1, e.g.,  $\hat{\tau}_j = \hat{\gamma}_j \hat{\beta}_E \hat{\theta}_j$ ,  $j = 1, \dots, p$  for the strong heredity sail model.

## B.1. Least-Squares sail with Strong Heredity

A more detailed algorithm for fitting the least-squares sail model with strong heredity is given in Algorithm 3.

**Algorithm 3** Blockwise Coordinate Descent for Least-Squares sail with Strong Heredity

---

1: **function** sail( $X, Y, X_E, \text{basis}, \lambda, \alpha, w_j, w_E, w_{jE}, \epsilon$ ) ▷ Algorithm for solving (27)  
2:    $\Psi_j \leftarrow \text{basis}(X_j), \tilde{\Psi}_j \leftarrow X_E \circ \Psi_j$  for  $j = 1, \dots, p$   
3:   Center all variables by their sample means  
4:   Initialize:  $\beta_E^{(0)} = \theta_j^{(0)} = \gamma_j^{(0)} \leftarrow 0$  for  $j = 1, \dots, p$ .  
5:   Set iteration counter  $k \leftarrow 0$   
6:    $R^* \leftarrow Y - \beta_E^{(k)} X_E - \sum_j (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)}$   
7:   **repeat**  
8:     • To update  $\gamma = (\gamma_1, \dots, \gamma_p)$   
9:        $\tilde{X}_j \leftarrow \beta_E^{(k)} \tilde{\Psi}_j \theta_j^{(k)}$  for  $j = 1, \dots, p$   
10:       $R \leftarrow R^* + \sum_{j=1}^p \gamma_j^{(k)} \tilde{X}_j$   
11:  
12:       $\gamma^{(k)(new)} \leftarrow \arg \min_{\gamma} \frac{1}{2n} \left\| R - \sum_j \gamma_j \tilde{X}_j \right\|_2^2 + \lambda \alpha \sum_j w_{jE} |\gamma_j|$   
13:       $\Delta = \sum_j (\gamma_j^{(k)} - \gamma_j^{(k)(new)}) \tilde{X}_j$   
14:       $R^* \leftarrow R^* + \Delta$   
15:      • To update  $\theta = (\theta_1, \dots, \theta_p)$   
16:        $\tilde{X}_j \leftarrow \Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j$  for  $j = 1, \dots, p$   
17:       **for**  $j = 1, \dots, p$  **do**  
18:          $R \leftarrow R^* + \tilde{X}_j \theta_j^{(k)}$   
19:  
20:          $\Delta = \tilde{X}_j (\theta_j^{(k)} - \theta_j^{(k)(new)})$   
21:          $R^* \leftarrow R^* + \Delta$   
22:         • To update  $\beta_E$   
23:          $\tilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \tilde{\Psi}_j \theta_j^{(k)}$   
24:          $R \leftarrow R^* + \beta_E^{(k)} \tilde{X}_E$   
25:          $\beta_E^{(k)(new)} \leftarrow \frac{1}{\tilde{X}_E^\top \tilde{X}_E} S \left( \frac{1}{n \cdot w_E} \tilde{X}_E^\top R, \lambda(1 - \alpha) \right)$   
26:          $\Delta = (\beta_E^{(k)} - \beta_E^{(k)(new)}) \tilde{X}_E$   
27:          $R^* \leftarrow R^* + \Delta$   
28:  
29:         **until** convergence criterion is satisfied:  $|Q(\Phi^{(k-1)}) - Q(\Phi^{(k)})| / Q(\Phi^{(k-1)}) < \epsilon$   
30:         Compute the intercept  $\beta_0$   
31:          $\beta_0 \leftarrow \bar{Y} - \sum_{j=1}^p \bar{\Psi}_j \hat{\theta}_j - \hat{\beta}_E \bar{X}_E - \sum_{j=1}^p \hat{\gamma}_j \hat{\beta}_E (\bar{X}_E \circ \bar{\Psi}_j) \hat{\theta}_j$   


---

**B.2. Details on Update for  $\theta$** 

Here we discuss a computational speedup in the updates for the  $\theta$  parameter. The partial residual ( $R_s$ ) used for updating  $\theta_s$  ( $s \in 1, \dots, p$ ) at the  $k$ th iteration is given by

$$R_s = Y - \tilde{Y}_{(-s)}^{(k)} \quad (35)$$

where  $\tilde{Y}_{(-s)}^{(k)}$  is the fitted value at the  $k$ th iteration excluding the contribution from  $\Psi_s$ :

$$\tilde{Y}_{(-s)}^{(k)} = \beta_E^{(k)} X_E + \sum_{\ell \neq s} \Psi_\ell \theta_\ell^{(k)} + \sum_{\ell \neq s} \gamma_\ell^{(k)} \beta_E^{(k)} \tilde{\Psi}_\ell \theta_\ell^{(k)} \quad (36)$$

Using (36), (35) can be re-written as

$$\begin{aligned} R_s &= Y - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \\ &= R^* + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \end{aligned} \quad (37)$$

where

$$R^* = Y - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)} \quad (38)$$

Denote  $\theta_s^{(k)(new)}$  the solution for predictor  $s$  at the  $k$ th iteration, given by:

$$\theta_s^{(k)(new)} = \arg \min_{\theta_j} \frac{1}{2n} \|R_s - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_j\|_2^2 + \lambda(1-\alpha) w_s \|\theta_j\|_2 \quad (39)$$

Now we want to update the parameters for the next predictor  $\theta_{s+1}$  ( $s+1 \in 1, \dots, p$ ) at the  $k$ th iteration. The partial residual used to update  $\theta_{s+1}$  is given by

$$R_{s+1} = R^* + (\Psi_{s+1} + \gamma_{s+1}^{(k)} \beta_E^{(k)} \tilde{\Psi}_{s+1}) \theta_{s+1}^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) (\theta_s^{(k)} - \theta_s^{(k)(new)}) \quad (40)$$

where  $R^*$  is given by (38),  $\theta_s^{(k)}$  is the parameter value prior to the update, and  $\theta_s^{(k)(new)}$  is the updated value given by (39). Taking the difference between (37) and (40) gives

$$\begin{aligned} \Delta &= R_t - R_s \\ &= (\Psi_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\Psi}_t) \theta_t^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) (\theta_s^{(k)} - \theta_s^{(k)(new)}) - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \\ &= (\Psi_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\Psi}_t) \theta_t^{(k)} - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)(new)} \end{aligned} \quad (41)$$

Therefore  $R_t = R_s + \Delta$ , and the partial residual for updating the next predictor can be computed by updating the previous partial residual by  $\Delta$ , given by (41). This formulation can lead to computational speedups especially when  $\Delta = 0$ , meaning the partial residual does not need to be re-calculated.

### B.3. Maximum penalty parameter ( $\lambda_{max}$ ) for strong heredity

The subgradient equations (28)–(30) can be used to determine the largest value of  $\lambda$  such that all coefficients are 0. From the subgradient Equation (28), we see that  $\beta_E = 0$  is a solution if

$$\frac{1}{w_E} \left| \frac{1}{n} \left( X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^\top R_{(-E)} \right| \leq \lambda(1-\alpha) \quad (42)$$

From the subgradient Equation (29), we see that  $\theta_j = \mathbf{0}$  is a solution if

$$\frac{1}{w_j} \left\| \frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-jE)} \right\|_2 \leq \lambda(1-\alpha) \quad (43)$$

From the subgradient Equation (30), we see that  $\gamma_j = 0$  is a solution if

$$\frac{1}{w_{jE}} \left| \frac{1}{n} (\beta_E (X_E \circ \Psi_j) \theta_j)^\top R_{(-jE)} \right| \leq \lambda \alpha \quad (44)$$

Due to the strong heredity property, the parameter vector  $(\beta_E, \theta_1, \dots, \theta_p, \gamma_1, \dots, \gamma_p)$  will be entirely equal to  $\mathbf{0}$  if  $(\beta_E, \theta_1, \dots, \theta_p) = \mathbf{0}$ . Therefore, the smallest value of  $\lambda$  for which the entire parameter vector (excluding the intercept) is  $\mathbf{0}$  is:

$$\lambda_{max} = \frac{1}{n(1-\alpha)} \max \left\{ \frac{1}{w_E} \left( X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^T R_{(-E)}, \max_j \frac{1}{w_j} \|(\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^T R_{(-j)}\|_2 \right\} \quad (45)$$

which reduces to

$$\lambda_{max} = \frac{1}{n(1-\alpha)} \max \left\{ \frac{1}{w_E} (X_E)^T R_{(-E)}, \max_j \frac{1}{w_j} \|(\Psi_j)^T R_{(-j)}\|_2 \right\}$$

#### B.4. Least-Squares sail with Weak Heredity

We assume the same centering constraints as in Section B.1. The least-squares sail model with weak heredity has the form

$$\hat{Y} = \sum_{j=1}^p \Psi_j \theta_j + \beta_E X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \quad (46)$$

The objective function is given by

$$Q(\Phi) = \frac{1}{2n} \|Y - \hat{Y}\|_2^2 + \lambda(1-\alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (47)$$

Denote the  $n$ -dimensional residual column vector  $R = Y - \hat{Y}$ . The subgradient equations are given by

$$\frac{\partial Q}{\partial \beta_E} = -\frac{1}{n} \left( X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^T R + \lambda(1-\alpha) w_E s_1 = 0 \quad (48)$$

$$\frac{\partial Q}{\partial \theta_j} = -\frac{1}{n} (\Psi_j + \gamma_j (X_E \circ \Psi_j))^T R + \lambda(1-\alpha) w_j s_2 = \mathbf{0} \quad (49)$$

$$\frac{\partial Q}{\partial \gamma_j} = -\frac{1}{n} \left( (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \right)^T R + \lambda\alpha w_{jE} s_3 = 0 \quad (50)$$

where  $s_1$  is in the subgradient of the  $\ell_1$  norm:

$$s_1 \in \begin{cases} \text{sign}(\beta_E) & \text{if } \beta_E \neq 0 \\ [-1, 1] & \text{if } \beta_E = 0, \end{cases}$$

$s_2$  is in the subgradient of the  $\ell_2$  norm:

$$s_2 \in \begin{cases} \frac{\theta_j}{\|\theta_j\|_2} & \text{if } \theta_j \neq \mathbf{0} \\ u \in \mathbb{R}^{m_j} : \|u\|_2 \leq 1 & \text{if } \theta_j = \mathbf{0}, \end{cases}$$

and  $s_3$  is in the subgradient of the  $\ell_1$  norm:

$$s_3 \in \begin{cases} \text{sign}(\gamma_j) & \text{if } \gamma_j \neq 0 \\ [-1, 1] & \text{if } \gamma_j = 0. \end{cases}$$

Define the partial residuals, without the  $j$ th predictor for  $j = 1, \dots, p$ , as

$$R_{(-j)} = Y - \sum_{\ell \neq j} \Psi_\ell \theta_\ell - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell (X_E \circ \Psi_\ell) (\beta_E \cdot \mathbf{1}_{m_\ell} + \theta_\ell)$$

the partial residual without  $X_E$  as

$$R_{(-E)} = Y - \sum_{j=1}^p \Psi_j \theta_j - \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j$$

and the partial residual without the  $j$ th interaction for  $j = 1, \dots, p$

$$R_{(-jE)} = Y - \sum_{j=1}^p \Psi_j \theta_j - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell (X_E \circ \Psi_\ell) (\beta_E \cdot \mathbf{1}_{m_\ell} + \theta_\ell)$$

From the subgradient Equation (48), we see that  $\beta_E = 0$  is a solution if

$$\frac{1}{w_E} \left| \frac{1}{n} \left( X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^\top R_{(-E)} \right| \leq \lambda(1-\alpha) \quad (51)$$

From the subgradient Equation (49), we see that  $\theta_j = \mathbf{0}$  is a solution if

$$\frac{1}{w_j} \left\| \frac{1}{n} \left( \Psi_j + \gamma_j (X_E \circ \Psi_j) \right)^\top R_{(-j)} \right\|_2 \leq \lambda(1-\alpha) \quad (52)$$

From the subgradient Equation (50), we see that  $\gamma_j = 0$  is a solution if

$$\frac{1}{w_{jE}} \left| \frac{1}{n} \left( (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \right)^\top R_{(-jE)} \right| \leq \lambda\alpha \quad (53)$$

From the subgradient equations we see that

$$\hat{\beta}_E = \frac{S \left( \frac{1}{n \cdot w_E} \left( X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^\top R_{(-E)}, \lambda(1-\alpha) \right)}{\left( X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^\top \left( X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)} \quad (54)$$

$$\lambda(1-\alpha) w_j \frac{\theta_j}{\|\theta_j\|_2} = \frac{1}{n} \left( \Psi_j + \gamma_j (X_E \circ \Psi_j) \right)^\top R_{(-j)} \quad (55)$$

$$\hat{\gamma}_j = \frac{S \left( \frac{1}{n \cdot w_{jE}} \left( (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \right)^\top R_{(-jE)}, \lambda\alpha \right)}{\left( (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \right)^\top \left( (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \right)} \quad (56)$$

where  $S(x, t) = \text{sign}(x)(|x| - t)$  is the soft-thresholding operator. As was the case in the strong heredity sail model, there is a closed form solution for  $\beta_E$ , each  $\gamma_j$  also has a closed form solution and can be solved efficiently for  $j = 1, \dots, p$  using the coordinate descent procedure implemented in the `glmnet` package (Friedman et al., 2010), while we use the quadratic majorization technique implemented in the `gglasso` package (Yang and Zou, 2015) to solve (55). Algorithm 4 details the procedure used to fit the least-squares weak heredity sail model.

**Algorithm 4** Coordinate descent for least-squares sail with weak heredity

---

1: **function** sail( $X, Y, X_E, \text{basis}, \lambda, \alpha, w_j, w_E, w_{jE}, \epsilon$ ) ▷ Algorithm for solving (47)  
2:    $\Psi_j \leftarrow \text{basis}(X_j)$ ,  $\widetilde{\Psi}_j \leftarrow X_E \circ \Psi_j$  for  $j = 1, \dots, p$   
3:   Center all variables by their sample means  
4:   Initialize:  $\beta_E^{(0)} = \theta_j^{(0)} = \gamma_j^{(0)} \leftarrow 0$  for  $j = 1, \dots, p$ .  
5:   Set iteration counter  $k \leftarrow 0$   
6:    $R^* \leftarrow Y - \beta_E^{(k)} X_E - \sum_j \Psi_j \theta_j^{(k)} - \sum_j \gamma_j^{(k)} \widetilde{\Psi}_j (\beta_E^{(k)} \cdot \mathbf{1}_{m_j} + \theta_j^{(k)})$   
7:   **repeat**  
8:     • To update  $\gamma = (\gamma_1, \dots, \gamma_p)$   
9:        $\widetilde{X}_j \leftarrow \widetilde{\Psi}_j (\beta_E^{(k)} \cdot \mathbf{1}_{m_j} + \theta_j^{(k)})$  for  $j = 1, \dots, p$   
10:       $R \leftarrow R^* + \sum_{j=1}^p \gamma_j^{(k)} \widetilde{X}_j$   
11:  
12:       $\gamma^{(k)(new)} \leftarrow \arg \min_{\gamma} \frac{1}{2n} \left\| R - \sum_j \gamma_j \widetilde{X}_j \right\|_2^2 + \lambda \alpha \sum_j w_{jE} |\gamma_j|$   
13:       $\Delta = \sum_j (\gamma_j^{(k)} - \gamma_j^{(k)(new)}) \widetilde{X}_j$   
14:       $R^* \leftarrow R^* + \Delta$   
15:     • To update  $\theta = (\theta_1, \dots, \theta_p)$   
16:        $\widetilde{X}_j \leftarrow \Psi_j + \gamma_j^{(k)} \widetilde{\Psi}_j$  for  $j = 1, \dots, p$   
17:       **for**  $j = 1, \dots, p$  **do**  
18:          $R \leftarrow R^* + \widetilde{X}_j \theta_j^{(k)}$   
19:  
20:          $\Delta = \widetilde{X}_j (\theta_j^{(k)} - \theta_j^{(k)(new)})$   
21:          $R^* \leftarrow R^* + \Delta$   
22:     • To update  $\beta_E$   
23:        $\widetilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \widetilde{\Psi}_j \mathbf{1}_{m_j}$   
24:        $R \leftarrow R^* + \beta_E^{(k)} \widetilde{X}_E$   
25:  
26:          $\Delta = (\beta_E^{(k)} - \beta_E^{(k)(new)}) \widetilde{X}_E$   
27:          $R^* \leftarrow R^* + \Delta$   
28:  
29:          $k \leftarrow k + 1$   
30:  
31:         **until** convergence criterion is satisfied:  $|Q(\Phi^{(k-1)}) - Q(\Phi^{(k)})| / Q(\Phi^{(k-1)}) < \epsilon$   
32:         Compute the intercept  $\beta_0$   
33:          $\beta_0 \leftarrow \bar{Y} - \sum_{j=1}^p \bar{\Psi}_j \hat{\theta}_j - \hat{\beta}_E \bar{X}_E - \sum_{j=1}^p \hat{\gamma}_j \hat{\beta}_E (\bar{X}_E \circ \bar{\Psi}_j) \hat{\theta}_j$   


---

▷  $S(x, t) = \text{sign}(x)(|x| - t)_+$

**B.4.1. Maximum penalty parameter ( $\lambda_{max}$ ) for weak heredity**

The smallest value of  $\lambda$  for which the entire parameter vector  $(\beta_E, \theta_1, \dots, \theta_p, \gamma_1, \dots, \gamma_p)$  is  $\mathbf{0}$  is:

$$\lambda_{max} = \frac{1}{n} \max \left\{ \frac{1}{(1-\alpha)w_E} \left( X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^T R_{(-E)}, \right.$$

$$\max_j \frac{1}{(1-\alpha)w_j} \left\| (\Psi_j + \gamma_j(X_E \circ \Psi_j))^\top R_{(-j)} \right\|_2, \\ \max_j \frac{1}{\alpha w_{jE}} \left( (X_E \circ \Psi_j)(\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \right)^\top R_{(-jE)} \quad (57)$$

which reduces to

$$\lambda_{max} = \frac{1}{n(1-\alpha)} \max \left\{ \frac{1}{w_E} (X_E)^\top R_{(-E)}, \max_j \frac{1}{w_j} \left\| (\Psi_j)^\top R_{(-j)} \right\|_2 \right\}$$

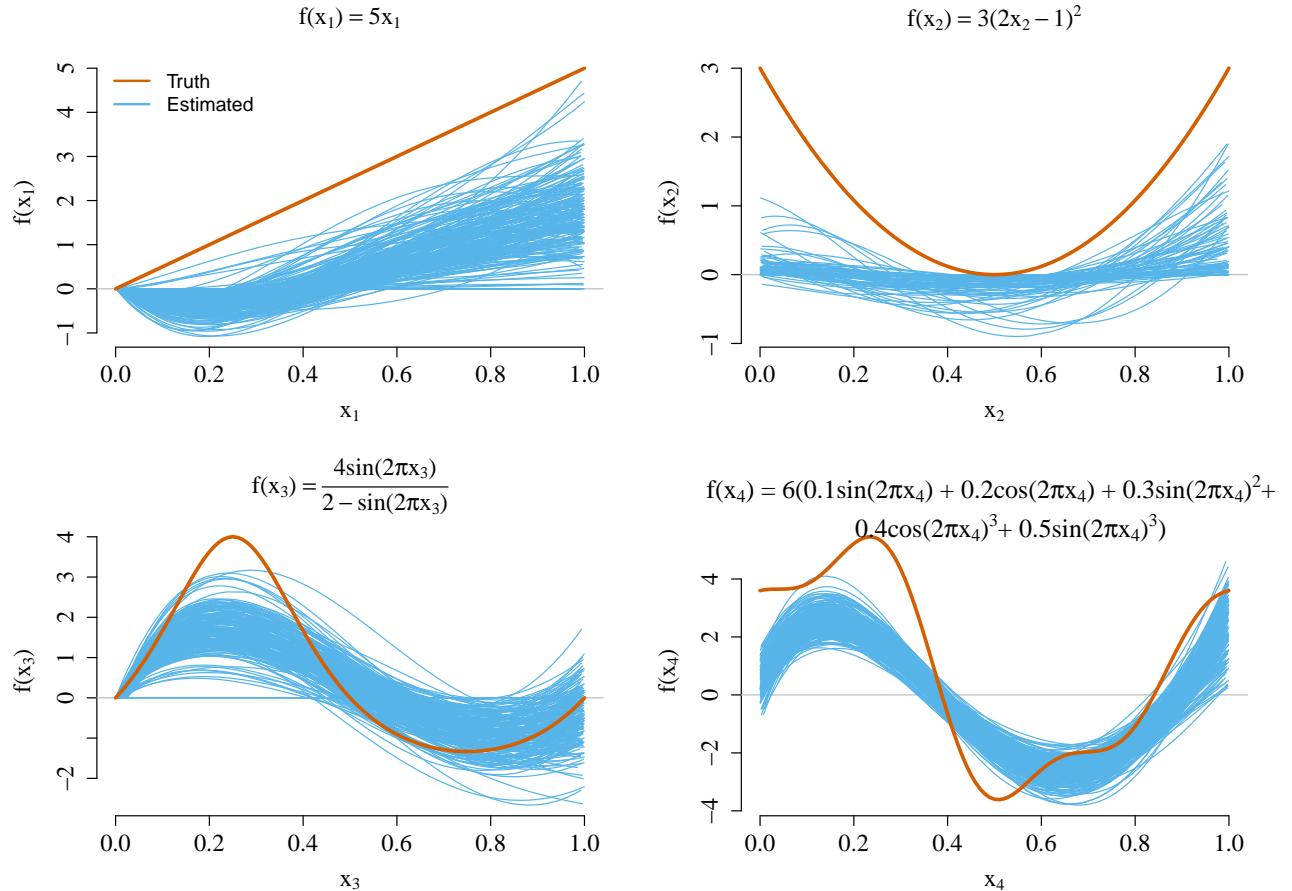
This is the same  $\lambda_{max}$  as the least-squares strong heredity `sail` model.

## C. Additional Simulation Results

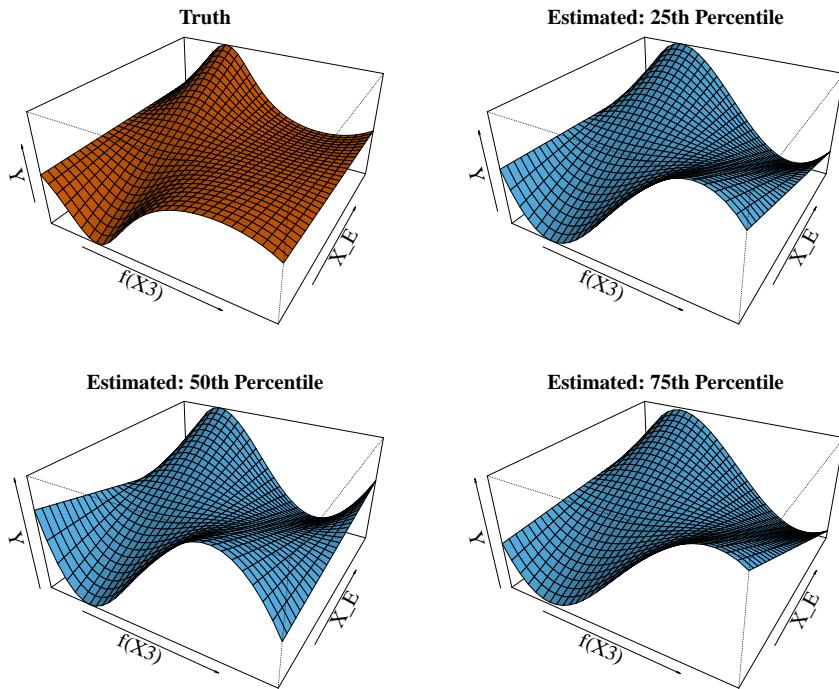
We visually inspected whether our method could correctly capture the shape of the association between the predictors and the response for both main and interaction effects. To do so, we plotted the true and predicted curves for scenario 1a) only. Figure 5 shows each of the four main effects with the estimated curves from each of the 200 simulations along with the true curve. We can see the effect of the penalty on the parameters, i.e., decreasing prediction variance at the cost of increased bias. This is particularly well illustrated in the bottom right panel where `sail` smooths out the very wiggly component function  $f_4(x)$ . Nevertheless, the primary shapes are clearly being captured.

To visualize the estimated interaction effects, we ordered the 200 simulation runs by the Euclidean distance between the estimated and true regression functions. Following Radchenko et al. (Radchenko and James, 2010), we then identified the 25th, 50th, and 75th best simulations and plotted, in Figures 6 and 7, the interaction effects of  $X_E$  with  $f_3(X_3)$  and  $f_4(X_4)$ , respectively. We see that `sail` does a good job at capturing the true interaction surface for  $X_E \cdot f_3(X_3)$ . Again, the smoothing and shrinkage effect is apparent when looking at the interaction surfaces for  $X_E \cdot f_4(X_4)$ .

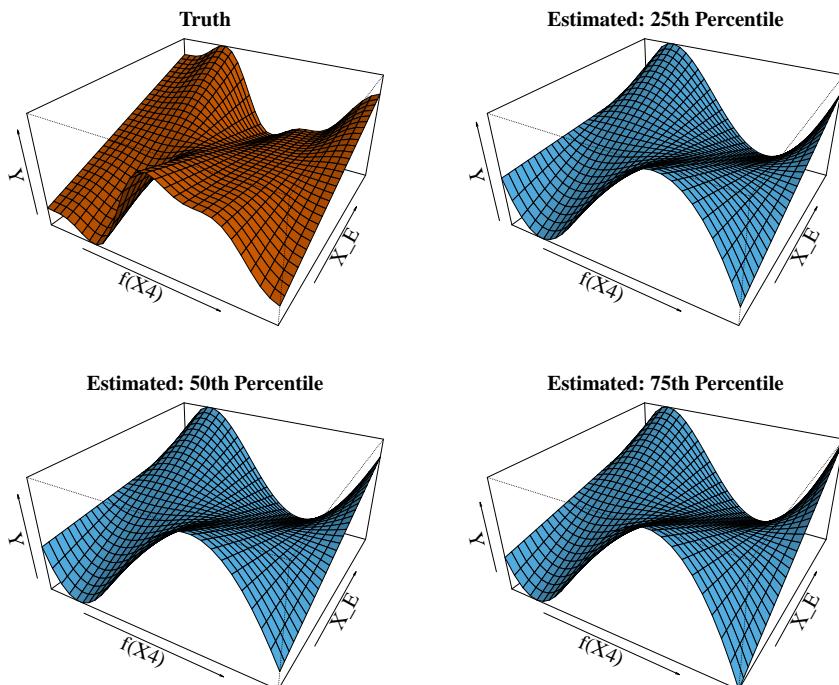
In Figure 8 we visualize the variable selection results from 210 replications of the simulation study for strong hierarchy `sail` using UpSet plots (Conway, Lex and Gehlenborg, 2017). Shown are the selected models and their frequencies. We can see that the environment variable is always selected across all simulation scenarios and replications.



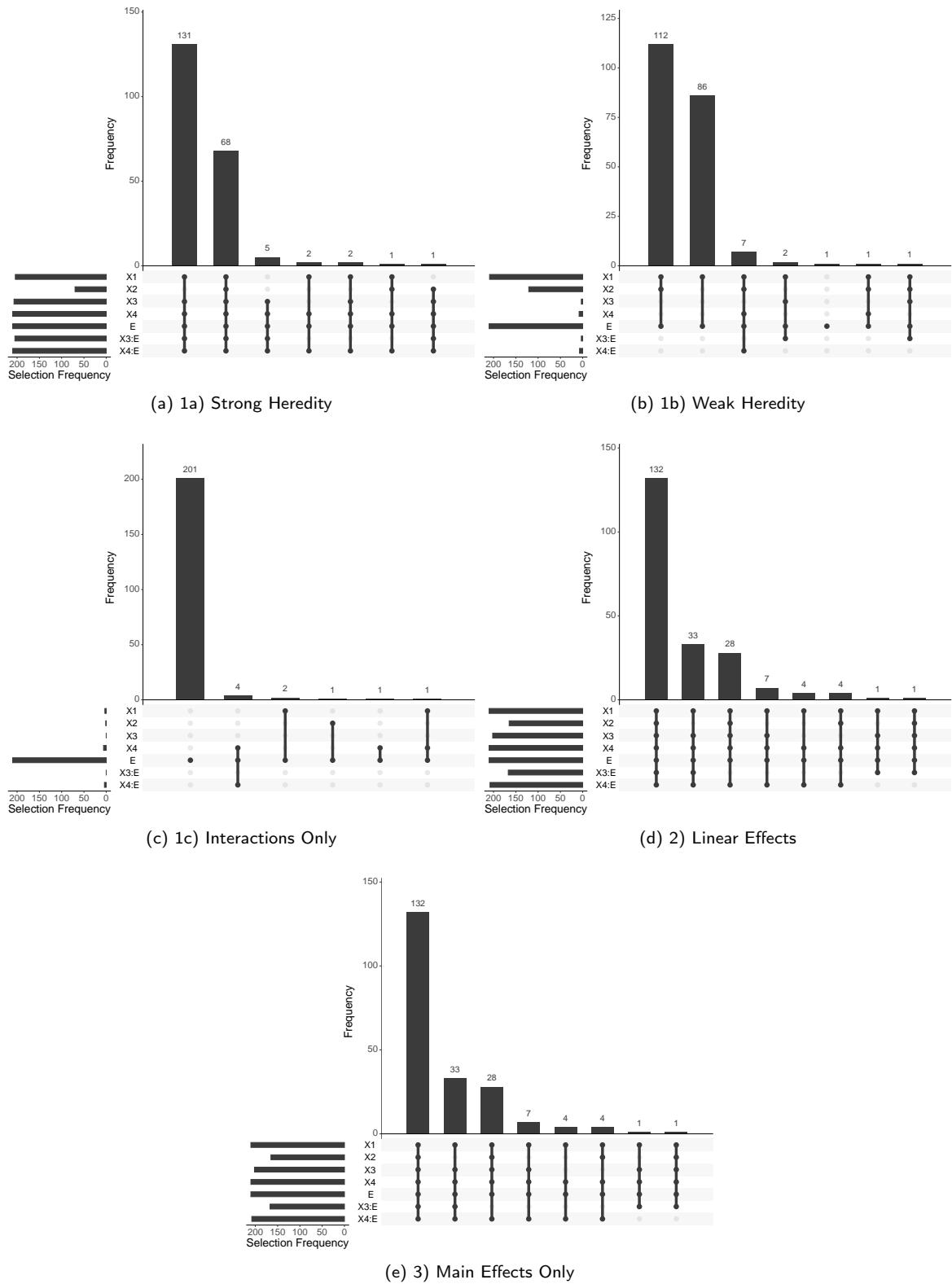
**Figure 5:** True and estimated main effect component functions for scenario 1a). The estimated curves represent the results from each one of the 200 replications conducted.



**Figure 6:** True and estimated interaction effects for  $X_E \cdot f_3(X_3)$  in simulation scenario 1a).

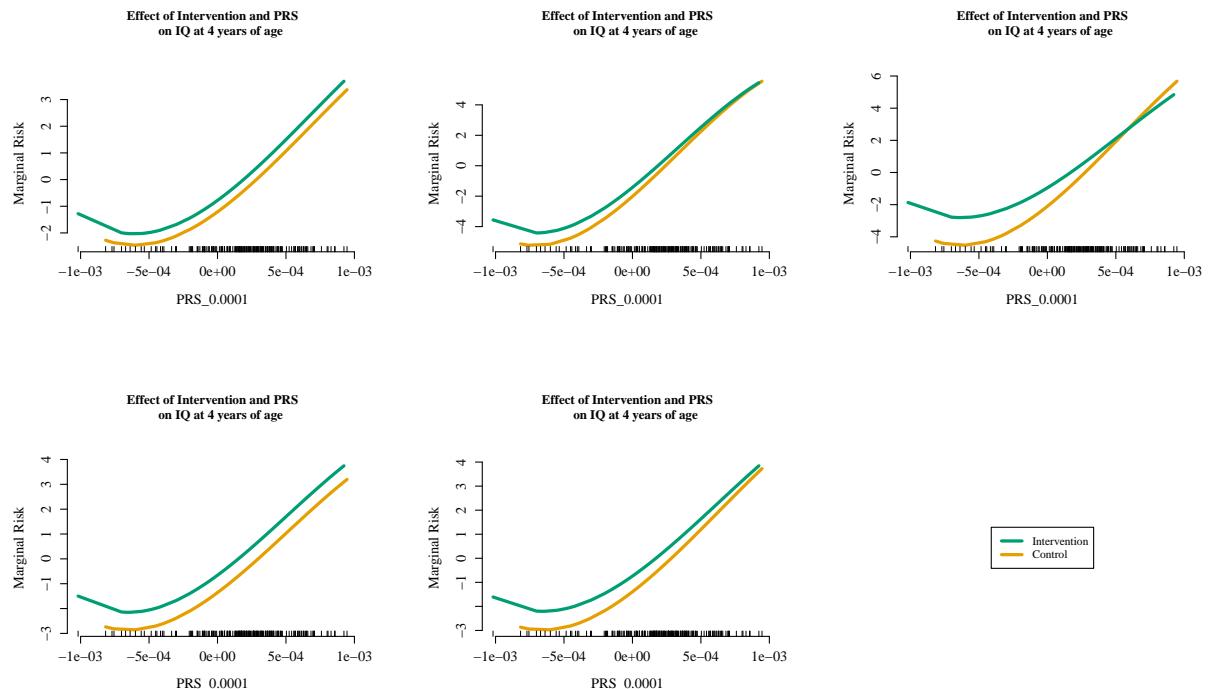


**Figure 7:** True and estimated interaction effects for  $X_E \cdot f_4(X_4)$  in simulation scenario 1a).

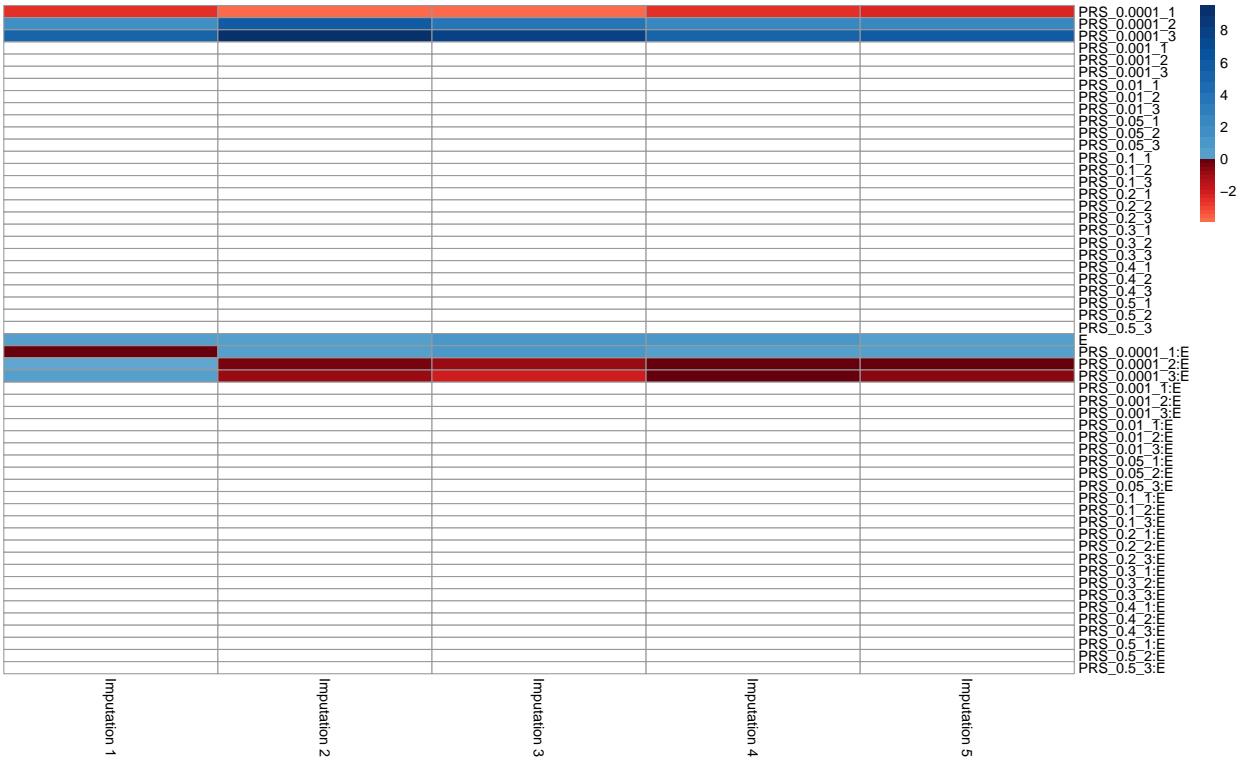


**Figure 8:** Variable selection results from 210 replications of the simulation study for strong hierarchy sail visualized using UpSet plots (Conway et al., 2017). Shown are the selected models and their frequencies. We can see that the environment variable is always selected across all simulation scenarios and replications.

## D. Additional Results on PRS for Educational Attainment



**Figure 9:** Estimated interaction effect identified by the weak heredity sail using cubic B-splines and  $\alpha = 0.1$  for the Nurse Family Partnership data for the 5 imputed datasets. Of the 189 subjects, 19 IQ scores were imputed using mice (Buuren and Groothuis-Oudshoorn, 2010). The selected model, chosen via 10-fold cross-validation, contained three variables: the main effects for the intervention and the PRS for educational attainment using genetic variants significant at the 0.0001 level, as well as their interaction.



**Figure 10:** Coefficient estimates obtained by the weak heredity sail using cubic B-splines and  $\alpha = 0.1$  for the Nurse Family Partnership data for the 5 imputed datasets. Of the 189 subjects, 19 IQ scores were imputed using mice (Buuren and Groothuis-Oudshoorn, 2010). The selected model, chosen via 10-fold cross-validation, contained three variables: the main effects for the intervention and the PRS for educational attainment using genetic variants significant at the 0.0001 level, as well as their interaction. This result was consistent across all 5 imputed datasets. The white boxes indicate a coefficient estimate of 0.

## E. Data Availability and Code to Reproduce Results

The R scripts used to simulate the data for the simulation studies in Section 4 are provided along with the code for each of the methods being compared. The data used for the two real data analyses in Section 5 are publicly available. The first dataset from the Nurse Family Partnership program is provided by one of the authors of the manuscript (David Olds). The second dataset from the Study to Understand Prognoses Preferences Outcomes and Risks of Treatment (SUPPORT) is publicly available from the Vanderbilt University Department of Biostatistics website.

### E.1. Datasets

The datasets are available at [https://github.com/sahirbhatnagar/sail/tree/master/manuscript/raw\\_data](https://github.com/sahirbhatnagar/sail/tree/master/manuscript/raw_data)

1. Nurse Family Partnership program data consists of three files. They are merged together using the script [https://github.com/sahirbhatnagar/sail/blob/master/manuscript/bin/PRS\\_bootstrap.R](https://github.com/sahirbhatnagar/sail/blob/master/manuscript/bin/PRS_bootstrap.R)
  - Gen\_3PC\_scores.txt
  - IQ\_and\_mental\_development\_variables\_for\_Sahir\_with\_study\_ID.txt
  - NFP\_170614\_INFO08\_nodup\_hard09\_noambi\_GWAS\_EduYears\_Pooled\_beta\_withaf\_5000pruned\_noambi\_16Jan2018.score
2. The SUPPORT data consists of a single file:
  - [https://github.com/sahirbhatnagar/sail/blob/master/manuscript/raw\\_data/support2.csv](https://github.com/sahirbhatnagar/sail/blob/master/manuscript/raw_data/support2.csv)

All datasets are in .txt format. Code used to read in the datasets are provided in the section below. All output from this project published online is available according to the conditions of the Creative Commons License (<https://creativecommons.org/licenses/by-nc-sa/2.0/>)

## E.2. Code

The software which implements our algorithm is available in an R package published on CRAN (<https://cran.r-project.org/package=sail>) version 0.1.0 with MIT license. The paper itself is written in knitr format, and therefore includes both the code and text in the same .Rnw file.

The scripts and data used to produce the results in the manuscript are available at <https://github.com/sahirbhatnagar/sail/tree/master/manuscript>.

The knitr file which contains both the main text and code is available at: [https://github.com/sahirbhatnagar/sail/blob/master/manuscript/source/sail\\_manuscript\\_v2.Rnw](https://github.com/sahirbhatnagar/sail/blob/master/manuscript/source/sail_manuscript_v2.Rnw)

The manuscript was compiled using R version 3.6.1 with knitr version 1.25.

The bootstrap analysis was run in parallel on a compute cluster with 40 cores. Though this is not necessary to reproduce the results, it definitely speeds up the computation time.

### E.2.1. Instructions for Use

All tables and figures from the paper can be reproduced by compiling the knitr file. The easiest way to reproduce the results is to download the GitHub repository and compile the knitr file from within an R session as follows:

1. Download the GitHub repository <https://github.com/sahirbhatnagar/sail/archive/master.zip>
2. From within an R session, run the command: `knitr::knit2pdf('sail_manuscript_v2.Rnw')`

Note that to speed up compilation time, we have saved the simulation and bootstrap results in .RData files available at <https://github.com/sahirbhatnagar/sail/tree/master/manuscript/results>. These .RData files are called directly by the knitr file.

Note also that the R scripts used to generate the results are called from the knitr file using the ‘code externalization’ functionality of knitr (<https://yihui.org/knitr/demo/externalization/>). That is, the actual R code is stored in R scripts and not within the knitr file. These R scripts are available at <https://github.com/sahirbhatnagar/sail/tree/master/manuscript/bin>.

The expected run time to compile the manuscript is about 5 minutes on a standard desktop machine, assuming that you are using the pre-run simulation and bootstrap results.

### E.2.2. R Package Vignette

A website with two vignettes has been created for our sail package available at <https://sahirbhatnagar.com/sail/>

The 2 vignettes are:

1. <https://sahirbhatnagar.com/sail/articles/introduction-to-sail.html>
2. <https://sahirbhatnagar.com/sail/articles/user-defined-design.html>