

Sparse Additive Interaction Learning

Sahir R Bhatnagar^{1,2}, Yi Yang⁴, and Celia MT Greenwood^{1,2,5}

¹Department of Epidemiology, Biostatistics and Occupational Health, McGill
University

²Lady Davis Institute, Jewish General Hospital, Montréal, QC

⁴Department of Mathematics and Statistics, McGill University

⁵Departments of Oncology and Human Genetics, McGill University

May 2, 2018

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 are several applications where interaction models can reflect biological phenomena and improve statistical power. For example, genome wide association studies (GWAS) have been unable to explain a large proportion of heritability (the variance in phenotype attributable to genetic variants) and it has been suggested that this missing heritability may in part be due to gene-environment interactions [9]. Furthermore, diseases are now thought to be the result of entire biological networks whose states are af-

affected by environmental factors, and these systemic changes can induce or eliminate strong correlations between elements in a network. Therefore, we propose a multivariate penalization procedure for detecting interactions between high dimensional data ($p \gg n$) and an environmental factor, where the effect of this environmental factor on the high dimensional data is widespread and plays a role in predicting the response. Our approach improves on existing procedures for detecting such interactions in three ways; 1) it automatically enforces the strong heredity property, i.e., an interaction term can only be included in the model if the corresponding main effects are in the model 2) it reduces the dimensionality of the problem and leverages the high correlations by transforming the input feature space using network connectivity measures and 3) it leads to interpretable models which are biologically meaningful. This thesis is motivated by three studies; one that looks at the impact of maternal care on child development, another that characterizes normal brain development as a function of intelligence scores and one that looks at the long-lasting effects of diet in subsequent generations.

We consider a regression model for an outcome variable $Y = (Y_1, \dots, Y_n)$ where n is the number of subjects. Let $E = (E_1, \dots, E_n)$ be a binary or continuous environment vector and $\mathbf{X} = (X_1, \dots, X_n)^\top$ be the $n \times p$ matrix of high-dimensional data where $X_i = (X_{i1}, \dots, X_{ij}, \dots, X_{ip})$, and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)$ a vector of errors. Consider the regression model with main effects and their interactions with E :

$$g(\boldsymbol{\mu}) = \beta_0 + \sum_{j=1}^p \beta_j X_j + \beta_E E + \sum_{j=1}^p \alpha_j E X_j, \quad (1)$$

where $g(\cdot)$ is a known link function, $\boldsymbol{\mu} = \mathbb{E}[Y|\mathbf{X}, E, \boldsymbol{\beta}, \boldsymbol{\alpha}]$, and $\beta_0, \beta_j, \beta_E, \alpha_j$ are the true unknown model parameters for $j = 1, \dots, p$.

Due to the large number of parameters to estimate with respect to the number of observations, one commonly-used approach is to shrink the regression coefficients by placing a constraint on the values of $\boldsymbol{\beta}$ and $\boldsymbol{\alpha}$. For example, the LASSO [12] penalizes the squared

loss of the data with the L_1 -norm of the regression coefficients resulting in a method that performs both model selection and estimation. A natural extension of the LASSO to the interaction model (1) is given by:

$$\arg \min_{\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}} \frac{1}{2} \|Y - g(\boldsymbol{\mu})\|^2 + \lambda (\|\boldsymbol{\beta}\|_1 + \|\boldsymbol{\alpha}\|_1) \quad (2)$$

where $\|Y - g(\boldsymbol{\mu})\|^2 = \sum_i (y_i - g(\mu_i))^2$, $\|\boldsymbol{\beta}\|_1 = \sum_j |\beta_j|$, $\|\boldsymbol{\alpha}\|_1 = \sum_j |\alpha_j|$ and $\lambda \geq 0$ is a data driven tuning parameter that can set some of the coefficients to zero when sufficiently large.

However, since no constraint is placed on the structure of the model in 2, it is possible that the estimated main effects are zero while the interaction term is not. This has motivated methods that produce structured sparsity [1]. For example, Bien *et al.* [2] propose a strong hierarchical lasso which forces the main effects to be included if the interaction term is non-zero. However this method and related ones are restricted to all pairwise interactions between p measured variables. Here we concern ourselves with methods that impose a strong hierarchy in the context of gene environment interactions. We are interested in imposing the strong heredity principle [3]:

$$\hat{\alpha}_j \neq 0 \quad \Rightarrow \quad \hat{\beta}_j \neq 0 \quad \text{and} \quad \hat{\beta}_E \neq 0 \quad (3)$$

In words, the interaction term will only have a non-zero estimate if its corresponding main effects are estimated to be non-zero. One benefit brought by hierarchy is that the number of measured variables can be reduced, referred to as practical sparsity [2, 11]. For example, a model involving $X_1, E, X_1 \cdot E$ is more parsimonious than a model involving $X_1, E, X_2 \cdot E$, because in the first model a researcher would only have to measure two variables compared to three in the second model. In order to address these issues, we propose to extend the model of [4] to simultaneously perform variable selection, estimation and impose the strong heredity

principle in the context of high dimensional interactions with the environment ($\text{HD} \times E$). To do so, we follow Choi and reparametrize the coefficients for the interaction terms as $\alpha_j = \gamma_j \beta_j \beta_E$. Plugging this into (1):

$$g(\boldsymbol{\mu}) = \beta_0 + \sum_{j=1}^p \beta_j X_j + \beta_E E + \sum_{j=1}^p \gamma_j \beta_j \beta_E E X_j \quad (4)$$

This reparametrization directly enforces the strong heredity principle (3), i.e., if either main effect estimates are 0, then $\hat{\alpha}_j$ will be zero and a non-zero interaction coefficient implies non-zero $\hat{\beta}_j$ and $\hat{\beta}_E$. To perform variable selection in this new parametrization, we follow [4] and penalize $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$ instead of penalizing $\boldsymbol{\alpha}$ as in (2), leading to the following penalized least squares criterion:

$$\arg \min_{\beta_0, \boldsymbol{\beta}, \boldsymbol{\gamma}} \frac{1}{2} \|Y - g(\boldsymbol{\mu})\|^2 + \lambda_\beta \sum_{j=1}^p w_j |\beta_j| + \lambda_\gamma \sum_{j=1}^p w_{jE} |\gamma_{jE}| \quad (5)$$

where $g(\boldsymbol{\mu})$ is from (10), λ_β and λ_γ are tuning parameters and $\mathbf{w} = (w_1, \dots, w_p, w_{1E}, \dots, w_{pE})$ are prespecified adaptive weights. The λ_β tuning parameter controls the amount of shrinkage applied to the main effects, while λ_γ controls the interaction estimates and allows for the possibility of excluding the interaction term from the model even if the corresponding main effects are non-zero. The adaptive weights serve as a way of allowing parameters to be penalized differently. Furthermore, adaptive weighting [13] has been shown to construct oracle procedures [5], i.e., asymptotically, it performs as well as if the true model were given in advance. The oracle property is achieved when the weights are a function of any root- n consistent estimator of the true parameters e.g. maximum likelihood (MLE) or ridge regression estimates. It can be shown that the procedure in (5) asymptotically possesses the oracle property [4], even when the number of parameters tends to ∞ as the sample size

increases, if the weights are chosen such that

$$w_j = \left| \frac{1}{\hat{\beta}_j} \right|, \quad w_{jE} = \left| \frac{\hat{\beta}_j \hat{\beta}_E}{\hat{\alpha}_{jE}} \right| \quad \text{for } j = 1, \dots, q \quad (6)$$

where $\hat{\beta}_j$ and $\hat{\alpha}_j$ are the MLEs, from (1) or the ridge regression estimates when $p > n$. The rationale behind the data-dependent $\hat{\mathbf{w}}$ is that as the sample size grows, the weights for the truly zero predictors go to ∞ (which translates to a large penalty), whereas the weights for the truly non-zero predictors converge to a finite constant [13].

1.1 Toy example

We begin with a toy example to better illustrate our method. We sample $p = 20$ covariates independently from a $N(0, 1)$ truncated to the interval $[0, 1]$ and sample size $N = 100$. We generated data from the model

$$Y = f_1(X_1) + f_2(X_2) + 1.75E + 1.5E \cdot f_2(X_2) + \varepsilon \quad (7)$$

where $f_1(x) = -3x$, $f_2(x) = 2(2x - 1)^3$ and the error term ε is generated from a normal distribution with variance chosen such that the signal-to-noise ratio (SNR) is 2. We run the `sail` method with cubic b-splines and 10-fold CV to choose the optimal value of λ . Default values were used for all other arguments. We plot the solution path for both main effects and interactions in Figure ?? and the estimated functions \hat{f}_1 and \hat{f}_2 in Figure ??.

Type	Model	Software
Linear	CAP (Zhao et al. 2009, <i>Ann. Stat</i>)	X
	SHIM (Choi et al. 2009, <i>JASA</i>)	X
	hiernet (Bien et al. 2013, <i>Ann. Stat</i>)	hierNet(x, y)
	GRESH (She and Jiang 2014, <i>JASA</i>)	X
	FAMILY (Haris et al. 2014, <i>JCGS</i>)	FAMILY(x, z, y)
	glinternet (Lim and Hastie 2015, <i>JCGS</i>)	glinternet(x, y)
	RAMP (Hao et al. 2016, <i>JASA</i>)	RAMP(x, y)
Non-linear	VANISH (Radchenko and James 2010, <i>JASA</i>)	X
	sail (Bhatnagar et al. 2017+)	sail(x, e, y)

1.2 Existing Literature

2 Extension to Additive Models

Let $Y = (Y_1, \dots, Y_n) \in \mathbb{R}^n$ be a continuous outcome variable, $X_E = (E_1, \dots, E_n) \in \mathbb{R}^n$ a binary or continuous environment vector, $\mathbf{X} = (X_1, \dots, X_p) \in \mathbb{R}^{n \times p}$ a matrix of predictors, and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$ a vector of i.i.d random variables with mean 0. 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:

$$f_j(X_j) = \sum_{\ell=1}^{m_j} \psi_{j\ell}(X_j) \beta_{j\ell} \quad (8)$$

Here, the $\{\psi_{j\ell}\}_1^{m_j}$ are a family of basis functions in X_j [6]. Let Ψ_j be the $n \times m_j$ matrix of evaluations of the $\psi_{j\ell}$ and $\boldsymbol{\theta}_j = (\beta_{j1}, \dots, \beta_{jm_j}) \in \mathbb{R}^{m_j}$ for $j = 1, \dots, p$, i.e., $\boldsymbol{\theta}_j$ 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

$$Y = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j + \beta_E X_E + \sum_{j=1}^p (X_E \circ \Psi_j) \boldsymbol{\alpha}_j + \varepsilon \quad (9)$$

where β_0 is the intercept, β_E is the coefficient for the environment variable, $\boldsymbol{\alpha}_j = (\alpha_{j1}, \dots, \alpha_{jm_j}) \in \mathbb{R}^{m_j}$ are the basis coefficients for the j th interaction term and $(X_E \circ \Psi_j)$ is the $n \times m_j$ matrix formed by the component-wise multiplication of the column vector X_E by each column of Ψ_j . To enforce the strong heredity property, we reparametrize the coefficients for the interaction terms in (9) as $\boldsymbol{\alpha}_j = \gamma_j \beta_E \boldsymbol{\theta}_j$:

$$Y = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j + \beta_E X_E + \sum_{j=1}^p \gamma_j \beta_E (X_E \circ \Psi_j) \boldsymbol{\theta}_j + \varepsilon \quad (10)$$

For a continuous response, we use the squared-error loss:

$$\mathcal{L}(Y; \boldsymbol{\Theta}) = \frac{1}{2n} \left\| Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j - \beta_E X_E - \sum_{j=1}^p \gamma_j \beta_E (X_E \circ \Psi_j) \boldsymbol{\theta}_j \right\|_2^2 \quad (11)$$

where $\boldsymbol{\Theta} := (\beta_0, \beta_E, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p, \gamma_1, \dots, \gamma_p)$.

We consider the following penalized least squares criterion for this problem:

$$\arg \min_{\boldsymbol{\Theta}} \mathcal{L}(Y; \boldsymbol{\Theta}) + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda \alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (12)$$

where $\lambda > 0$ and $\alpha \in (0, 1)$ are tuning parameters and w_E, w_j, w_{jE} are adaptive weights for $j = 1, \dots, p$. These weights serve as a way of allowing parameters to be penalized differently. Furthermore, adaptive weighting [13] has been shown to construct oracle procedures [5], i.e., asymptotically, it performs as well as if the true model were given in advance. These weights

are given by

$$w_E = \left| \frac{1}{\hat{\beta}_E} \right|, \quad w_j = \frac{1}{\|\hat{\boldsymbol{\theta}}_j\|_2}, \quad w_{jE} = \left| \frac{\hat{\beta}_E \|\hat{\boldsymbol{\theta}}_j\|_2}{\|\hat{\boldsymbol{\alpha}}_j\|_2} \right| \quad \text{for } j = 1, \dots, p \quad (13)$$

where $\hat{\beta}_E$, $\hat{\boldsymbol{\theta}}_j$ and $\hat{\boldsymbol{\alpha}}_j$ are the MLEs, from (9) or the ridge regression estimates when $p > n$.

3 Regularization Path

The **sail** model has the form

$$\hat{Y} = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \boldsymbol{\Psi}_j \boldsymbol{\theta}_j + \beta_E X_E + \sum_{j=1}^p \gamma_j \beta_E (X_E \circ \boldsymbol{\Psi}_j) \boldsymbol{\theta}_j \quad (14)$$

The objective function is given by

$$Q(\boldsymbol{\Theta}) = \frac{1}{2n} \|Y - \hat{Y}\|_2^2 + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (15)$$

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

$$\frac{\partial Q}{\partial \beta_0} = \frac{1}{n} \left(Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \boldsymbol{\Psi}_j \boldsymbol{\theta}_j - \beta_E X_E - \sum_{j=1}^p \gamma_j \beta_E (X_E \circ \boldsymbol{\Psi}_j) \boldsymbol{\theta}_j \right)^\top \mathbf{1} = 0 \quad (16)$$

$$\frac{\partial Q}{\partial \beta_E} = -\frac{1}{n} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \boldsymbol{\Psi}_j) \boldsymbol{\theta}_j \right)^\top R + \lambda(1 - \alpha) w_E s_1 = 0 \quad (17)$$

$$\frac{\partial Q}{\partial \boldsymbol{\theta}_j} = -\frac{1}{n} (\boldsymbol{\Psi}_j + \gamma_j \beta_E (X_E \circ \boldsymbol{\Psi}_j))^\top R + \lambda(1 - \alpha) w_j s_2 = \mathbf{0} \quad (18)$$

$$\frac{\partial Q}{\partial \gamma_j} = -\frac{1}{n} (\beta_E (X_E \circ \boldsymbol{\Psi}_j) \boldsymbol{\theta}_j)^\top R + \lambda\alpha w_{jE} s_3 = 0 \quad (19)$$

where s_1 is in the subgradient of the ℓ_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 ℓ_2 norm:

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

and s_3 is in the subgradient of the ℓ_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

$$R_{(-j)} = Y - \beta_0 \cdot \mathbf{1} - \sum_{\ell \neq j} \boldsymbol{\Psi}_\ell \boldsymbol{\theta}_\ell - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell \beta_E (X_E \circ \boldsymbol{\Psi}_\ell) \boldsymbol{\theta}_\ell$$

the partial residual without X_E as

$$R_{(-E)} = Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \boldsymbol{\Psi}_j \boldsymbol{\theta}_j$$

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

$$R_{(-jE)} = Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \boldsymbol{\Psi}_j \boldsymbol{\theta}_j - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell \beta_E (X_E \circ \boldsymbol{\Psi}_\ell) \boldsymbol{\theta}_\ell$$

From the subgradient Equation (17), we see that $\beta_E = 0$ is a solution if

$$\frac{1}{w_E} \left| \frac{1}{n} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^\top R_{(-E)} \right| \leq \lambda(1 - \alpha) \quad (20)$$

From the subgradient Equation (18), we see that $\theta_j = \mathbf{0}$ is a solution if

$$\frac{1}{w_j} \left\| \frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-j)} \right\|_2 \leq \lambda(1 - \alpha) \quad (21)$$

From the subgradient Equation (19), we see that $\gamma_j = 0$ is a solution if

$$\frac{1}{w_{jE}} \left| \frac{1}{n} (\beta_E (X_E \circ \Psi_j) \theta_j)^\top R_{(-jE)} \right| \leq \lambda \alpha \quad (22)$$

3.1 Lambda Max

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

$$\lambda_{max} = \frac{1}{n(1 - \alpha)} \max \left\{ \frac{1}{w_E} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^\top R_{(-E)}, \max_j \frac{1}{w_j} \left\| (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-j)} \right\|_2 \right\} \quad (23)$$

which reduces to

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

3.2 Optimization of Parameters

From the subgradient equations we see that

$$\hat{\beta}_0 = \left(Y - \sum_{j=1}^p \Psi_j \hat{\theta}_j - \hat{\beta}_E X_E - \sum_{j=1}^p \hat{\gamma}_j \hat{\beta}_E (X_E \circ \Psi_j) \hat{\theta}_j \right)^\top \mathbf{1} \quad (24)$$

$$\hat{\beta}_E = S \left(\frac{1}{n \cdot w_E} \left(X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \hat{\theta}_j \right)^\top R_{(-E)}, \lambda(1 - \alpha) \right) \quad (25)$$

$$\lambda(1 - \alpha) w_j \frac{\theta_j}{\|\theta_j\|_2} = \frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-j)} \quad (26)$$

$$\hat{\gamma}_j = S \left(\frac{1}{n \cdot w_{jE}} (\beta_E (X_E \circ \Psi_j) \theta_j)^\top R_{(-jE)}, \lambda \alpha \right) \quad (27)$$

where $S(x, t) = \text{sign}(x)(|x| - t)$ is the soft-thresholding operator

4 Algorithm

For each function f_j , we use a cubic B-spline parameterization with 5 degrees of freedom implemented in the `bs` function in R [10].

4.1 Details on update for θ

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

$$R_s = Y - \tilde{Y}_{(-s)}^{(k)} \quad (28)$$

Algorithm 1 Coordinate descent for least-squares **sail** with strong heredity

```

1: function sail( $Y, \mathbf{X}, X_E, \text{df}, \text{degree}, \epsilon$ ) ▷ Algorithm for solving (15)
2:    $\Psi_j \leftarrow \text{splines}::\text{bs}(X_j, \text{df}, \text{degree})$  for  $j = 1, \dots, p$ 
3:    $\tilde{\Psi}_j \leftarrow X_E \circ \Psi_j$  for  $j = 1, \dots, p$ 
4:   Initialize:  $\beta_0^{(0)} \leftarrow \bar{Y}$ ,  $\beta_E^{(0)} = \boldsymbol{\theta}_j^{(0)} \leftarrow 0$  for  $j = 1, \dots, p$ .
5:   Set iteration counter  $k \leftarrow 0$ 
6:    $R^* \leftarrow Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_j (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \boldsymbol{\theta}_j^{(k)}$ 
7:   repeat
8:     • To update  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$ 
9:        $\tilde{X}_j \leftarrow \beta_E^{(k)} \tilde{\Psi}_j \boldsymbol{\theta}_j^{(k)}$  for  $j = 1, \dots, p$ 
10:       $R \leftarrow R^* + \sum_{j=1}^p \gamma_j^{(k)} \tilde{X}_j$ 
11:
12:      
$$\boldsymbol{\gamma}^{(k)(new)} \leftarrow \arg \min_{\boldsymbol{\gamma}} \frac{1}{2n} \left\| R - \sum_j \gamma_j \tilde{X}_j \right\|_2^2 + \lambda \alpha \sum_j w_{jE} |\gamma_j|$$

13:
14:       $\Delta = \sum_j (\gamma_j^{(k)} - \gamma_j^{(k)(new)}) \tilde{X}_j$ 
15:       $R^* \leftarrow R^* + \Delta$ 
16:     • To update  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)$ 
17:        $\tilde{X}_j \leftarrow \Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j$  for  $j = 1, \dots, p$ 
18:       for  $j = 1, \dots, p$  do
19:          $R \leftarrow R^* + \tilde{X}_j \boldsymbol{\theta}_j^{(k)}$ 
20:
21:         
$$\boldsymbol{\theta}_j^{(k)(new)} \leftarrow \arg \min_{\boldsymbol{\theta}_j} \frac{1}{2n} \left\| R - \tilde{X}_j \boldsymbol{\theta}_j \right\|_2^2 + \lambda (1 - \alpha) w_j \|\boldsymbol{\theta}_j\|_2$$

22:
23:          $\Delta = \tilde{X}_j (\boldsymbol{\theta}_j^{(k)} - \boldsymbol{\theta}_j^{(k)(new)})$ 
24:          $R^* \leftarrow R^* + \Delta$ 
25:     • To update  $\beta_E$ 
26:        $\tilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \tilde{\Psi}_j \boldsymbol{\theta}_j^{(k)}$ 
27:        $R \leftarrow R^* + \beta_E^{(k)} \tilde{X}_E$ 
28:
29:       
$$\beta_E^{(k)(new)} \leftarrow S \left( \frac{1}{n \cdot w_E} \tilde{X}_E^\top R, \lambda (1 - \alpha) \right)$$

30:       ▷  $S(x, t) = \text{sign}(x)(|x| - t)_+$ 
31:
32:        $\Delta = (\beta_E^{(k)} - \beta_E^{(k)(new)}) \tilde{X}_E$ 
33:        $R^* \leftarrow R^* + \Delta$ 
34:     • To update  $\beta_0$ 
35:        $R \leftarrow R^* + \beta_0^{(k)}$ 
36:
37:       
$$\beta_0^{(k)(new)} \leftarrow \frac{1}{n} R^* \cdot \mathbf{1}$$

38:
39:        $\Delta = \beta_0^{(k)} - \beta_0^{(k)(new)}$ 
40:        $R^* \leftarrow R^* + \Delta$ 
41:        $k \leftarrow k + 1$ 
42:   until convergence criterion is satisfied:  $\left\| \boldsymbol{\Theta}^{(k)} - \boldsymbol{\Theta}^{(k-1)} \right\|_2^2 < \epsilon$ 

```

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

$$\tilde{Y}_{(-s)}^{(k)} = \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_{\ell \neq s} \Psi_\ell \theta_\ell^{(k)} - \sum_{\ell \neq s} \gamma_\ell^{(k)} \beta_E^{(k)} \tilde{\Psi}_\ell \theta_\ell^{(k)} \quad (29)$$

Using (29), (28) can be re-written as

$$\begin{aligned} R_s &= Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \\ &= R^* + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \end{aligned} \quad (30)$$

where

$$R^* = Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)} \quad (31)$$

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

$$\theta_s^{(k)(new)} = \arg \min_{\theta_j} \frac{1}{2n} \left\| R_s - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_j \right\|_2^2 + \lambda(1 - \alpha) w_s \|\theta_j\|_2 \quad (32)$$

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

$$R_{s+1} = R^* + (\Psi_{s+1} + \gamma_{s+1}^{(k)} \beta_E^{(k)} \tilde{\Psi}_{s+1}) \theta_{s+1}^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) (\theta_s^{(k)} - \theta_s^{(k)(new)}) \quad (33)$$

where R^* is given by (31), $\theta_s^{(k)}$ is the parameter value prior to the update, and $\theta_s^{(k)(new)}$ is the updated value given by (32). Taking the difference between (30) and (33) gives

$$\begin{aligned} \Delta &= R_t - R_s \\ &= (\Psi_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\Psi}_t) \theta_t^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) (\theta_s^{(k)} - \theta_s^{(k)(new)}) - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \\ &= (\Psi_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\Psi}_t) \theta_t^{(k)} - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)(new)} \end{aligned} \quad (34)$$

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 Δ , given by (34). This formulation can lead to computational speedups especially when $\Delta = 0$, meaning the partial residual does not need to be re-calculated.

5 Weak Heredity

We can also enforce the weak heredity property:

$$\hat{\alpha}_{jE} \neq 0 \quad \Rightarrow \quad \hat{\beta}_j \neq 0 \quad \text{or} \quad \hat{\beta}_E \neq 0 \quad (35)$$

That is, an interaction term can only be present if at least one of it's corresponding main effects is nonzero. To do so, we reparametrize the coefficients for the interaction terms in (9) as $\alpha_j = \gamma_j(\beta_E \cdot \mathbf{1}_{m_j} + \theta_j)$, where $\mathbf{1}_{m_j}$ is a vector of ones with dimension m_j (i.e. the length of θ_j).

5.1 Regularization Path

The `sail` model with weak heredity has the form

$$\hat{Y} = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \Psi_j \theta_j + \beta_E X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \quad (36)$$

The objective function is given by

$$Q(\Theta) = \frac{1}{2n} \|Y - \hat{Y}\|_2^2 + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (37)$$

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

$$\frac{\partial Q}{\partial \beta_0} = \frac{1}{n} \left(Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j - \beta_E X_E - \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j) \right)^\top \mathbf{1} = 0 \quad (38)$$

$$\frac{\partial Q}{\partial \beta_E} = -\frac{1}{n} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^\top R + \lambda(1 - \alpha) w_E s_1 = 0 \quad (39)$$

$$\frac{\partial Q}{\partial \boldsymbol{\theta}_j} = -\frac{1}{n} (\Psi_j + \gamma_j (X_E \circ \Psi_j))^\top R + \lambda(1 - \alpha) w_j s_2 = \mathbf{0} \quad (40)$$

$$\frac{\partial Q}{\partial \gamma_j} = -\frac{1}{n} ((X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j))^\top R + \lambda \alpha w_{jE} s_3 = 0 \quad (41)$$

where s_1 is in the subgradient of the ℓ_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 ℓ_2 norm:

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

and s_3 is in the subgradient of the ℓ_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

$$R_{(-j)} = Y - \beta_0 \cdot \mathbf{1} - \sum_{\ell \neq j} \Psi_\ell \boldsymbol{\theta}_\ell - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell (X_E \circ \Psi_\ell) (\beta_E \cdot \mathbf{1}_{m_\ell} + \boldsymbol{\theta}_\ell)$$

the partial residual without X_E as

$$R_{(-E)} = Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j - \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \boldsymbol{\theta}_j$$

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

$$R_{(-jE)} = Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell (X_E \circ \Psi_\ell) (\beta_E \cdot \mathbf{1}_{m_\ell} + \boldsymbol{\theta}_\ell)$$

From the subgradient Equation (39), we see that $\beta_E = 0$ is a solution if

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

From the subgradient Equation (40), we see that $\boldsymbol{\theta}_j = \mathbf{0}$ is a solution if

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

From the subgradient Equation (41), we see that $\gamma_j = 0$ is a solution if

$$\frac{1}{w_{jE}} \left| \frac{1}{n} ((X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j))^\top R_{(-jE)} \right| \leq \lambda\alpha \quad (44)$$

5.2 Lambda Max

The smallest value of λ for which the entire parameter vector $(\beta_E, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p, \gamma_1, \dots, \gamma_p)$ is $\mathbf{0}$ is:

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

which reduces to

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

5.3 Optimization of Parameters

From the subgradient equations we see that

$$\hat{\beta}_0 = \left(Y - \sum_{j=1}^p \Psi_j \hat{\theta}_j - \hat{\beta}_E X_E - \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) (\hat{\beta}_E \cdot \mathbf{1}_{m_j} + \hat{\theta}_j) \right)^\top \mathbf{1} \quad (46)$$

$$\hat{\beta}_E = S \left(\frac{1}{n \cdot w_E} \left(X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^\top R_{(-E)}, \lambda(1-\alpha) \right) \quad (47)$$

$$\lambda(1-\alpha)w_j \frac{\theta_j}{\|\theta_j\|_2} = \frac{1}{n} (\Psi_j + \gamma_j (X_E \circ \Psi_j))^\top R_{(-j)} \quad (48)$$

$$\hat{\gamma}_j = S \left(\frac{1}{n \cdot w_{jE}} \left((X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j) \right)^\top R_{(-jE)}, \lambda\alpha \right) \quad (49)$$

where $S(x, t) = \text{sign}(x)(|x| - t)$ is the soft-thresholding operator

Algorithm 2 Coordinate descent for least-squares **sail** with weak heredity

```

1: function sail( $Y, \mathbf{X}, X_E, \text{df}, \text{degree}, \epsilon$ ) ▷ Algorithm for solving (37)
2:    $\Psi_j \leftarrow \text{splines::bs}(X_j, \text{df}, \text{degree})$  for  $j = 1, \dots, p$ 
3:    $\tilde{\Psi}_j \leftarrow X_E \circ \Psi_j$  for  $j = 1, \dots, p$ 
4:   Initialize:  $\beta_0^{(0)} \leftarrow \bar{Y}$ ,  $\beta_E^{(0)} = \boldsymbol{\theta}_j^{(0)} \leftarrow 0$  for  $j = 1, \dots, p$ .
5:   Set iteration counter  $k \leftarrow 0$ 
6:    $R^* \leftarrow Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_j \Psi_j \boldsymbol{\theta}_j^{(k)} - \sum_j \gamma_j^{(k)} \tilde{\Psi}_j(\beta_E^{(k)}) \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j^{(k)}$ 
7:   repeat
8:     • To update  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$ 
9:        $\tilde{X}_j \leftarrow \tilde{\Psi}_j(\beta_E^{(k)}) \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j^{(k)}$  for  $j = 1, \dots, p$ 
10:       $R \leftarrow R^* + \sum_{j=1}^p \gamma_j^{(k)} \tilde{X}_j$ 
11:
12:      
$$\boldsymbol{\gamma}^{(k)(new)} \leftarrow \arg \min_{\boldsymbol{\gamma}} \frac{1}{2n} \left\| R - \sum_j \gamma_j \tilde{X}_j \right\|_2^2 + \lambda \alpha \sum_j w_{jE} |\gamma_j|$$

13:       $\Delta = \sum_j (\gamma_j^{(k)} - \gamma_j^{(k)(new)}) \tilde{X}_j$ 
14:       $R^* \leftarrow R^* + \Delta$ 
15:     • To update  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)$ 
16:        $\tilde{X}_j \leftarrow \Psi_j + \gamma_j^{(k)} \tilde{\Psi}_j$  for  $j = 1, \dots, p$ 
17:       for  $j = 1, \dots, p$  do
18:          $R \leftarrow R^* + \tilde{X}_j \boldsymbol{\theta}_j^{(k)}$ 
19:
20:         
$$\boldsymbol{\theta}_j^{(k)(new)} \leftarrow \arg \min_{\boldsymbol{\theta}_j} \frac{1}{2n} \left\| R - \tilde{X}_j \boldsymbol{\theta}_j \right\|_2^2 + \lambda(1 - \alpha) w_j \|\boldsymbol{\theta}_j\|_2$$

21:          $\Delta = \tilde{X}_j (\boldsymbol{\theta}_j^{(k)} - \boldsymbol{\theta}_j^{(k)(new)})$ 
22:          $R^* \leftarrow R^* + \Delta$ 
23:     • To update  $\beta_E$ 
24:        $\tilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \tilde{\Psi}_j \mathbf{1}_{m_j}$ 
25:        $R \leftarrow R^* + \beta_E^{(k)} \tilde{X}_E$ 
26:
27:       
$$\beta_E^{(k)(new)} \leftarrow S \left( \frac{1}{n \cdot w_E} \tilde{X}_E^\top R, \lambda(1 - \alpha) \right)$$

28:       ▷  $S(x, t) = \text{sign}(x)(|x| - t)_+$ 
29:
30:        $\Delta = (\beta_E^{(k)} - \beta_E^{(k)(new)}) \tilde{X}_E$ 
31:        $R^* \leftarrow R^* + \Delta$ 
32:     • To update  $\beta_0$ 
33:        $R \leftarrow R^* + \beta_0^{(k)}$ 
34:
35:       
$$\beta_0^{(k)(new)} \leftarrow \frac{1}{n} R^* \cdot \mathbf{1}$$

36:
37:        $\Delta = \beta_0^{(k)} - \beta_0^{(k)(new)}$ 
38:        $R^* \leftarrow R^* + \Delta$ 
39:        $k \leftarrow k + 1$ 
40:   until convergence criterion is satisfied:  $\left\| \boldsymbol{\Theta}^{(k)} - \boldsymbol{\Theta}^{(k-1)} \right\|_2^2 < \epsilon$ 

```

5.4 Algorithm

6 Simulations

The covariates are simulated as follows. First, we generate w_1, \dots, w_p, u, v independently from a standard normal distribution truncated to the interval $[0,1]$ for $i = 1, \dots, n$. Then we set $x_j = (w_j + t \cdot u)/(1+t)$ for $j = 1, \dots, 4$ and $x_j = (w_j + t \cdot v)/(1+t)$ for $j = 5, \dots, p$, where the parameter t controls the amount of correlation among predictors. This leads to a compound symmetry correlation structure where $\text{Corr}(x_j, x_k) = t^2/(1+t^2)$, for $1 \leq j \leq 4, 1 \leq k \leq 4$, and $\text{Corr}(x_j, x_k) = t^2/(1+t^2)$, for $5 \leq j \leq p, 5 \leq k \leq p$, but the covariates of the nonzero and zero components are independent [7, 8]

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.

1. Hierarchy

(a) Truth obeys strong hierarchy.

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

(b) Truth obeys weak hierarchy.

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

(c) Truth only has interactions.

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

2. Non-linearity

(a) Truth is linear

$$Y = \sum_{j=1}^4 \beta_j X_j + \beta_E \cdot X_E + X_E \cdot X_3 + X_E \cdot X_4 + \varepsilon$$

3. Interactions

(a) Truth only has main effects

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

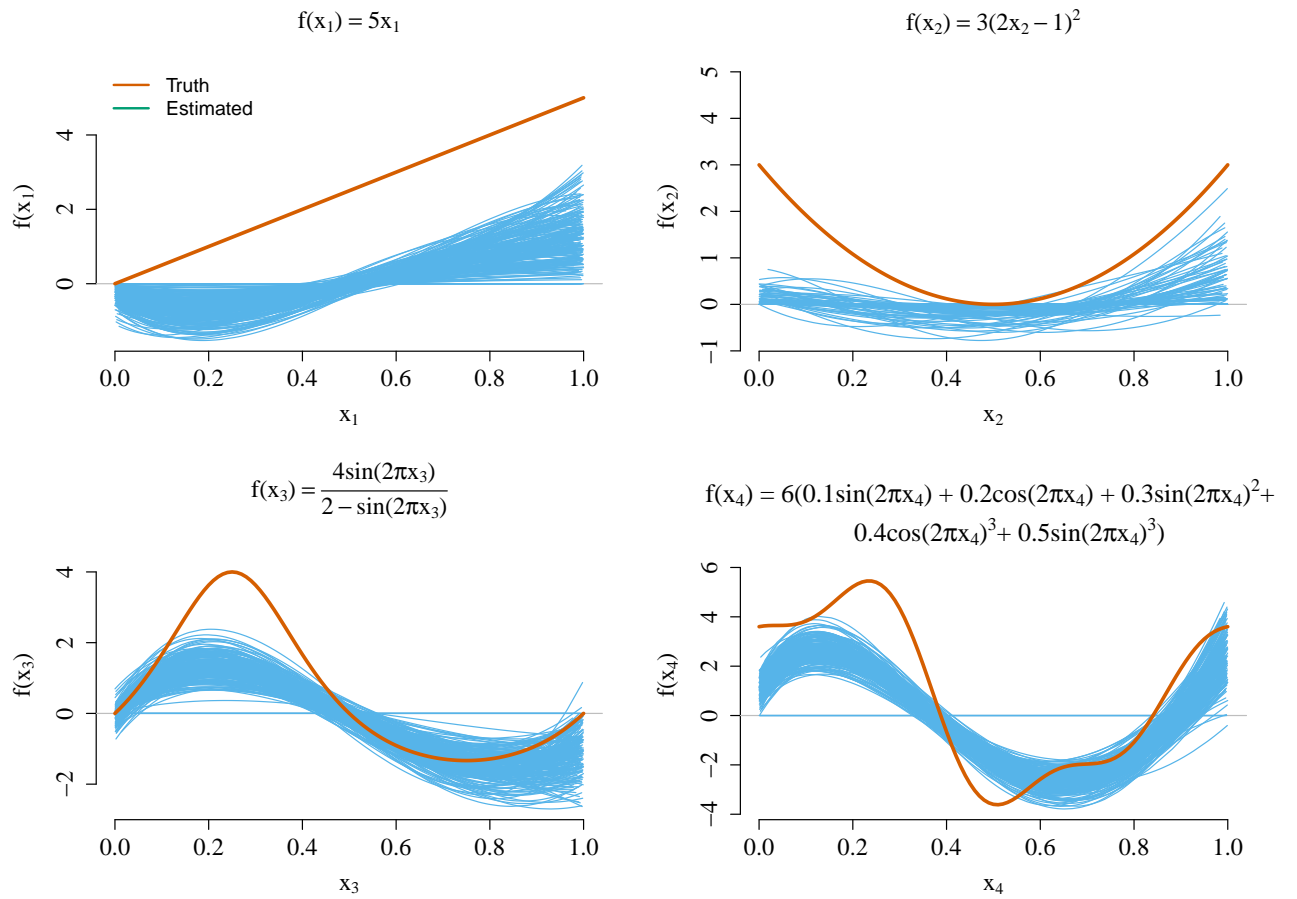


Figure 1: text

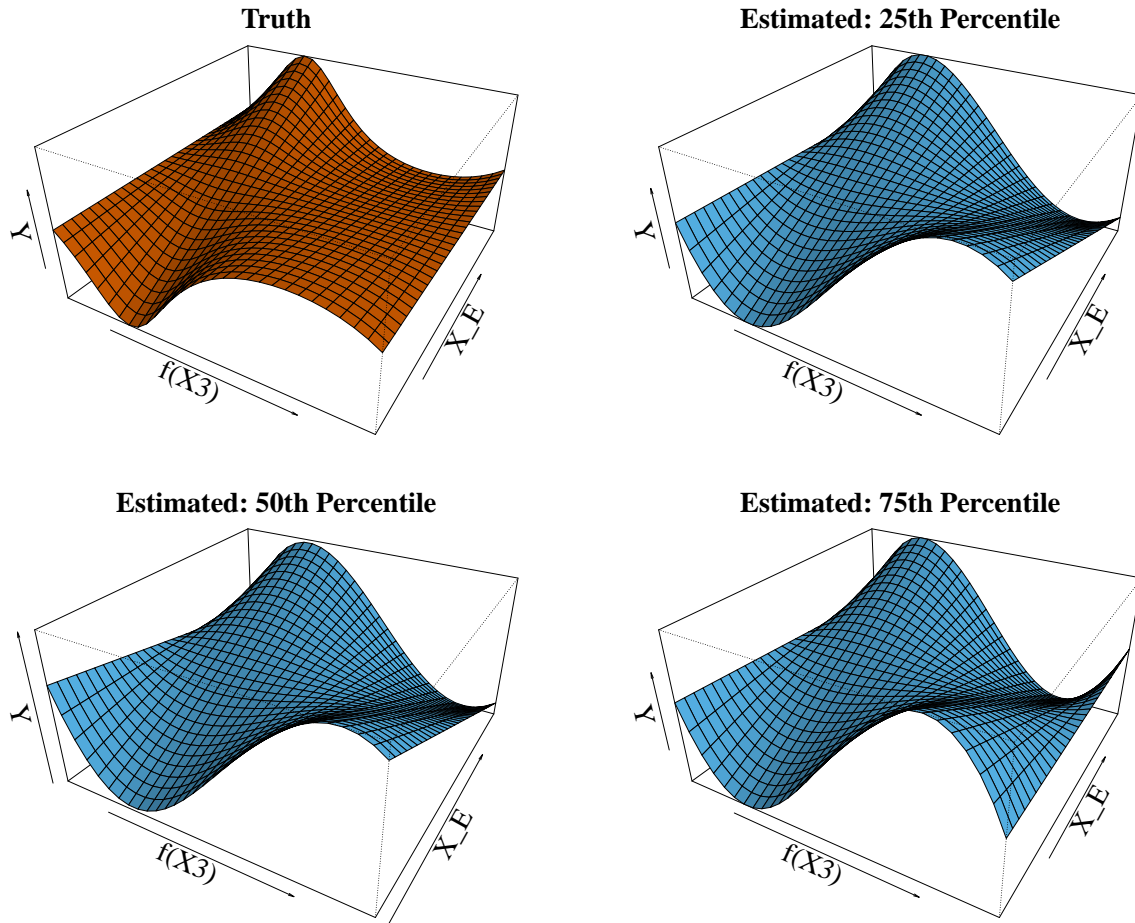


Figure 2: True and estimated interaction effects for $E \cdot f(X_3)$ in simulation scenario 1a).

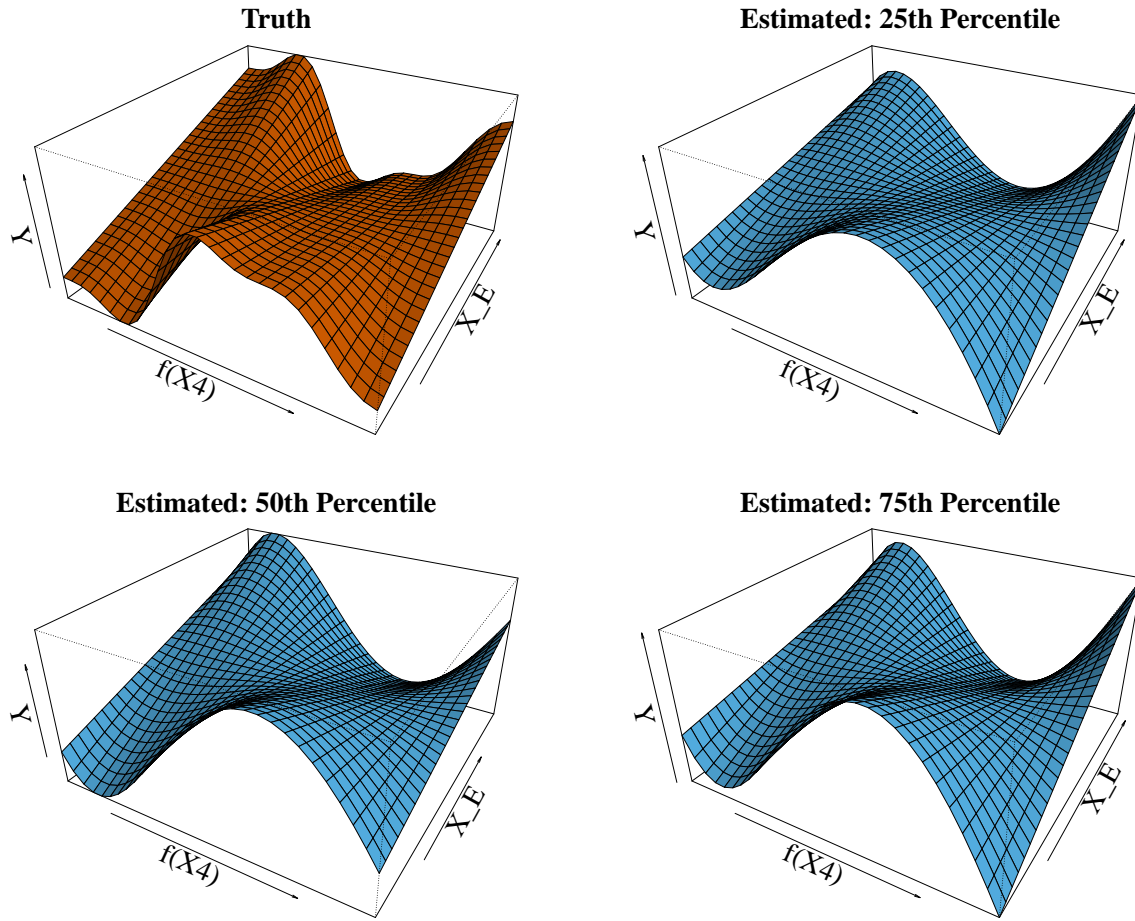


Figure 3: True and estimated interaction effects for $E \cdot f(X_4)$ in simulation scenario 1a).

7 Real Data Application

References

- [1] F. Bach, R. Jenatton, J. Mairal, G. Obozinski, et al. Structured sparsity through convex optimization. *Statistical Science*, 27(4):450–468, 2012.
- [2] J. Bien, J. Taylor, R. Tibshirani, et al. A lasso for hierarchical interactions. *The Annals of Statistics*, 41(3):1111–1141, 2013.
- [3] H. Chipman. Bayesian variable selection with related predictors. *Canadian Journal of Statistics*, 24(1):17–36, 1996.
- [4] N. H. Choi, W. Li, and J. Zhu. Variable selection with the strong heredity constraint and its oracle property. *Journal of the American Statistical Association*, 105(489):354–364, 2010.
- [5] J. Fan and R. Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001.
- [6] T. Hastie, R. Tibshirani, and M. Wainwright. *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press, 2015.
- [7] J. Huang, J. L. Horowitz, and F. Wei. Variable selection in nonparametric additive models. *Annals of statistics*, 38(4):2282, 2010.
- [8] Y. Lin, H. H. Zhang, et al. Component selection and smoothing in multivariate non-parametric regression. *The Annals of Statistics*, 34(5):2272–2297, 2006.
- [9] T. A. Manolio, F. S. Collins, N. J. Cox, D. B. Goldstein, L. A. Hindorff, D. J. Hunter, M. I. McCarthy, E. M. Ramos, L. R. Cardon, A. Chakravarti, et al. Finding the missing heritability of complex diseases. *Nature*, 461(7265):747–753, 2009.
- [10] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2017.

- [11] Y. She and H. Jiang. Group regularized estimation under structural hierarchy. *arXiv preprint arXiv:1411.4691*, 2014.
- [12] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.
- [13] H. Zou. The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006.