

Sparse Additive Interaction Learning

Sahir R Bhatnagar^{1,2}, Amanda Lovato², 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

August 20, 2019

Abstract

Diseases are now thought to be the result of changes in entire biological networks whose states are affected by a complex interaction of genetic and environmental factors. In general, power to estimate interactions is low, the number of possible interactions could be enormous and their effects may be non-linear. Existing approaches such as the lasso might keep an interaction but remove a main effect, which is problematic for interpretation. In this work, we introduce a sparse additive interaction learning model called **sail** for detecting non-linear interactions with a key environmental or exposure variable in high-dimensional settings. Our method can accommodate either the strong or weak heredity constraints. 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 error, sensitivity and specificity when there are non-linear interactions with an exposure variable. We apply **sail** to the Alzheimer’s Disease Neuroimaging Initiative (ADNI) data to select non-linear interactions between clinical diagnosis and $A\beta$ protein in 96 brain regions on mini-mental state examination. Our algorithms are available in an R package (<https://github.com/greenwoodlab>).

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 [1]. Accurate capture of interactions may hold the potential for better understanding biological phenomena and improving prediction accuracy. For example, a model that considered interactions between brain imaging data and genetic features had better classification accuracy compared to a model that considered the main effects only [2]. Furthermore, the manifestations of disease are 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 [3]. However, there is a general deficit of such replicated interactions in the literature [4]. Indeed, power to detect interactions is always lower than for main effects, and in high-dimensional settings ($p \gg n$), this lack of power to detect interactions is exacerbated, since the number of possible interactions could be enormous and their effects may be non-linear. Hence, analytic methods that may improve power are essential.

Interactions may occur in numerous types and of varying complexities. In this paper, we consider one specific type of interaction models, where one (exposure) variable 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 \mathbf{X} and E .

1.1 A Sparse additive interaction model

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, and $\mathbf{X} = (X_1, \dots, X_p) \in \mathbb{R}^{n \times p}$ a matrix of predictors, possibly high-dimensional. 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 $\{\psi_{j\ell}\}_1^{m_j}$ are a family of basis functions in X_j [5]. 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 $j = 1, \dots, p$ ($\boldsymbol{\theta}_j$ 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

$$Y = \beta_0 \cdot \mathbf{1}_n + \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j + \beta_E X_E + \sum_{j=1}^p (X_E \circ \Psi_j) \boldsymbol{\tau}_j + \varepsilon \quad (2)$$

where $\beta_0 \in \mathbb{R}$ is the intercept, $\beta_E \in \mathbb{R}$ is the coefficient for the environment variable, $\boldsymbol{\tau}_j = (\tau_{j1}, \dots, \tau_{jm_j}) \in \mathbb{R}^{m_j}$ 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 Ψ_j , and $\varepsilon \in \mathbb{R}^n$ is a vector of iid 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 is to shrink the regression coefficients by placing a constraint on the values of $(\beta_E, \boldsymbol{\theta}_j, \boldsymbol{\tau}_j)$. 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 [6]. 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 [7]. In light of these goals, we consider the following penalized objective function:

$$Q(\boldsymbol{\Theta}) = -L(\boldsymbol{\Theta}) + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} \|\boldsymbol{\tau}_j\|_2 \quad (3)$$

where $\boldsymbol{\Theta} = (\beta_0, \beta_E, \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p, \boldsymbol{\tau}_1, \dots, \boldsymbol{\tau}_p)$, $L(\boldsymbol{\Theta})$ is the log-likelihood function of the observations $\mathbf{V}_i = (Y_i, \boldsymbol{\Psi}_i, X_{iE})$ for $i = 1, \dots, n$, $\|\boldsymbol{\theta}_j\|_2 = \sqrt{\sum_{k=1}^{m_j} \beta_{jk}^2}$, $\|\boldsymbol{\tau}_j\|_2 = \sqrt{\sum_{k=1}^{m_j} \tau_{jk}^2}$, $\lambda > 0$ 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. The first term in the penalty penalizes the main effects while the second term penalizes the interactions. The parameter α 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 nonzero 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 [8]. This is known as the strong heredity principle [9]. Indeed, large main effects are more likely to lead to detectable interactions [10]. In the next section we discuss how a simple reparametrization of the model (3) can lead to this desirable property.

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. 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 [9]. In the context of penalized regression methods, these principles can be formulated as structured sparsity [11] problems. Several authors have proposed to modify the type of penalty in order to achieve the heredity principle [12, 13, 14, 15]. We take an alternative approach. Following Choi et al. [16], we introduce a new set of parameters $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p) \in \mathbb{R}^p$ and reparametrize the coefficients for the interaction terms $\boldsymbol{\tau}_j$ in (2) as a function of γ_j and the main effect parameters $\boldsymbol{\theta}_j$ and β_E . This reparametrization for both strong and weak heredity is summarized in Table 1.

Table 1: Reparametrization for strong and weak heredity principle for **sail** model

Type	Feature	Reparametrization
Strong heredity	$\hat{\tau}_j \neq 0$ only if $\hat{\boldsymbol{\theta}}_j \neq 0$ and $\hat{\beta}_E \neq 0$	$\boldsymbol{\tau}_j = \gamma_j \beta_E \boldsymbol{\theta}_j$
Weak heredity	$\hat{\tau}_j \neq 0$ only if $\hat{\boldsymbol{\theta}}_j \neq 0$ or $\hat{\beta}_E \neq 0$	$\boldsymbol{\tau}_j = \gamma_j (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j)$

To perform variable selection in this new parametrization, we penalize $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$ instead of penalizing $\boldsymbol{\tau}$ as in (3), leading to the following penalized objective function:

$$Q(\boldsymbol{\Theta}) = -L(\boldsymbol{\Theta}) + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (4)$$

An estimate of the regression parameters is given by $\hat{\boldsymbol{\Theta}} = \arg \min_{\boldsymbol{\Theta}} Q(\boldsymbol{\Theta})$. 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 α would lead to more interactions being included in the final model while values approaching 1 would favor main effects. Similar to the elastic net [17], we fix α and obtain a solution path over a sequence

of λ values.

1.3 Toy example

We present here a toy example to better illustrate our method. 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 :

$$Y = f_1(X_1) + f_2(X_2) + 1.75E + 1.5E \cdot f_2(X_2) + \varepsilon \quad (5)$$

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 ε 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 with cubic B-splines 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 were used for all other arguments. We plot the solution path for both main effects and interactions in Figure ??, 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 nonzero if the main effects E and $f_2(X_2)$ (black and orange lines respectively in the top panel) are nonzero, while nonzero main effects doesn't necessarily imply a nonzero interaction.

In Figure ??, 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, but the means are not well aligned with the data. Lack-of-fit for $f_1(X_1)$ can be partially explained by acknowledging that `sail` is trying to fit a cubic spline to a linear function. Nevertheless, this example demonstrates that `sail` can still identify trends reasonably well.

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} [13, 16, 18, 19] which can be computationally prohibitive when p is large. The computational limitation can be perceived through the relatively small number of variables used in simulations and real data analysis in [13, 16, 18, 19]. More recent proposals for selection of interactions allow the user to restrict the search space to interaction candidates [14, 15]. 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 [20, 21] in the linear setting. There are many fewer methods available for estimating non-linear interactions. For example, Radchenko and James (2010) [12] 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.

While working on this paper, we were made aware of the recently proposed pliable lasso [22] which considers the interactions between $\mathbf{X}_{n \times p}$ and another matrix $\mathbf{Z}_{n \times K}$ and takes the

form

$$Y = \beta_0 + \sum_{j=1}^p \beta_j X_j + \sum_{j=1}^K \theta_j Z_j + \sum_{j=1}^p (X_j \circ \mathbf{Z}) \boldsymbol{\alpha}_j + \varepsilon \quad (6)$$

where $\boldsymbol{\alpha}_j$ is a K -dimensional vector. Our proposal is most closely related to this method with \mathbf{Z} being a single column matrix; the key difference being the non-linearity effects of our predictor variables. As pointed out by the authors of the pliable lasso, either their or ours can be seen as a varying coefficient model, i.e., the effect of X varies as a function of the exposure variable E or \mathbf{Z} in equation 6.

The main contributions of this paper are fourfold. 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 over existing methods that only consider linear interactions or additive main effects. Third, we show that our method possesses the oracle property [23], i.e., it performs as well as if the true model were known in advance. Fourth, all of our algorithms are implemented in the **sail** R package hosted on GitHub with extensive documentation (<http://sahirbhatnagar.com/sail/>). 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 [24] idea.

The rest of the paper is organized as follows. Section 3 describes our optimization procedure and some details about the algorithm used to fit the **sail** model for the least squares case. 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 some real data examples and Section 6 discusses some limitations and future directions.

2 Theory

In this section we study the asymptotic behaviour of the **sail** estimator $\widehat{\Theta}$, 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 [24]. The regularity conditions and proofs are given in the Supplementary Material.

2.1 Problem setup

Let $\Theta^* = (\beta_E^*, \theta_1^*, \dots, \theta_p^*, \gamma_1^*, \dots, \gamma_p^*)$ denote the unknown vector of true coefficients in (4), $\mathcal{H} = \{h : \Theta_h^* \neq 0\}$ the unknown sparsity pattern of Θ^* , and $\widehat{\mathcal{H}} = \{h : \widehat{\Theta}_h \neq 0\}$ the estimated **sail** model selector. Here, we assume Θ^* is ordered such that the index $h = 1$ corresponds to β_E^* , the indices $h = 2, \dots, p+1$ correspond to $\theta_1^*, \dots, \theta_p^*$, and the indices $h = p+2, \dots, 2p+1$ corresponds to $\gamma_1^*, \dots, \gamma_p^*$. To simplify the notation, we redefine the penalty terms in (4) and parameter indices, and consider the **sail** estimates $\widehat{\Theta}_n$ given by

$$\widehat{\Theta}_n = \arg \min_{\Theta} Q_n(\Theta) = -L_n(\Theta) + n\lambda_h^\beta |\beta_h| \cdot \mathbf{I}_{\{h=1\}} + n \sum_{h=2}^{p+1} \lambda_h^\theta \|\theta_h\|_2 + n \sum_{h=p+2}^{2p+1} \lambda_h^\gamma |\gamma_h| \quad (7)$$

where $\lambda_h^\beta = \frac{1}{n}\lambda(1-\alpha)w_h^\beta$ for $h = 1$, $\lambda_h^\theta = \frac{1}{n}\lambda(1-\alpha)w_h^\theta$ for $h = 2, \dots, p+1$, and $\lambda_h^\gamma = \frac{1}{n}\lambda\alpha w_h^\gamma$ for $h = p+2, \dots, 2p+1$. The superscripts β, θ, γ on the tuning parameters (λ_h) and adaptive weights (w_h) indicate which parameter (environment, main effect, interaction) they are being applied to. Let $a_n = \max \left\{ \lambda_h^\beta, \lambda_h^\theta, \lambda_h^\gamma : h \in \mathcal{H} \right\}$ and $b_n = \min \left\{ \lambda_h^\beta, \lambda_h^\theta, \lambda_h^\gamma : h \notin \mathcal{H} \right\}$. In words, a_n is a sequence containing the largest tuning parameters such that all causal covariates are included in the model, and b_n is a sequence containing the smallest tuning parameters such that all noise variables are excluded from the model. 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 β_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). *Assume $a_n = O_P(\frac{1}{\sqrt{n}})$. Let $\eta_n = \frac{1}{\sqrt{n}} + a_n$ and $\{\Theta^* + \eta_n \delta : \|\delta\|_2 \leq C\}$ be the ball around Θ^* for $\delta \in \mathbb{R}^d$, where d is the dimension of the design matrix and C is some constant. Under the regularity assumptions, there exists a local minimizer $\hat{\Theta}_n$ of $Q_n(\Theta)$ such that $\|\hat{\Theta}_n - \Theta^*\|_2 = O_P(\eta_n)$.*

The proof is in Appendix A. Lemma (1) states that if the tuning parameters corresponding to the nonzero coefficients converge to 0 at a speed faster than $\frac{1}{\sqrt{n}}$, then there exists a local minimizer of $Q_n(\Theta)$ which is \sqrt{n} -consistent [16, 25]. The following Theorem shows that this estimator is model selection consistent.

Theorem 1 (Model selection consistency). *Assume $\sqrt{nb_n} \rightarrow \infty$ and $\|\hat{\Theta}_n - \Theta^*\|_2 = O_P(\frac{1}{\sqrt{n}})$. Then*

$$\mathbb{P} \left\{ \hat{\mathcal{H}} = \mathcal{H} \right\} \rightarrow 1, \quad \text{as } n \rightarrow \infty \quad (8)$$

The proof is in Appendix B. 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 (8) [26]. Next, we obtain the asymptotic distribution of the **sail** estimator.

Theorem 2 (Asymptotic normality). *Assume that $a_n = O_P(\frac{1}{\sqrt{n}})$ and $\sqrt{nb_n} \rightarrow \infty$. Then, under the regularity conditions,*

$$\sqrt{n} \left(\hat{\Theta} - \Theta^* \right) \xrightarrow{d} Z, \quad (9)$$

where $Z_{\mathcal{H}} \sim N(0, \mathbf{I}^{-1}(\Theta_{\mathcal{H}}^*))$, $Z_{\mathcal{H}^c} = 0$, and $\mathbf{I}(\Theta_{\mathcal{H}}^*)$ is the Fisher information matrix for the causal predictors only.

The proof is in Appendix C. Together, Theorems (1) and (2) establish that if the tuning parameters satisfy the conditions $a_n = O_P(\frac{1}{\sqrt{n}})$ and $\sqrt{n}b_n \rightarrow \infty$, then as the sample size grows large, `sail` has the oracle property [23]. In order for the conditions on the tuning parameters to be satisfied, we follow the strategies outlined for the adaptive lasso [24], the adaptive group lasso [26] and the adaptive elastic-net [27]. That is, we define the adaptive weights as $w_h = \|\hat{\Theta}_h^{ini} + 1/n\|_2^{-\xi}$ for $h = 1, \dots, 2p+1$, where ξ is a positive constant and $\hat{\Theta}_h^{ini}$ is an initial \sqrt{n} -consistent estimate of Θ^* . Here, the $1/n$ is to avoid division by zero.

2.2 Lemma 1 proof

For this proof, we adopt the approaches outlined in [16, 23, 25, 26] and extend it to our situation. Let $\eta_n = \frac{1}{\sqrt{n}} + a_n$ and $\{\Theta^* + \eta_n \delta : \|\delta\|_2 \leq C\}$ be the ball around Θ^* for $\delta = (u_1, \mathbf{u}_2, \dots, \mathbf{u}_{p+1}, v_{p+2}, \dots, v_{2p+1})^\top \in \mathbb{R}^d$, where d is the dimension of the design matrix and C is some constant. The objective function is given by

$$Q_n(\Theta) = -L_n(\Theta) + n\lambda_h^\beta |\beta_h| \cdot \mathbf{I}_{\{h=1\}} + n \sum_{h=2}^{p+1} \lambda_h^\theta \|\theta_h\|_2 + n \sum_{h=p+2}^{2p+1} \lambda_h^\gamma |\gamma_h|$$

Define

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

Let $\mathcal{H}_1 = \{h : \Theta_h^* \neq 0 \text{ for } h = 1\}$, $\mathcal{H}_2 = \{h : \Theta_h^* \neq 0 \text{ for } h = 2, \dots, p+1\}$

and $\mathcal{H}_3 = \{h : \Theta_h^* \neq 0 \text{ for } h = p+2, \dots, 2p+1\}$ where $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \mathcal{H}_3$. Then for δ that satisfies $\|\delta\|_2 = C$, we have

$$\begin{aligned}
D_n(\boldsymbol{\delta}) &= -L_n(\boldsymbol{\Theta}^* + \eta_n \boldsymbol{\delta}) + L_n(\boldsymbol{\Theta}^*) + n\lambda_h^\beta(|\beta_h^* + \eta_n u_h| - |\beta_h^*|) \cdot \mathbf{I}_{\{h=1\}} + \sum_{h=2}^{p+1} \lambda_h^\theta(\|\boldsymbol{\theta}_h^* + \eta_n \mathbf{u}_h\|_2 - \|\boldsymbol{\theta}_h^*\|_2) \\
&\quad + n \sum_{h=p+2}^{2p+1} \lambda_h^\gamma(|\gamma_h^* + \eta_n v_h| - |\gamma_h^*|) \\
&\stackrel{(a)}{\geq} -L_n(\boldsymbol{\Theta}^* + \eta_n \boldsymbol{\delta}) + L_n(\boldsymbol{\Theta}^*) + \lambda_h^\beta(|\beta_h^* + \eta_n u_h| - |\beta_h^*|) \cdot \mathbf{I}_{\{h \in \mathcal{H}_1\}} + n \sum_{h \in \mathcal{H}_2} \lambda_h^\theta(\|\boldsymbol{\theta}_h^* + \eta_n \mathbf{u}_h\|_2 - \|\boldsymbol{\theta}_h^*\|_2) \\
&\quad + n \sum_{h \in \mathcal{H}_3} \lambda_h^\gamma(|\gamma_h^* + \eta_n v_h| - |\gamma_h^*|) \\
&\stackrel{(b)}{\geq} -L_n(\boldsymbol{\Theta}^* + \eta_n \boldsymbol{\delta}) + L_n(\boldsymbol{\Theta}^*) - n\eta_n \lambda_h^\beta |u_h| \cdot \mathbf{I}_{\{h \in \mathcal{H}_1\}} - n\eta_n \sum_{h \in \mathcal{H}_2} \lambda_h^\theta \|\mathbf{u}_h\|_2 - n\eta_n \sum_{h \in \mathcal{H}_3} \lambda_h^\gamma |v_h| \\
&\stackrel{(c)}{\geq} -L_n(\boldsymbol{\Theta}^* + \eta_n \boldsymbol{\delta}) + L_n(\boldsymbol{\Theta}^*) - n\eta_n^2 |u_h| \cdot \mathbf{I}_{\{h \in \mathcal{H}_1\}} - n\eta_n^2 \sum_{h \in \mathcal{H}_2} \|\mathbf{u}_h\|_2 - n\eta_n^2 \sum_{h \in \mathcal{H}_3} |v_h| \\
&\stackrel{(d)}{\geq} -L_n(\boldsymbol{\Theta}^* + \eta_n \boldsymbol{\delta}) + L_n(\boldsymbol{\Theta}^*) - n\eta_n^2(|\mathcal{H}_1| + |\mathcal{H}_2| + |\mathcal{H}_3|)C \\
&\stackrel{(e)}{=} -[\nabla L_n(\boldsymbol{\Theta}^*)]^\top (\eta_n \boldsymbol{\delta}) - \frac{1}{2}(\eta_n \boldsymbol{\delta})^\top [\nabla^2 L_n(\boldsymbol{\Theta}^*)](\eta_n \boldsymbol{\delta})(1 + o_p(1)) \\
&\quad - n\eta_n^2(|\mathcal{H}_1| + |\mathcal{H}_2| + |\mathcal{H}_3|)C
\end{aligned} \tag{10}$$

Inequality (a) is by the fact that $\sum_{h \notin \mathcal{H}_1} |\boldsymbol{\Theta}_h^*| = 0$ and $\sum_{h \notin \mathcal{H}_2} |\boldsymbol{\Theta}_h^*| = 0$. Inequality (b) is due to the reverse triangle inequality $|a| - |b| \geq -|a - b|$. Inequality (c) is by $\lambda_h^\beta \leq \eta_n$, $\lambda_h^\theta \leq \eta_n$ and $\lambda_h^\gamma \leq \eta_n$. Inequality (d) is by the Cauchy-Schwarz inequality:

$$\sum_{j \in \mathcal{H}} |u_j| = |\mathbf{1}_{\mathcal{H}}^\top \mathbf{u}_{\mathcal{H}}| \leq \|\mathbf{1}_{\mathcal{H}}\|_2 \|\mathbf{u}_{\mathcal{H}}\|_2 = \sqrt{|\mathcal{H}|} \|\mathbf{u}_{\mathcal{H}}\|_2 \leq |\mathcal{H}|C.$$

Equality (e) is by the standard argument on the Taylor expansion of the loss function:

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

We split (10) into three parts:

$$\begin{aligned} A_1 &= -[\nabla L_n(\boldsymbol{\Theta}^*)]^\top (\eta_n \boldsymbol{\delta}) \\ A_2 &= -\frac{1}{2} (\eta_n \boldsymbol{\delta})^\top [\nabla^2 L_n(\boldsymbol{\Theta}^*)] (\eta_n \boldsymbol{\delta}) (1 + o_p(1)) \\ A_3 &= -n\eta_n^2 (|\mathcal{H}_1| + |\mathcal{H}_2| + |\mathcal{H}_3|) C \end{aligned}$$

$$\begin{aligned} A_1 &= -\eta_n [\nabla L_n(\boldsymbol{\Theta}^*)]^\top \boldsymbol{\delta} \\ &= -\sqrt{n}\eta_n \left(\frac{1}{\sqrt{n}} \nabla L_n(\boldsymbol{\Theta}^*) \right)^\top \boldsymbol{\delta} \\ &= -\sqrt{n}\eta_n \left(\sqrt{n} \frac{1}{n} \sum_{i=1}^n \nabla \log f(\mathbf{V}_i, \boldsymbol{\Theta})|_{\boldsymbol{\Theta}=\boldsymbol{\Theta}^*} \right)^\top \boldsymbol{\delta} \\ &= -\sqrt{n}\eta_n \left(\sqrt{n} \left[\frac{1}{n} \sum_{i=1}^n \nabla \log f(\mathbf{V}_i, \boldsymbol{\Theta})|_{\boldsymbol{\Theta}=\boldsymbol{\Theta}^*} - \mathbf{0} \right] \right)^\top \boldsymbol{\delta} \\ &= -\sqrt{n}\eta_n \left(\sqrt{n} \left[\frac{1}{n} \sum_{i=1}^n \nabla \log f(\mathbf{V}_i, \boldsymbol{\Theta})|_{\boldsymbol{\Theta}=\boldsymbol{\Theta}^*} - E_{\boldsymbol{\Theta}^*} \nabla L(\boldsymbol{\Theta}^*) \right] \right)^\top \boldsymbol{\delta} \\ &= -\sqrt{n}\eta_n O_P(1) \boldsymbol{\delta} \\ &= -O_P(\sqrt{n}\eta_n \|\boldsymbol{\delta}\|_2) \end{aligned} \tag{11}$$

$$\begin{aligned} A_2 &= \frac{1}{2} n\eta_n^2 \left\{ \boldsymbol{\delta}^\top \left[-\frac{1}{n} \nabla^2 L_n(\boldsymbol{\Theta}^*) \right] \boldsymbol{\delta} \right\} (1 + o_p(1)) \\ &= \frac{1}{2} n\eta_n^2 \left\{ \boldsymbol{\delta}^\top [I(\boldsymbol{\Theta}^*)] \boldsymbol{\delta} \right\} (1 + o_p(1)) \text{ by the weak law of large numbers.} \\ &= O(n\eta_n^2 \|\boldsymbol{\delta}\|_2^2) \end{aligned} \tag{12}$$

Combining (11) and (12) with (10) gives:

$$\begin{aligned} D_n(\boldsymbol{\delta}) &\geq A_1 + A_2 + A_3 \\ &= -O_P(\sqrt{n}\eta_n \|\boldsymbol{\delta}\|_2) + O(n\eta_n^2 \|\boldsymbol{\delta}\|_2^2) - n\eta_n^2 (|\mathcal{H}_1| + |\mathcal{H}_2| + |\mathcal{H}_3|) C \end{aligned}$$

We can conclude that for a large enough constant $C = \|\boldsymbol{\delta}\|_2$, the positive term $O(n\eta_n^2\|\boldsymbol{\delta}\|_2^2)$ dominates all the others. Note that this is a positive term since $I(\boldsymbol{\Theta})$ is positive definite at $\boldsymbol{\Theta} = \boldsymbol{\Theta}^*$ by regularity condition (C2). Therefore, for each $\varepsilon > 0$, there exists a large enough constant C such that, for large enough n

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

This implies with probability at least $1 - \varepsilon$ that the empirical likelihood Q_n has a local minimizer in the ball $\{\boldsymbol{\Theta}^* + \eta_n \boldsymbol{\delta} : \|\boldsymbol{\delta}\|_2 \leq C\}$ (since Q_n is bounded and $\{\boldsymbol{\Theta}^* + \alpha_n \boldsymbol{\delta} : \|\boldsymbol{\delta}\|_2 \leq C\}$ is closed). In other words, there exists a solution $\hat{\boldsymbol{\Theta}}_n$ such that $\|\hat{\boldsymbol{\Theta}}_n - \boldsymbol{\Theta}^*\| \leq \eta_n \|\boldsymbol{\delta}\|_2 \leq \eta_n C = O_P(\eta_n) = O_P(n^{-1/2} + a_n)$. Hence, $\|\hat{\boldsymbol{\Theta}}_n - \boldsymbol{\Theta}^*\| = O_P(\eta_n)$. \square

2.3 Theorem 1 proof

We first recall the definitions of the active sets from Lemma 1. Let $\mathcal{H}_1 = \{h : \boldsymbol{\Theta}_h^* \neq 0 \text{ for } h = 1\}$, $\mathcal{H}_2 = \{h : \boldsymbol{\Theta}_h^* \neq 0 \text{ for } h = 2, \dots, p+1\}$ and $\mathcal{H}_3 = \{h : \boldsymbol{\Theta}_h^* \neq 0 \text{ for } h = p+2, \dots, 2p+1\}$ where $\mathcal{H} = \mathcal{H}_1 \cup \mathcal{H}_2 \cup \mathcal{H}_3$. We first consider consistency for the main effects:

$$\begin{aligned} P(\hat{\boldsymbol{\Theta}}_{\mathcal{H}_1^c} &= 0) \rightarrow 1 \\ P(\hat{\boldsymbol{\Theta}}_{\mathcal{H}_2^c} &= 0) \rightarrow 1 \end{aligned}$$

Following [16, 23], it is sufficient to show for any $h \in \mathcal{H}_1^c \cup \mathcal{H}_2^c$

$$\frac{\partial Q_n(\hat{\boldsymbol{\Theta}}_n)}{\partial \boldsymbol{\Theta}_h} < 0 \quad \text{for } -\varepsilon_n < \hat{\boldsymbol{\Theta}}_h < 0 \quad (13)$$

$$\frac{\partial Q_n(\hat{\boldsymbol{\Theta}}_n)}{\partial \boldsymbol{\Theta}_h} > 0 \quad \text{for } 0 < \hat{\boldsymbol{\Theta}}_h < \varepsilon_n \quad (14)$$

with probability tending to 1 where $\varepsilon_n = Cn^{-1/2}$ and $C > 0$ is any constant. To show (14), notice that

for $h = 1$

$$\begin{aligned} \frac{\partial Q_n(\hat{\Theta}_n)}{\partial \Theta_h} &= -\frac{\partial L_n(\hat{\Theta}_n)}{\partial \Theta_h} + n\lambda_h^\beta \text{sgn}(\hat{\Theta}_h) \\ &= -\frac{\partial L_n(\Theta^*)}{\partial \Theta_h} - \sum_{k=2}^{2p+1} \frac{\partial^2 L_n(\Theta^*)}{\partial \Theta_h \partial \Theta_k} (\hat{\Theta}_k - \Theta_k^*) \\ &\quad - \sum_{k=2}^{2p+1} \sum_{l=2}^{2p+1} (2p+1) \frac{\partial^3 L_n(\tilde{\Theta})}{\partial \Theta_h \partial \Theta_k \partial \Theta_l} (\hat{\Theta}_k - \Theta_k^*) (\hat{\Theta}_l - \Theta_l^*) + n\lambda_h^\beta \text{sgn}(\hat{\Theta}_h) \end{aligned}$$

for $h = 2, \dots, p+1$

$$\begin{aligned} \frac{\partial Q_n(\hat{\Theta}_n)}{\partial \Theta_h} &= -\frac{\partial L_n(\hat{\Theta}_n)}{\partial \Theta_h} + n\lambda_h^\theta \frac{\hat{\Theta}_h}{\|\hat{\Theta}_h\|_2} \\ &= -\frac{\partial L_n(\Theta^*)}{\partial \Theta_h} - \sum_{k=2}^{2p+1} \frac{\partial^2 L_n(\Theta^*)}{\partial \Theta_h \partial \Theta_k} (\hat{\Theta}_k - \Theta_k^*) \\ &\quad - \sum_{k=2}^{2p+1} \sum_{l=2}^{2p+1} (2p+1) \frac{\partial^3 L_n(\tilde{\Theta})}{\partial \Theta_h \partial \Theta_k \partial \Theta_l} (\hat{\Theta}_k - \Theta_k^*) (\hat{\Theta}_l - \Theta_l^*) + n\lambda_h^\theta \frac{\hat{\Theta}_h}{\|\hat{\Theta}_h\|_2} \end{aligned}$$

where $\tilde{\Theta}$ lies between $\hat{\Theta}_n$ and Θ^* . In both cases above, the first term is

$$\frac{\partial L_n(\Theta^*)}{\partial \Theta_h} = \sqrt{n} \sqrt{n} \frac{1}{n} \frac{\partial L_n(\Theta^*)}{\partial \Theta_h} = \sqrt{n} O_P(1) = O_P(n^{1/2}).$$

Since

$$\frac{1}{n} \frac{\partial^2 L_n(\Theta^*)}{\partial \Theta_h \partial \Theta_k} = E \left\{ \frac{\partial^2 \log f(\Theta^*)}{\partial \Theta_h \partial \Theta_k} \right\} + o_P(1),$$

then

$$\frac{\partial^2 L_n(\Theta^*)}{\partial \Theta_h \partial \Theta_k} = nE \left\{ \frac{\partial^2 \log f(\Theta^*)}{\partial \Theta_h \partial \Theta_k} \right\} + o_P(n),$$

Then the second term

$$\begin{aligned} \sum_{k=2}^{2p+1} \frac{\partial^2 L_n(\Theta^*)}{\partial \Theta_h \partial \Theta_k} (\hat{\Theta}_k - \Theta_k^*) &= \left[nE \left\{ \frac{\partial^2 \log f(\Theta_*)}{\partial \Theta_h \partial \Theta_k} \right\} + o_P(n) \right] O_p(n^{-1/2}) \\ &= O_P(n^{1/2}) + o_P(n^{1/2}) \end{aligned}$$

The third term is

$$- \sum_{k=2}^{2p+1} \sum_{l=2}^{2p+1} (2p+1) \frac{\partial^3 L_n(\tilde{\Theta})}{\partial \Theta_h \partial \Theta_k \partial \Theta_l} (\hat{\Theta}_k - \Theta_k^*) (\hat{\Theta}_l - \Theta_l^*) = O_p(n^{-1/2}) O_p(n^{-1/2}) = O_p(n^{-1})$$

By the regularity conditions and the condition $\|\hat{\Theta}_n - \Theta^*\| = O_p\left(\frac{1}{\sqrt{n}}\right)$,

for $h = 1$

$$\frac{\partial Q_n(\hat{\Theta}_n)}{\partial \Theta_h} = \sqrt{n} \left\{ O_p(1) + \sqrt{n} \lambda_j^\beta \text{sgn}(\hat{\beta}_j) \right\}. \quad (15)$$

for $h = 2, \dots, p+1$

$$\frac{\partial Q_n(\hat{\Theta}_n)}{\partial \Theta_h} = \sqrt{n} \left\{ O_p(1) + \sqrt{n} \lambda_h^\theta \frac{\hat{\Theta}_h}{\|\hat{\Theta}_h\|_2} \right\}. \quad (16)$$

As $\sqrt{n} \lambda_h^\beta \rightarrow \infty$ for $h \in \mathcal{H}_1^c$ from the assumption, the sign of $\frac{\partial Q_n(\hat{\Theta}_n)}{\partial \Theta_h}$ in (15) is dominated by $\text{sgn}(\hat{\Theta}_h)$. Similarly, as $\sqrt{n} \lambda_h^\theta \rightarrow \infty$ for $h \in \mathcal{H}_2^c$, the sign of $\frac{\partial Q_n(\hat{\Theta}_n)}{\partial \Theta_h}$ in (16) is dominated by the norm of $\hat{\Theta}_h$.

Therefore,

$$P \left[\frac{\partial Q_n(\hat{\Theta}_n)}{\partial \Theta_h} > 0 \quad \text{for } 0 < \hat{\Theta}_h < \varepsilon_n \right] \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

(13) can be shown in the same way.

Due to the strong heredity property and the proof above for the consistency of the main effects, it follows that $P\left(\widehat{\Theta}_{\mathcal{H}_3^c} = 0\right) \rightarrow 1$. \square

2.4 Theorem 2 proof

By Lemma (1) and Theorem (1), there exists a $\widehat{\Theta}_{\mathcal{H}}$ that is a \sqrt{n} -consistent local minimizer of $Q(\Theta_{\mathcal{H}})$, and that satisfies (with probability tending to 1):

$$\left. \frac{\partial Q_n(\Theta_{\mathcal{H}})}{\partial \Theta_h} \right|_{\Theta_{\mathcal{H}} = \widehat{\Theta}_{\mathcal{H}}} = 0, \quad \forall h \in \mathcal{H} \quad (17)$$

where

$$\begin{aligned} Q_n(\Theta_{\mathcal{H}}) &= -L_n(\Theta_{\mathcal{H}}) + n \underbrace{\lambda_h^\beta |\beta_h| \cdot \mathbf{I}_{\{h \in \mathcal{H}_1\}} + n \sum_{h \in \mathcal{H}_2} \lambda_h^\theta \|\theta_h\|_2 + n \sum_{h \in \mathcal{H}_3} \lambda_h^\gamma |\gamma_h|}_{\equiv nP(\Theta_{\mathcal{H}})} \\ &= -L_n(\Theta_{\mathcal{H}}) + nP(\Theta_{\mathcal{H}}) \end{aligned} \quad (18)$$

From (17) and (18) we have

$$\nabla_{\mathcal{H}} Q_n(\widehat{\Theta}_{\mathcal{H}}) = -\nabla_{\mathcal{H}} L_n(\widehat{\Theta}_{\mathcal{H}}) + n \nabla_{\mathcal{H}} P(\widehat{\Theta}_{\mathcal{H}}) = \mathbf{0}, \quad (19)$$

with probability tending to 1.

We then expand $-\nabla_{\mathcal{H}} L_n(\Theta_{\mathcal{H}})$ at $\Theta_{\mathcal{H}} = \Theta_{\mathcal{H}}^*$ in (19):

$$\begin{aligned} -\nabla_{\mathcal{H}} L_n(\widehat{\Theta}_{\mathcal{H}}) &= -\nabla_{\mathcal{H}} L_n(\Theta_{\mathcal{H}}^*) - [\nabla_{\mathcal{H}}^2 L_n(\Theta_{\mathcal{H}}^*) + o_p(1)] (\widehat{\Theta}_{\mathcal{H}} - \Theta_{\mathcal{H}}^*) \\ &= \sqrt{n} \left[-\frac{1}{\sqrt{n}} \nabla_{\mathcal{H}} L_n(\Theta_{\mathcal{H}}^*) + \left(-\frac{1}{n} \nabla_{\mathcal{H}}^2 L_n(\Theta_{\mathcal{H}}^*) - o_p(1) \right) \sqrt{n} (\widehat{\Theta}_{\mathcal{H}} - \Theta_{\mathcal{H}}^*) \right] \\ &= \sqrt{n} \left[-\frac{1}{\sqrt{n}} \nabla_{\mathcal{H}} L_n(\Theta_{\mathcal{H}}^*) + I(\Theta_{\mathcal{H}}^*) \sqrt{n} (\widehat{\Theta}_{\mathcal{H}} - \Theta_{\mathcal{H}}^*) + o_p(1) \right] \end{aligned}$$

The third line follows by

$$\frac{1}{n} \nabla_{\mathcal{H}}^2 L_n(\boldsymbol{\Theta}_{\mathcal{H}}^*) = E \left\{ \nabla_{\mathcal{H}}^2 L(\boldsymbol{\Theta}_{\mathcal{H}}^*) \right\} + o_P(1) = -\mathbf{I}(\boldsymbol{\Theta}_{\mathcal{H}}^*) + o_P(1)$$

We also expand $n \nabla_{\mathcal{H}} P(\boldsymbol{\Theta}_{\mathcal{H}})$ at $\boldsymbol{\Theta}_{\mathcal{H}} = \boldsymbol{\Theta}_{\mathcal{H}}^*$ in (19):

$$\begin{aligned} n \nabla_{\mathcal{H}} P(\hat{\boldsymbol{\Theta}}_{\mathcal{H}}) &= n \left\{ \begin{bmatrix} \lambda_h^\beta \text{sgn}(\beta_h) \\ \lambda_h^\theta \frac{\boldsymbol{\theta}_h}{\|\boldsymbol{\theta}_h\|_2} \\ \lambda_h^\gamma \text{sgn}(\gamma_h) \end{bmatrix}_{h \in \mathcal{H}} + o_p(1) (\hat{\boldsymbol{\Theta}}_{\mathcal{H}} - \boldsymbol{\Theta}_{\mathcal{H}}^*) \right\} \\ &= \sqrt{n} o_p(1) \end{aligned}$$

because $\sqrt{n} a_n = o(1)$ and $\|\hat{\boldsymbol{\Theta}}_{\mathcal{H}} - \boldsymbol{\Theta}_{\mathcal{H}}^*\| = O_p(n^{-1/2})$.

Thus,

$$0 = \sqrt{n} \left[-\frac{1}{\sqrt{n}} \nabla_{\mathcal{H}} L_n(\boldsymbol{\Theta}_{\mathcal{H}}^*) + \mathbf{I}(\boldsymbol{\Theta}_{\mathcal{H}}^*) \sqrt{n} (\hat{\boldsymbol{\Theta}}_{\mathcal{H}} - \boldsymbol{\Theta}_{\mathcal{H}}^*) + o_p(1) \right].$$

It follows

$$\sqrt{n} (\hat{\boldsymbol{\Theta}}_{\mathcal{H}} - \boldsymbol{\Theta}_{\mathcal{H}}^*) = \mathbf{I}(\boldsymbol{\Theta}_{\mathcal{H}}^*)^{-1} \sqrt{n} \frac{1}{n} \sum_{i=1}^n \nabla_{\mathcal{H}} \log f(\mathbf{V}_i, \boldsymbol{\Theta}_{\mathcal{H}}) + o_p(1).$$

Therefore, by the central limit theorem,

$$\sqrt{n} (\hat{\boldsymbol{\Theta}}_{\mathcal{H}} - \boldsymbol{\Theta}_{\mathcal{H}}^*) \rightarrow_d N(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\Theta}_{\mathcal{H}}^*)).$$

□

3 Algorithm and Computational Details

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 α and minimize the objective function over a decreasing sequence of λ values ($\lambda_{max} > \dots > \lambda_{min}$). We use the subgradient equations to determine the maximal value λ_{max} such that all estimates are zero. Due to the heredity principle, this reduces to finding the largest λ such that all main effects ($\beta_E, \theta_1, \dots, \theta_p$) are zero. Following Friedman et al. [28], we construct a λ -sequence of 100 values decreasing from λ_{max} to $0.001\lambda_{max}$ on the log scale, and use the warm start strategy where the solution for λ_ℓ is used as a starting value for $\lambda_{\ell+1}$.

3.1 Blockwise coordinate descent for least-squares loss

The strong heredity `sail` model with least-squares loss has the form

$$\hat{Y} = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \Psi_j \theta_j + \beta_E X_E + \sum_{j=1}^p \gamma_j \beta_E (X_E \circ \Psi_j) \theta_j \quad (20)$$

and the objective function is given by

$$Q(\Theta) = \frac{1}{2n} \|Y - \hat{Y}\|_2^2 + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (21)$$

Solving (21) in a blockwise manner allows us to leverage computationally fast algorithms for ℓ_1 and ℓ_2 norm penalized regression. We show in Supplemental Section A that by careful construction of pseudo responses and pseudo design matrices, existing efficient algorithms can be used to estimate the parameters. Indeed, the objective function simplifies to a modified lasso problem when holding all θ_j fixed, and a modified group lasso problem when holding β_E and all γ_j fixed.

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

$$\frac{\partial Q}{\partial \beta_0} = \frac{1}{n} \left(Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \Psi_j \theta_j - \beta_E X_E - \sum_{j=1}^p \gamma_j \beta_E (X_E \circ \Psi_j) \theta_j \right)^\top \mathbf{1} = 0 \quad (22)$$

$$\frac{\partial Q}{\partial \beta_E} = -\frac{1}{n} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^\top R + \lambda(1 - \alpha) w_E s_1 = 0 \quad (23)$$

$$\frac{\partial Q}{\partial \theta_j} = -\frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R + \lambda(1 - \alpha) w_j s_2 = \mathbf{0} \quad (24)$$

$$\frac{\partial Q}{\partial \gamma_j} = -\frac{1}{n} (\beta_E (X_E \circ \Psi_j) \theta_j)^\top R + \lambda \alpha w_{jE} s_3 = 0 \quad (25)$$

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:

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

and s_3 is in the subgradient of the 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

$$R_{(-j)} = Y - \beta_0 \cdot \mathbf{1} - \sum_{\ell \neq j} \Psi_\ell \theta_\ell - \beta_E X_E - \sum_{\ell \neq j} \gamma_\ell \beta_E (X_E \circ \Psi_\ell) \theta_\ell$$

the partial residual without X_E as

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

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

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

From the subgradient equations (22)–(25) we see that

$$\hat{\beta}_0 = \left(Y - \sum_{j=1}^p \Psi_j \hat{\boldsymbol{\theta}}_j - \hat{\beta}_E X_E - \sum_{j=1}^p \hat{\gamma}_j \hat{\beta}_E (X_E \circ \Psi_j) \hat{\boldsymbol{\theta}}_j \right)^\top \mathbf{1} \quad (26)$$

$$\hat{\beta}_E = S \left(\frac{1}{n \cdot w_E} \left(X_E + \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) \hat{\boldsymbol{\theta}}_j \right)^\top R_{(-E)}, \lambda(1 - \alpha) \right) \quad (27)$$

$$\lambda(1 - \alpha) w_j \frac{\boldsymbol{\theta}_j}{\|\boldsymbol{\theta}_j\|_2} = \frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-j)} \quad (28)$$

$$\hat{\gamma}_j = S \left(\frac{1}{n \cdot w_{jE}} (\beta_E (X_E \circ \Psi_j) \boldsymbol{\theta}_j)^\top R_{(-jE)}, \lambda \alpha \right) \quad (29)$$

where $S(x, t) = \text{sign}(x)(|x| - t)$ is the soft-thresholding operator. We see from (26) and (27) that there are closed form solutions for the intercept and β_E . From (29), each γ_j also has a closed form solution and can be solved efficiently for $j = 1, \dots, p$ using a coordinate descent procedure [28]. Since there is no closed form solution for β_j , we use a quadratic majorization technique [29] to solve (28). Furthermore, we update each $\boldsymbol{\theta}_j$ in a coordinate wise fashion and leverage this to implement further computational speedups which are detailed in Supplemental Section A.2. From these estimates, we compute the interaction effects using the reparametrizations presented in Table 1, e.g., $\hat{\boldsymbol{\tau}}_j = \hat{\gamma}_j \hat{\beta}_E \hat{\boldsymbol{\theta}}_j$, $j = 1, \dots, p$ for the strong heredity **sail** model. We provide an overview of the computations in Algorithm 1. A more detailed version of this algorithm is given in Section A.1 of the Appendix.

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

For a decreasing sequence $\lambda = \lambda_{max}, \dots, \lambda_{min}$ and fixed α :

1. Initialize $\beta_0^{(0)}, \beta_E^{(0)}, \boldsymbol{\theta}_j^{(0)}, \gamma_j^{(0)}$ for $j = 1, \dots, p$ and set iteration counter $k \leftarrow 0$.
2. Repeat the following until convergence:
 - (a) update $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$
 - i. Compute the pseudo design: $\tilde{X}_j \leftarrow \beta_E^{(k)}(X_E \circ \boldsymbol{\Psi}_j)\boldsymbol{\theta}_j^{(k)}$ for $j = 1, \dots, p$
 - ii. Compute the pseudo response \tilde{Y} by removing the contribution of every term not involving $\boldsymbol{\gamma}$ from Y
 - iii. Solve:

$$\boldsymbol{\gamma}^{(k)(new)} \leftarrow \arg \min_{\boldsymbol{\gamma}} \frac{1}{2n} \left\| \tilde{Y} - \sum_j \gamma_j \tilde{X}_j \right\|_2^2 + \lambda \alpha \sum_j w_{jE} |\gamma_j| \quad (30)$$

- iv. Set $\boldsymbol{\gamma}^{(k)} = \boldsymbol{\gamma}^{(k)(new)}$
- (b) update $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p)$
 - for $j = 1, \dots, p$
 - i. Compute the pseudo design: $\tilde{X}_j \leftarrow \boldsymbol{\Psi}_j + \gamma_j^{(k)} \beta_E^{(k)}(X_E \circ \boldsymbol{\Psi}_j)$
 - ii. Compute the pseudo response (\tilde{Y}) by removing the contribution of every term not involving $\boldsymbol{\theta}_j$ from Y
 - iii. Solve:

$$\boldsymbol{\theta}_j^{(k)(new)} \leftarrow \arg \min_{\boldsymbol{\theta}_j} \frac{1}{2n} \left\| \tilde{Y} - \tilde{X}_j \boldsymbol{\theta}_j \right\|_2^2 + \lambda(1 - \alpha) w_j \|\boldsymbol{\theta}_j\|_2 \quad (31)$$

- iv. Set $\boldsymbol{\theta}_j^{(k)} \leftarrow \boldsymbol{\theta}_j^{(k)(new)}$
- (c) update β_E
 - i. Compute the pseudo design: $\tilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \tilde{\boldsymbol{\Psi}}_j \boldsymbol{\theta}_j^{(k)}$
 - ii. Compute the pseudo response (\tilde{Y}) by removing the contribution of every term not involving β_E from Y
 - iii. Soft-threshold update ($S(x, t) = \text{sign}(x)(|x| - t)_+$):

$$\beta_E^{(k)(new)} \leftarrow S \left(\frac{1}{n \cdot w_E} \tilde{X}_E^\top \tilde{Y}, \lambda(1 - \alpha) \right) \quad (32)$$

- iv. Set $\beta_E^{(k+1)} \leftarrow \beta_E^{(k)(new)}$, $k \leftarrow k + 1$
-

3.2 Maximum penalty parameter (Lambda max)

The subgradient equations (23)–(25) can be used to determine the largest value of λ such that all coefficients are 0. From the subgradient Equation (23), we see that $\beta_E = 0$ is a solution if

$$\frac{1}{w_E} \left| \frac{1}{n} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^\top R_{(-E)} \right| \leq \lambda(1 - \alpha) \quad (33)$$

From the subgradient Equation (24), we see that $\theta_j = \mathbf{0}$ is a solution if

$$\frac{1}{w_j} \left\| \frac{1}{n} (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-j)} \right\|_2 \leq \lambda(1 - \alpha) \quad (34)$$

From the subgradient Equation (25), we see that $\gamma_j = 0$ is a solution if

$$\frac{1}{w_{jE}} \left| \frac{1}{n} (\beta_E (X_E \circ \Psi_j) \theta_j)^\top R_{(-jE)} \right| \leq \lambda \alpha \quad (35)$$

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 λ for which the entire parameter vector (excluding the intercept) is $\mathbf{0}$ is:

$$\lambda_{max} = \frac{1}{n(1 - \alpha)} \max \left\{ \frac{1}{w_E} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \theta_j \right)^\top R_{(-E)}, \max_j \frac{1}{w_j} \left\| (\Psi_j + \gamma_j \beta_E (X_E \circ \Psi_j))^\top R_{(-j)} \right\|_2 \right\} \quad (36)$$

which reduces to

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

3.3 Weak Heredity

Our method can be easily adapted to enforce the weak heredity property:

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

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 (2) as $\boldsymbol{\alpha}_j = \gamma_j(\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j)$, where $\mathbf{1}_{m_j}$ is a vector of ones with dimension m_j (i.e. the length of $\boldsymbol{\theta}_j$). We defer the algorithm details for fitting the **sail** model with weak heredity in Section A.3 of the Appendix, as it is very similar to Algorithm 1 for the strong heredity **sail** model.

3.4 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 [24]

We describe the adaptive **sail** in Algorithm 2. This is a general procedure that can be applied to the weak and strong heredity settings, as well as both least squares and logistic loss functions. We provide this capability in the **sail** package using the `penalty.factor` argument and provide an example in Section C.6 of the Appendix.

Algorithm 2 Adaptive **sail** algorithm

1. For a decreasing sequence $\lambda = \lambda_{max}, \dots, \lambda_{min}$ and fixed α 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
$$w_E = \left(\left| \widehat{\beta}_E^{[opt]} \right| \right)^{-1}, w_j = \left(\left\| \widehat{\theta}_j^{[opt]} \right\|_2 \right)^{-1}, w_{jE} = \left(\left\| \widehat{\tau}_j^{[opt]} \right\|_2 \right)^{-1} \text{ for } j = 1, \dots, p$$
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 λ

3.5 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 \mathbf{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 provide such an example in the **sail** package showcase in Section C.7 of the Appendix.

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) allows at least one of linear effects, non-linear effects or interaction effects, and 3) has a software implementation in R. The selected methods can be grouped into three categories:

1. Linear main effects: `lasso` [30], `adaptive lasso` [24]
2. Linear interactions: `lassoBT` [21], `GLinternet` [14]
3. Non-linear main effects: `HierBasis` [31], `SPAM` [32], `gamsel` [33]

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 (\mathbf{X}, X_E) as the data matrix, 100 for the number of tuning parameters to fit, and used the default values otherwise¹. `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 nonzero, 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), `sail weak` which has the weak heredity property and `linear`

¹R code for each method available at https://github.com/sahirbhatnagar/sail/blob/master/my_sims/method_functions.R

`sail` as described in Section 3.5. For each function f_j , we use a B-spline basis matrix with `degree=5` implemented in the `bs` function in R [34]. 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 [35, 36]. We extend this framework to include interaction effects as well. The covariates are simulated as follows. First, we generate z_1, \dots, z_p, u, v independently from a standard normal distribution truncated to the interval $[0,1]$ for $i = 1, \dots, n$. Then we set $x_j = (z_j + t \cdot u)/(1 + t)$ for $j = 1, \dots, 4$ and $x_j = (z_j + t \cdot v)/(1 + t)$ for $j = 5, \dots, p$, where the parameter t controls the amount of correlation among predictors. The first four variables are nonzero (i.e. active in the response), while the rest of the variables are zero (i.e. are noise variables). 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. We consider the case when $p = 1000$ and $t = 0$. 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. Hierarchy simulation

Scenario (a) Truth obeys strong hierarchy. In this situation, the true model for

Y contains main effect terms for all covariates involved in interactions.

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

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

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

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

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

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.

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

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.

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

The true component functions are the same as in [35, 36] 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 ε 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 test set MSE results for each of the five simulation scenarios are shown in Figure ??, while Figure ?? shows the mean true positive rate (TPR) vs. the mean false positive rate (FPR) ± 1 standard deviation (SD).

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 hierarchy (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, our method is unable to capture any of the truly associated variables when only interactions are present (scenario 1c). However, the other methods also fail to capture any signal, with the exception of **GLinternet** which has a high TPR and FPR. When only linear

effects and interactions are present (scenario 2), we see that `linear sail` has a high TPR and low FPR as compared to the other linear interaction methods (`lassoBT` and `GLinternet`) though the test set MSE is not as good. The `lasso` and `adaptive lasso` have good test set MSE performance but poor sensitivity. Additional results are available in Section B of the Appendix. Specifically, in Figure ?? we plot the mean MSE against the mean number of active variables ± 1 standard deviation (SD). Figures ?? and ?? show the true positive and false positive rates, respectively. Figure ?? shows the number of active variables.

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 1 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. [12], we then identified the 25th, 50th, and 75th best simulations and plotted, in Figures 2 and 3, 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 $X_E \cdot f_4(X_4)$

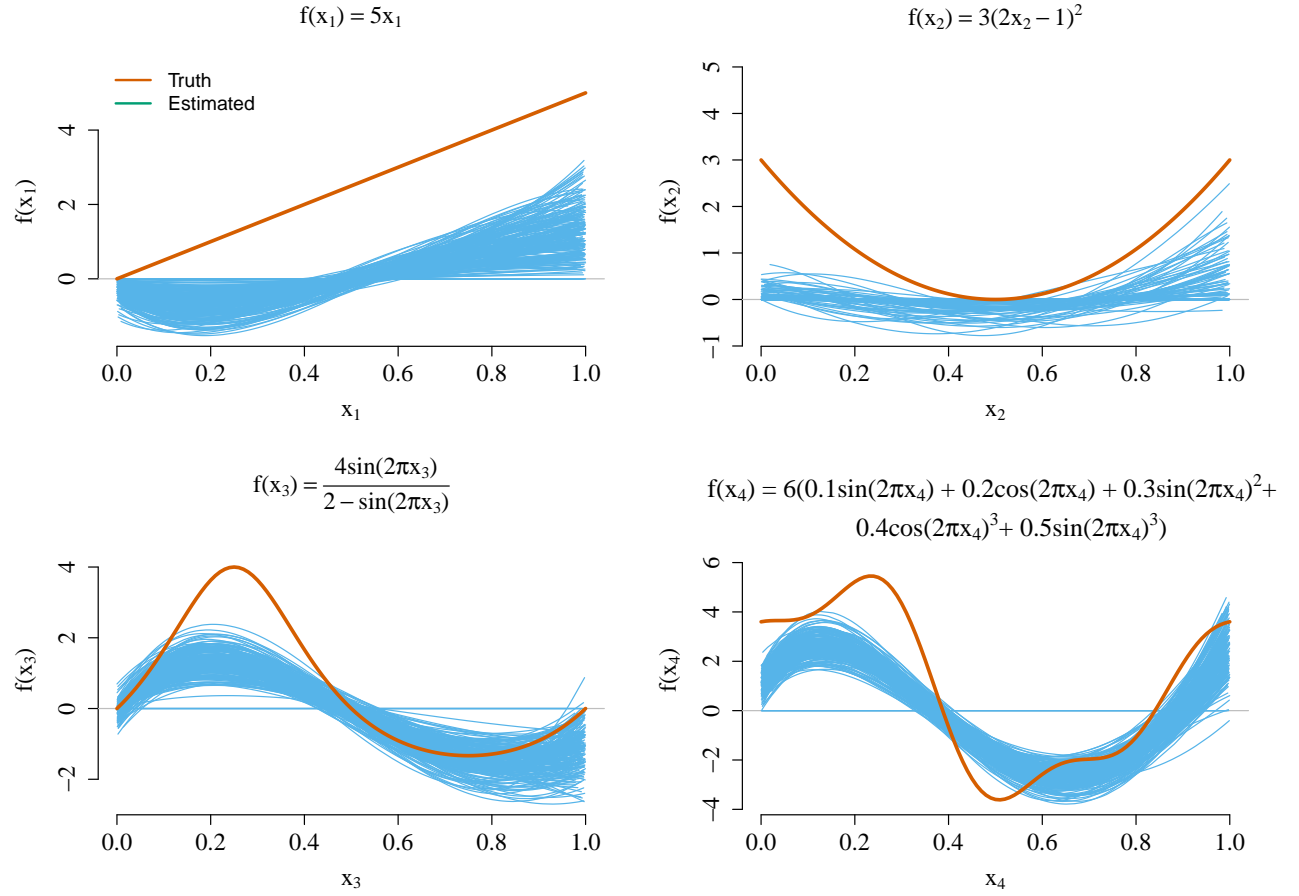


Figure 1: True and estimated main effect component functions for scenario 1a). The estimated curves represent the results from each one of the 200 simulations conducted.

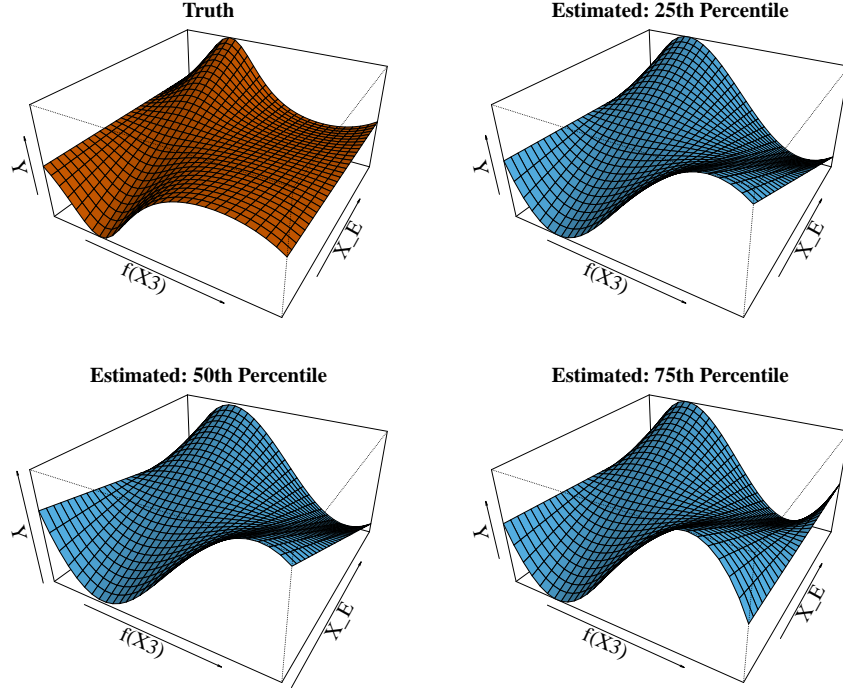


Figure 2: True and estimated interaction effects for $X_E \cdot f_3(X_3)$ in simulation scenario 1a).

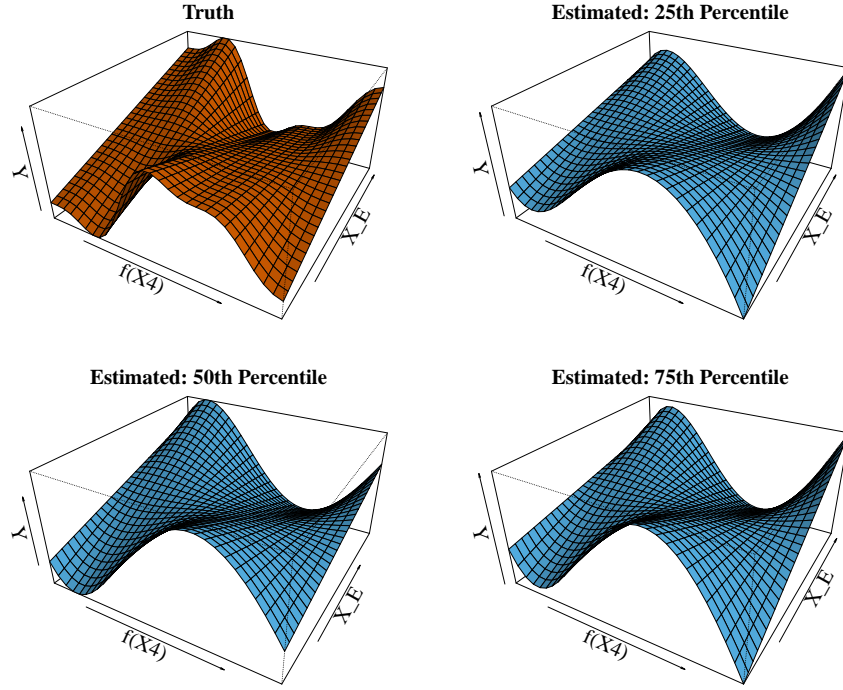


Figure 3: True and estimated interaction effects for $X_E \cdot f_4(X_4)$ in simulation scenario 1a).

5 Real Data Application

In this section we illustrate **sail** on several real data examples.

5.1 Alzheimer’s Disease Neuroimaging Initiative

Alzheimer’s is an irreversible neurodegenerative disease that results in a loss of mental function due to the deterioration of brain tissue. The overall goal of the Alzheimer’s Disease Neuroimaging Initiative (ADNI) is to validate biomarkers for use in Alzheimer’s disease clinical treatment trials [37]. The patients were selected into the study based on their clinical diagnosis: controls, mild cognitive impairment (MCI) or Alzheimer’s disease (AD). PET amyloid imaging was used to assess amyloid beta ($A\beta$) protein load in 96 brain regions. The response we use here is general cognitive decline, as measured by a continuous mini-mental state examination score. We applied **sail** to this data to see if there were any non-linear interactions between clinical diagnosis and $A\beta$ protein in the 96 brain regions on mini-mental state examination.

There were a total of 343 patients who we divided randomly into equal sized training/validation/test splits. We ran the strong heredity **sail** with cubic B-splines and $\alpha = 0.1$. We also applied the **lasso**, **lassoBT**, **HierBasis** and **GLinternet** to this data. Using the same default settings and strategy as the simulation study, we ran each method on the training data, determined the optimal tuning parameter on the validation data, and assessed MSE on the test data. We repeated this process 200 times.

In Figure ?? we plot the mean test set MSE vs. the mean number of active variables ± 1 SD. We see that **sail** produces the sparsest models but doesn’t perform as well as **HierBasis** and **GLinternet** in terms of MSE. **sail** achieves a similar MSE to both the **lasso** and **lassoBT** with fewer variables on average. **GLinternet** produces the largest models and seems to be sensitive to the train/validate/test split as evidenced by the large standard deviations.

To visualize the results from the **sail** method, we chose the train/validate/test split which led to the best test set MSE, and then plotted the interaction effects in Figure ?? . The left panel shows the middle occipital gyrus left region in the occipital lobe known for visual object perception. We see that more $A\beta$ protein loads leads to a worse cognitive score for the MCI and AD group but not for the controls. The right panel shows the cuneus region which is known to be involved in basic visual processing, and we see that more $A\beta$ proteins leads to better cognitive scores for the MCI and AD group and poorer scores for the controls.

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 developed a blockwise coordinate descent algorithm to solve the **sail** objective function for the least-squares loss function. All our algorithms are implemented in a computationally efficient, well-documented and freely available R package. 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 error, sensitivity and specificity when there are non-linear main effects only, as well as interactions with an exposure variable.

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.

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.

References

- [1] Sahir Rai Bhatnagar, Yi Yang, Budhachandra Khundrakpam, Alan C Evans, Mathieu Blanchette, Luigi Bouchard, and Celia MT Greenwood. An analytic approach for interpretable predictive models in high-dimensional data in the presence of interactions with exposures. *Genetic epidemiology*, 42(3):233–249, 2018. 2
- [2] Kaida Ning, Bo Chen, Fengzhu Sun, Zachary Hobel, Lu Zhao, Will Matloff, Arthur W Toga, Alzheimer’s Disease Neuroimaging Initiative, et al. Classifying alzheimer’s disease with brain imaging and genetic data using a neural network framework. *Neurobiology of aging*, 2018. 2
- [3] Eric E Schadt. Molecular networks as sensors and drivers of common human diseases. *Nature*, 461(7261):218–223, 2009. 2
- [4] Nicholas J Timpson, Celia MT Greenwood, Nicole Soranzo, Daniel J Lawson, and J Brent Richards. Genetic architecture: the shape of the genetic contribution to human traits and disease. *Nature Reviews Genetics*, 19(2):110, 2018. 2
- [5] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical Learning with Sparsity: The Lasso and Generalizations*. CRC Press, 2015. 3
- [6] Peter Bühlmann and Sara Van De Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011. 4
- [7] Jianqing Fan, Fang Han, and Han Liu. Challenges of big data analysis. *National science review*, 1(2):293–314, 2014. 4
- [8] Peter McCullagh and John A Nelder. *Generalized linear models*, volume 37. CRC press, 1989. 4
- [9] Hugh Chipman. Bayesian variable selection with related predictors. *Canadian Journal of Statistics*, 24(1):17–36, 1996. 4, 5

-
- [10] David R Cox. Interaction. *International Statistical Review/Revue Internationale de Statistique*, pages 1–24, 1984. 4
- [11] Francis Bach, Rodolphe Jenatton, Julien Mairal, Guillaume Obozinski, et al. Structured sparsity through convex optimization. *Statistical Science*, 27(4):450–468, 2012. 5
- [12] Peter Radchenko and Gareth M James. Variable selection using adaptive nonlinear interaction structures in high dimensions. *Journal of the American Statistical Association*, 105(492):1541–1553, 2010. 5, 7, 30
- [13] Jacob Bien, Jonathan Taylor, Robert Tibshirani, et al. A lasso for hierarchical interactions. *The Annals of Statistics*, 41(3):1111–1141, 2013. 5, 7
- [14] Michael Lim and Trevor Hastie. Learning interactions via hierarchical group-lasso regularization. *Journal of Computational and Graphical Statistics*, 24(3):627–654, 2015. 5, 7, 26
- [15] Asad Haris, Daniela Witten, and Noah Simon. Convex modeling of interactions with strong heredity. *Journal of Computational and Graphical Statistics*, 25(4):981–1004, 2016. 5, 7
- [16] Nam Hee Choi, William Li, and Ji Zhu. Variable selection with the strong heredity constraint and its oracle property. *Journal of the American Statistical Association*, 105(489):354–364, 2010. 5, 7, 10, 11, 14
- [17] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005. 5
- [18] Peng Zhao, Guilherme Rocha, and Bin Yu. The composite absolute penalties family for grouped and hierarchical variable selection. *The Annals of Statistics*, pages 3468–3497, 2009. 7

-
- [19] Yiyuan She and He Jiang. Group regularized estimation under structural hierarchy. *arXiv preprint arXiv:1411.4691*, 2014. 7
- [20] Ning Hao, Yang Feng, and Hao Helen Zhang. Model selection for high-dimensional quadratic regression via regularization. *Journal of the American Statistical Association*, pages 1–11, 2018. 7
- [21] Rajen D Shah. Modelling interactions in high-dimensional data with backtracking. *Journal of Machine Learning Research*, 17(207):1–31, 2016. 7, 26
- [22] Robert Tibshirani and Jerome Friedman. A pliable lasso. *arXiv preprint arXiv:1712.00484*, 2017. 7
- [23] Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American statistical Association*, 96(456):1348–1360, 2001. 8, 11, 14
- [24] Hui Zou. The adaptive lasso and its oracle properties. *Journal of the American statistical association*, 101(476):1418–1429, 2006. 8, 9, 11, 24, 26
- [25] Hansheng Wang, Guodong Li, and Chih-Ling Tsai. Regression coefficient and autoregressive order shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(1):63–78, 2007. 10, 11
- [26] Yuval Nardi, Alessandro Rinaldo, et al. On the asymptotic properties of the group lasso estimator for linear models. *Electronic Journal of Statistics*, 2:605–633, 2008. 10, 11
- [27] Hui Zou and Hao Helen Zhang. On the adaptive elastic-net with a diverging number of parameters. *Annals of statistics*, 37(4):1733, 2009. 11
- [28] Jerome Friedman, Trevor Hastie, and Rob Tibshirani. Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1):1, 2010. 19, 21, 46

-
- [29] Yi Yang and Hui Zou. A fast unified algorithm for solving group-lasso penalize learning problems. *Statistics and Computing*, 25(6):1129–1141, 2015. 21, 46
- [30] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996. 26
- [31] Asad Haris, Ali Shojaie, and Noah Simon. Nonparametric regression with adaptive truncation via a convex hierarchical penalty. *arXiv preprint arXiv:1611.09972*, 2016. 26
- [32] Pradeep Ravikumar, John Lafferty, Han Liu, and Larry Wasserman. Sparse additive models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(5):1009–1030, 2009. 26
- [33] Alexandra Chouldechova and Trevor Hastie. Generalized additive model selection. *arXiv preprint arXiv:1506.03850*, 2015. 26
- [34] R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2017. 27
- [35] Yi Lin, Hao Helen Zhang, et al. Component selection and smoothing in multivariate nonparametric regression. *The Annals of Statistics*, 34(5):2272–2297, 2006. 27, 29
- [36] Jian Huang, Joel L Horowitz, and Fengrong Wei. Variable selection in nonparametric additive models. *Annals of statistics*, 38(4):2282, 2010. 27, 29
- [37] Ronald Carl Petersen, PS Aisen, LA Beckett, MC Donohue, AC Gamst, DJ Harvey, CR Jack, WJ Jagust, LM Shaw, AW Toga, et al. Alzheimer’s disease neuroimaging initiative (adni) clinical characterization. *Neurology*, 74(3):201–209, 2010. 33
- [38] Yihui Xie. *Dynamic Documents with R and knitr*, volume 29. CRC Press, 2015. 48

A Algorithm Details

In this section we provide more specific details about the algorithms used to solve the **sail** objective function.

A.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( $\mathbf{X}, Y, X_E, \text{basis}, \lambda, \alpha, w_j, w_E, w_{jE}, \epsilon$ ) ▷ Algorithm for solving (21)
2:    $\Psi_j \leftarrow \text{basis}(X_j), \tilde{\Psi}_j \leftarrow X_E \circ \Psi_j$  for  $j = 1, \dots, p$ 
3:   Initialize:  $\beta_0^{(0)} \leftarrow \bar{Y}, \beta_E^{(0)} = \theta_j^{(0)} = \gamma_j^{(0)} \leftarrow 0$  for  $j = 1, \dots, p$ .
4:   Set iteration counter  $k \leftarrow 0$ 
5:    $R^* \leftarrow Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_j (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)}$ 
6:   repeat
7:     • To update  $\gamma = (\gamma_1, \dots, \gamma_p)$ 
8:        $\tilde{X}_j \leftarrow \beta_E^{(k)} \tilde{\Psi}_j \theta_j^{(k)}$  for  $j = 1, \dots, p$ 
9:        $R \leftarrow R^* + \sum_{j=1}^p \gamma_j^{(k)} \tilde{X}_j$ 
10:
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:
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:         
$$\theta_j^{(k)(new)} \leftarrow \arg \min_{\theta_j} \frac{1}{2n} \left\| R - \tilde{X}_j \theta_j \right\|_2^2 + \lambda (1 - \alpha) w_j \|\theta_j\|_2$$

21:
22:        $\Delta = \tilde{X}_j (\theta_j^{(k)} - \theta_j^{(k)(new)})$ 
23:        $R^* \leftarrow R^* + \Delta$ 
24:     • To update  $\beta_E$ 
25:        $\tilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \tilde{\Psi}_j \theta_j^{(k)}$ 
26:        $R \leftarrow R^* + \beta_E^{(k)} \tilde{X}_E$ 
27:
28:       
$$\beta_E^{(k)(new)} \leftarrow S \left( \frac{1}{n \cdot w_E} \tilde{X}_E^\top R, \lambda (1 - \alpha) \right)$$

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

37:
38:        $\Delta = \beta_0^{(k)} - \beta_0^{(k)(new)}$ 
39:        $R^* \leftarrow R^* + \Delta$ 
40:        $k \leftarrow k + 1$ 
41:   until convergence criterion is satisfied:  $|Q(\Theta^{(k-1)}) - Q(\Theta^{(k)})| / Q(\Theta^{(k-1)}) < \epsilon$ 

```

A.2 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

$$R_s = Y - \tilde{Y}_{(-s)}^{(k)} \quad (37)$$

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

$$\tilde{Y}_{(-s)}^{(k)} = \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_{\ell \neq s} \Psi_\ell \theta_\ell^{(k)} - \sum_{\ell \neq s} \gamma_\ell^{(k)} \beta_E^{(k)} \tilde{\Psi}_\ell \theta_\ell^{(k)} \quad (38)$$

Using (38), (37) can be re-written as

$$\begin{aligned} R_s &= Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \\ &= R^* + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_s^{(k)} \end{aligned} \quad (39)$$

where

$$R^* = Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_{j=1}^p (\Psi_j + \gamma_j^{(k)} \beta_E^{(k)} \tilde{\Psi}_j) \theta_j^{(k)} \quad (40)$$

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

$$\theta_s^{(k)(new)} = \arg \min_{\theta_j} \frac{1}{2n} \left\| R_s - (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) \theta_j \right\|_2^2 + \lambda(1 - \alpha) w_s \|\theta_j\|_2 \quad (41)$$

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

$$R_{s+1} = R^* + (\Psi_{s+1} + \gamma_{s+1}^{(k)} \beta_E^{(k)} \tilde{\Psi}_{s+1}) \theta_{s+1}^{(k)} + (\Psi_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\Psi}_s) (\theta_s^{(k)} - \theta_s^{(k)(new)}) \quad (42)$$

where R^* is given by (40), $\boldsymbol{\theta}_s^{(k)}$ is the parameter value prior to the update, and $\boldsymbol{\theta}_s^{(k)(new)}$ is the updated value given by (41). Taking the difference between (39) and (42) gives

$$\begin{aligned}\Delta &= R_t - R_s \\ &= (\boldsymbol{\Psi}_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\boldsymbol{\Psi}}_t) \boldsymbol{\theta}_t^{(k)} + (\boldsymbol{\Psi}_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\boldsymbol{\Psi}}_s) (\boldsymbol{\theta}_s^{(k)} - \boldsymbol{\theta}_s^{(k)(new)}) - (\boldsymbol{\Psi}_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\boldsymbol{\Psi}}_s) \boldsymbol{\theta}_s^{(k)} \\ &= (\boldsymbol{\Psi}_t + \gamma_t^{(k)} \beta_E^{(k)} \tilde{\boldsymbol{\Psi}}_t) \boldsymbol{\theta}_t^{(k)} - (\boldsymbol{\Psi}_s + \gamma_s^{(k)} \beta_E^{(k)} \tilde{\boldsymbol{\Psi}}_s) \boldsymbol{\theta}_s^{(k)(new)}\end{aligned}\quad (43)$$

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

A.3 Least-Squares sail with Weak Heredity

The least-squares **sail** model with weak heredity has the form

$$\hat{Y} = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \boldsymbol{\Psi}_j \boldsymbol{\theta}_j + \beta_E X_E + \sum_{j=1}^p \gamma_j (X_E \circ \boldsymbol{\Psi}_j) (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j) \quad (44)$$

The objective function is given by

$$Q(\boldsymbol{\Theta}) = \frac{1}{2n} \|Y - \hat{Y}\|_2^2 + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda \alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (45)$$

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

$$\frac{\partial Q}{\partial \beta_0} = \frac{1}{n} \left(Y - \beta_0 \cdot \mathbf{1} - \sum_{j=1}^p \Psi_j \boldsymbol{\theta}_j - \beta_E X_E - \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j) \right)^\top \mathbf{1} = 0 \quad (46)$$

$$\frac{\partial Q}{\partial \beta_E} = -\frac{1}{n} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^\top R + \lambda(1 - \alpha) w_E s_1 = 0 \quad (47)$$

$$\frac{\partial Q}{\partial \boldsymbol{\theta}_j} = -\frac{1}{n} (\Psi_j + \gamma_j (X_E \circ \Psi_j))^\top R + \lambda(1 - \alpha) w_j s_2 = \mathbf{0} \quad (48)$$

$$\frac{\partial Q}{\partial \gamma_j} = -\frac{1}{n} ((X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j))^\top R + \lambda \alpha w_{jE} s_3 = 0 \quad (49)$$

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

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

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

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

From the subgradient Equation (47), we see that $\beta_E = 0$ is a solution if

$$\frac{1}{w_E} \left| \frac{1}{n} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \Psi_j) \mathbf{1}_{m_j} \right)^\top R_{(-E)} \right| \leq \lambda(1 - \alpha) \quad (50)$$

From the subgradient Equation (48), we see that $\theta_j = \mathbf{0}$ is a solution if

$$\frac{1}{w_j} \left\| \frac{1}{n} (\Psi_j + \gamma_j (X_E \circ \Psi_j))^\top R_{(-j)} \right\|_2 \leq \lambda(1 - \alpha) \quad (51)$$

From the subgradient Equation (49), we see that $\gamma_j = 0$ is a solution if

$$\frac{1}{w_{jE}} \left| \frac{1}{n} ((X_E \circ \Psi_j) (\beta_E \cdot \mathbf{1}_{m_j} + \theta_j))^\top R_{(-jE)} \right| \leq \lambda\alpha \quad (52)$$

From the subgradient equations we see that

$$\hat{\beta}_0 = \left(Y - \sum_{j=1}^p \Psi_j \hat{\theta}_j - \hat{\beta}_E X_E - \sum_{j=1}^p \hat{\gamma}_j (X_E \circ \Psi_j) (\hat{\beta}_E \cdot \mathbf{1}_{m_j} + \hat{\theta}_j) \right)^\top \mathbf{1} \quad (53)$$

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

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

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

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 are closed form solutions for the intercept and β_E , each γ_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 [28], while we use the quadratic majorization technique implemented in the **gglasso** package [29] to solve (55). Algorithm 4 details the procedure used to fit the least-squares weak heredity **sail** model.

A.3.1 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 **0** is:

$$\lambda_{max} = \frac{1}{n} \max \left\{ \frac{1}{(1 - \alpha)w_E} \left(X_E + \sum_{j=1}^p \gamma_j (X_E \circ \boldsymbol{\Psi}_j) \mathbf{1}_{m_j} \right)^\top R_{(-E)}, \right. \\ \max_j \frac{1}{(1 - \alpha)w_j} \left\| (\boldsymbol{\Psi}_j + \gamma_j (X_E \circ \boldsymbol{\Psi}_j))^\top R_{(-j)} \right\|_2, \\ \left. \max_j \frac{1}{\alpha w_{jE}} \left((X_E \circ \boldsymbol{\Psi}_j) (\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j) \right)^\top R_{(-jE)} \right\} \quad (57)$$

which reduces to

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

This is the same λ_{max} as the least-squares strong heredity **sail** model.

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

```

1: function sail( $\mathbf{X}, Y, X_E, \text{basis}, \lambda, \alpha, w_j, w_E, w_{jE}, \epsilon$ ) ▷ Algorithm for solving (45)
2:    $\Psi_j \leftarrow \text{basis}(X_j), \tilde{\Psi}_j \leftarrow X_E \circ \Psi_j$  for  $j = 1, \dots, p$ 
3:   Initialize:  $\beta_0^{(0)} \leftarrow \bar{Y}, \beta_E^{(0)} = \boldsymbol{\theta}_j^{(0)} = \gamma_j^{(0)} \leftarrow 0$  for  $j = 1, \dots, p$ .
4:   Set iteration counter  $k \leftarrow 0$ 
5:    $R^* \leftarrow Y - \beta_0^{(k)} - \beta_E^{(k)} X_E - \sum_j \Psi_j \boldsymbol{\theta}_j^{(k)} - \sum_j \gamma_j^{(k)} \tilde{\Psi}_j (\beta_E^{(k)} \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j^{(k)})$ 
6:   repeat
7:     • To update  $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)$ 
8:        $\tilde{X}_j \leftarrow \tilde{\Psi}_j (\beta_E^{(k)} \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j^{(k)})$  for  $j = 1, \dots, p$ 
9:        $R \leftarrow R^* + \sum_{j=1}^p \gamma_j^{(k)} \tilde{X}_j$ 
10:
11:         
$$\boldsymbol{\gamma}^{(k)(new)} \leftarrow \arg \min_{\boldsymbol{\gamma}} \frac{1}{2n} \left\| R - \sum_j \gamma_j \tilde{X}_j \right\|_2^2 + \lambda \alpha \sum_j w_{jE} |\gamma_j|$$

12:
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:
22:        $\Delta = \tilde{X}_j (\boldsymbol{\theta}_j^{(k)} - \boldsymbol{\theta}_j^{(k)(new)})$ 
23:        $R^* \leftarrow R^* + \Delta$ 
24:     • To update  $\beta_E$ 
25:        $\tilde{X}_E \leftarrow X_E + \sum_j \gamma_j^{(k)} \tilde{\Psi}_j \mathbf{1}_{m_j}$ 
26:        $R \leftarrow R^* + \beta_E^{(k)} \tilde{X}_E$ 
27:
28:       
$$\beta_E^{(k)(new)} \leftarrow S \left( \frac{1}{n \cdot w_E} \tilde{X}_E^\top R, \lambda (1 - \alpha) \right)$$

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

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

```

B Simulation Results

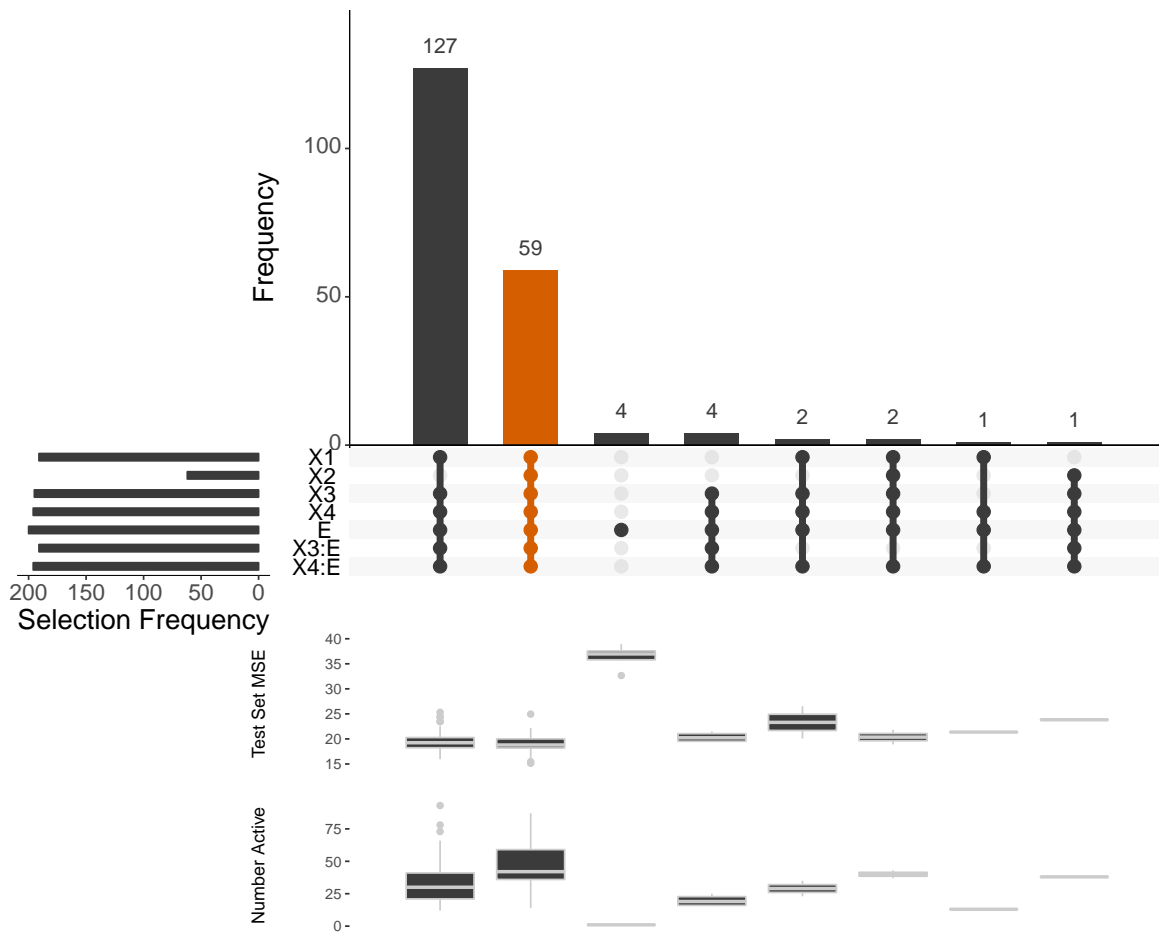


Figure B.1: Selection rates across 200 simulations of scenario 1a) for strong heredity `sail`.

C `sail` Package Showcase

In this section we briefly introduce the freely available and open source `sail` package in R. More comprehensive documentation is available at <https://sahirbhatnagar.com/sail>. Note that this entire section is reproducible; the code and text are combined in an `.Rnw`¹ file and compiled using `knitr` [38].

¹scripts available at <https://github.com/sahirbhatnagar/sail/tree/master/manuscript>

C.1 Installation

The package can be installed from [GitHub](#) via

```
install.packages("pacman")
pacman::p_load_gh('sahirbhatnagar/sail')
```

C.2 Quick Start

We give a quick overview of the main functions and go into details in other vignettes. We will use the simulated data which ships with the package and can be loaded via:

```
library(sail)
data("sailsim")
names(sailsim)
```

We first define a basis expansion. In this example we use B-splines with degree 5.

```
library(splines)
f.basis <- function(x) splines::bs(x, degree = 5)
```

Next we fit the model using the most basic call to `sail`

```
fit <- sail(x = sailsim$x, y = sailsim$y, e = sailsim$e, basis = f.basis)
```

`fit` is an object of class `sail` that contains all the relevant information of the fitted model including the estimated coefficients at each value of λ (by default the program chooses its own decreasing sequence of 100 λ values). There are `print`, `plot`, `coef` and `predict` methods of objects of class `sail`.

When `expand = TRUE` (i.e. the user did not provide their own design matrix), the `df_main` and `df_interaction` columns correspond to the number of non-zero predictors present in the model before basis expansion. This does not correspond to the number of non-zero coefficients in the model, but rather the number of unique variables. In this example we expanded each column of \mathbf{X} to five columns. If `df_main=4`, `df_interaction=2` and `df_environment=1`, then the total number of non-zero coefficients would be $5 \times (4 + 2) +$

1.

The entire solution path can be plotted via the `plot` method for objects of class `sail`. The y-axis is the value of the coefficient and the x-axis is the $\log(\lambda)$. Each line represents a coefficient in the model, and each color represents a variable (i.e. in this example a given variable will have 5 lines when it is non-zero). The numbers at the top of the plot represent the number of non-zero variables in the model: top panel (`df_main + df_environment`), bottom panel (`df_interaction`). The black line is the coefficient path for the environment variable.

```
plot(fit)
```

The estimated coefficients at each value of lambda is given by (matrix partially printed here for brevity)

```
coef(fit)[1:6,50:55]
```

The corresponding predicted response at each value of lambda (matrix partially printed here for brevity):

```
predict(fit)[1:5,50:55]
```

The predicted response at a specific value of lambda can be specified by the `s` argument:

```
predict(fit, s = 0.8)[1:5, ]
```

You can specify more than one value for ‘s’:

```
predict(fit, s = c(0.8, 0.2))[1:5, ]
```

You can also extract a list of active variables (i.e. variables with a non-zero estimated coefficient) for each value of lambda:

```
fit[["active"]][50:55]
```

C.3 Cross-Validation

`cv.sail` is the main function to do cross-validation along with `plot`, `predict`, and `coef` methods for objects of class `cv.sail`. We run it in parallel:

```
set.seed(432) # to reproduce results (randomness due to CV folds)
library(doMC)
registerDoMC(cores = 8)
cvfit <- cv.sail(x = sailsim$x, y = sailsim$y, e = sailsim$e, basis = f.basis,
nfolds = 5, parallel = TRUE)
```

We plot the cross-validated error curve which has the mean-squared error on the y-axis and $\log(\lambda)$ on the x-axis. It includes the cross-validation curve (red dotted line), and upper and lower standard deviation curves along the λ sequence (error bars). Two selected λ 's are indicated by the vertical dotted lines (see below). The numbers at the top of the plot represent the total number of non-zero variables at that value of λ (`df_main + df_environment + df_interaction`):

```
plot(cvfit)
```

`lambda.min` is the value of λ that gives minimum mean cross-validated error. The other λ saved is `lambda.1se`, which gives the most regularized model such that error is within one standard error of the minimum. We can view the selected λ 's and the corresponding coefficients:

```
cvfit[["lambda.min"]]
cvfit[["lambda.1se"]]
```

The estimated nonzero coefficients at `lambda.1se` and `lambda.min`:

```
predict(cvfit, type = "nonzero", s="lambda.1se") # lambda.1se is the default
predict(cvfit, type = "nonzero", s = "lambda.min")
```

C.4 Visualizing the Effect of the Non-linear Terms

B-splines are difficult to interpret. We provide a plotting function to visualize the effect of the non-linear function on the response.

C.4.1 Main Effects

Since we are using simulated data, we also plot the true curve:

```
plotMain(cvfit$sail.fit, x = sailsim$x, xvar = "X3",  
legend.position = "topright",  
s = cvfit$lambda.min, f.truth = sailsim$f3)
```

C.4.2 Interaction Effects

Again, since we are using simulated data, we also plot the true interaction:

```
plotInter(cvfit$sail.fit, x = sailsim$x, xvar = "X4",  
f.truth = sailsim$f4.inter,  
s = cvfit$lambda.min,  
title_z = "Estimated")
```

C.5 Linear Interactions

The `basis` argument in the `sail` function is very flexible in that it allows you to apply *any* basis expansion to the columns of \mathbf{X} . Of course, there might be situations where you do not expect any non-linear main effects or interactions to be present in your data. You can still use the `sail` method to search for linear main effects and interactions. This can be accomplished by specifying an identity map:

```
f.identity <- function(i) i
```

We then pass this function to the `basis` argument in `cv.sail`:

```
cvfit_linear <- cv.sail(x = sailsim$x, y = sailsim$y, e = sailsim$e,
basis = f.identity, nfolds = 5, parallel = TRUE)
```

Next we plot the cross-validated curve:

```
plot(cvfit_linear)
```

And extract the model at `lambda.min`:

```
predict(cvfit_linear, s = "lambda.min", type = "nonzero")
```

C.6 Applying a different penalty to each predictor

Recall that we consider the following penalized least squares criterion for this problem:

$$\arg \min_{\boldsymbol{\theta}} \mathcal{L}(Y; \boldsymbol{\theta}) + \lambda(1 - \alpha) \left(w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda \alpha \sum_{j=1}^p w_{jE} |\gamma_j| \quad (58)$$

The weights w_E, w_j, w_{jE} are by default set to 1 as specified by the `penalty.factor` argument. This argument allows users to apply separate penalty factors to each coefficient. In particular, any variable with `penalty.factor` equal to zero is not penalized at all. This feature can be applied mainly for 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](#)

In the following example, we want the environment variable to always be included so we set the first element of `p.fac` to zero. We also want to apply less of a penalty to the main effects for X_2, X_3, X_4 :

```
# the weights correspond to E, X1, X2, X3, ... X_p, X1:E, X2:E, ... X_p:E
p.fac <- c(0, 1, 0.4, 0.6, 0.7, rep(1, 2*ncol(sailsim$x) - 4))

fit_pf <- sail(x = sailsim$x, y = sailsim$y, e = sailsim$e, basis = f.basis,
penalty.factor = p.fac)
```

```
plot(fit_pf)
```

We see from the plot above that the black line (corresponding to the X_E variable with `penalty.factor` equal to zero) is always included in the model.

C.7 User-Defined Design Matrix

A limitation of the `sail` method is that the same basis expansion function $f(\cdot)$ is applied to all columns of the predictor matrix \mathbf{X} . Being able to automatically select linear vs. nonlinear components was not a focus of our paper, but is an active area of research for main effects only e.g. [gamsel](#) and [HierBasis](#).

However, if the user has some prior knowledge on possible effect relationships, then they can supply their own design matrix. This can be useful for example, when one has a combination of categorical (e.g. gender, race) and continuous variables, but would only like to apply $f(\cdot)$ on the continuous variables. We provide an example below to illustrate this functionality.

We use the simulated dataset `sailsim` provided in our package. We first add a categorical variable `race` to the data:

```
set.seed(1234)
library(sail)
x_df <- as.data.frame(sailsim$x)
x_df$race <- factor(sample(1:2, nrow(x_df), replace = TRUE))
table(x_df$race)
```

We then use the `model.matrix` function to create the design matrix. Note that the intercept should not be included, as this is computed internally in the `sail` function. This is why we add 0 to the formula. Notice also the flexibility we can have by including different basis expansions to each predictor:

```
library(splines)
x <- stats::model.matrix(~ 0 + bs(X1, degree = 5) + bs(X2, degree = 3) + ns(X3, df = 8) +
bs(X4, degree = 6) + X5 + poly(X6,2) + race, data = x_df)
head(x)
```

One benefit of using `stats::model.matrix` is that it returns the group membership as an attribute:

```
attr(x, "assign")
```

The group membership must be supplied to the `sail` function. This information is needed for the group lasso penalty, which will select the whole group as zero or non-zero.

C.7.1 Fit the `sail` Model

We need to set the argument `expand = FALSE` and provide the group membership. The first element of the group membership corresponds to the first column of `x`, the second element to the second column of `x`, and so on.

We can plot the solution path for both main effects and interactions using the `plot` method for objects of class `sail`:

In this instance, since we provided a user-defined design matrix and ‘`expand = FALSE`’, the numbers at the top of the plot represent the total number of non-zero coefficients.

C.7.2 Find the Optimal Value for λ

We can use cross-validation to find the optimal value of `lambda`:

We can plot the cross-validated mean squared error as a function of `lambda`:

The estimated non-zero coefficients at `lambda.1se`: