

# Sparse Additive Interaction Learning

Sahir Bhatnagar

Department of Epidemiology, Biostatistics and Occupational Health  
Department of Diagnostic Radiology

July 12, 2021

<https://sahirbhatnagar.com/sail>



# Outline

Betting on Sparsity

A Thought Experiment

Motivating Example: The Nurse Family Partnership

sail: Strong Additive Interaction Learning

Algorithm

Theory

Simulations

sail R package

Real Data Application

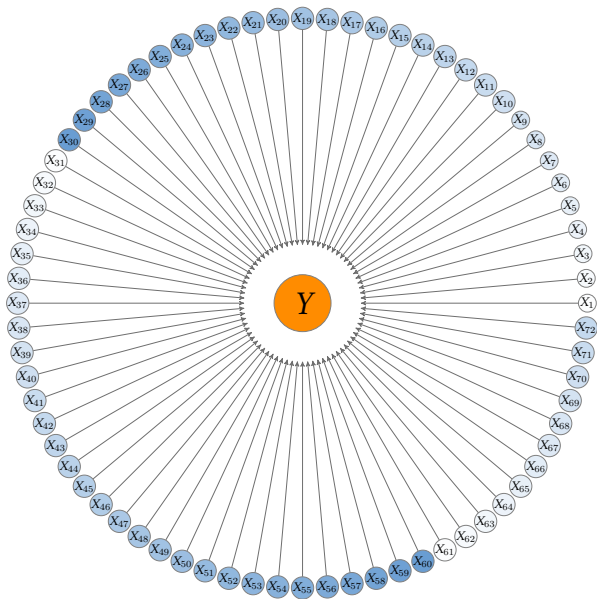
Discussion

Current and Future Work

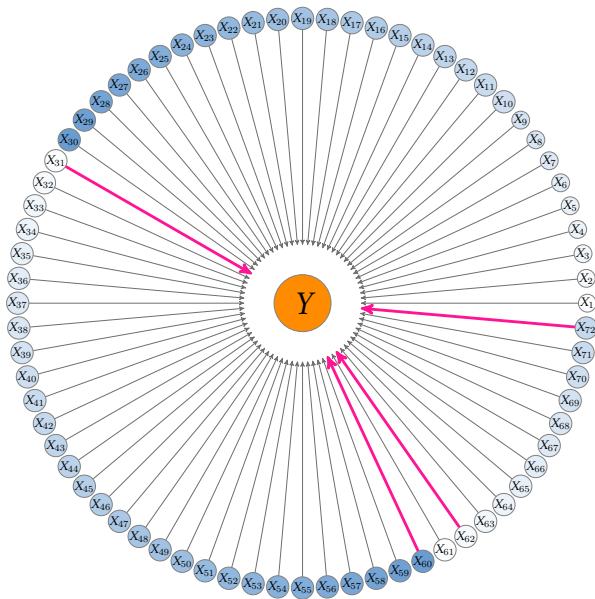
Acknowledgements



# Bet on Sparsity Principle



# Bet on Sparsity Principle



# Bet on Sparsity Principle

**Use a procedure that does well in sparse problems, since no procedure does well in dense problems.<sup>1</sup>**

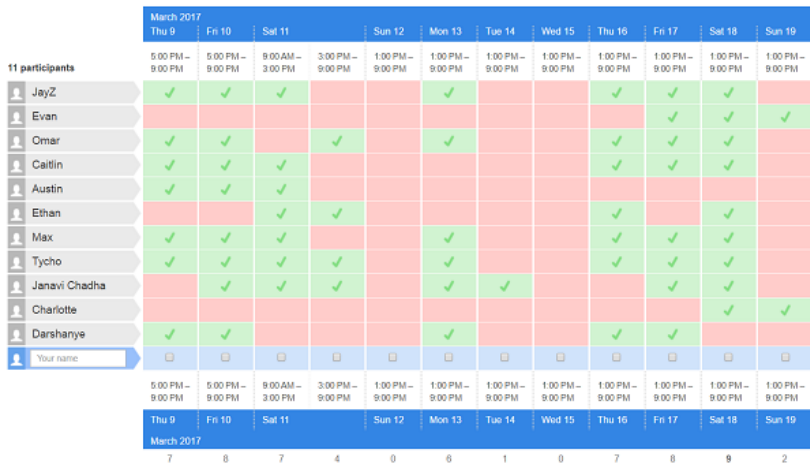
- We often don't have enough data to estimate so many parameters
- Even when we do, we might want to identify a **relatively small number of predictors** ( $k < N$ ) that play an important role
- Faster computation, easier to understand, and stable predictions on new datasets.

---

<sup>1</sup>The elements of statistical learning. Springer series in statistics, 2001.

How would you schedule a meeting of 20 people?

# How would you schedule a meeting of 20 people?





# Doctors Bet on Sparsity Also

## Doctors Bet on Sparsity Also



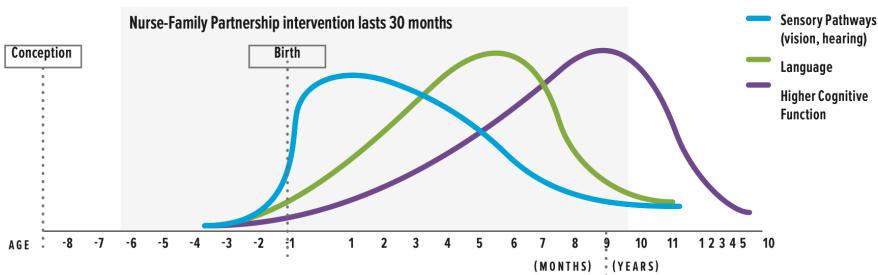




**Nurse-Family Partnership is an evidence-based, community health program with over 40 years of evidence showing significant improvements in the health and lives of first-time moms and their children living in poverty.**

## Human Brain Development

Synapse formation dependent on early experiences



Source: Nelson, C.A., In *Neurons to Neighborhoods* (2000).



## ALL CAUSE MORTALITY OVER 20 YEAR FOLLOW-UP

**3x**

Mothers who did not receive nurse home visits were nearly **3 times more likely to die** from all causes of death than nurse visited mothers (3.7% versus 1.3%)<sup>1</sup>

**8x**

Mothers that did not receive nurse home visits were **8 times more likely to die** from external causes – including unintentional injuries, suicide, drug overdose and homicide – than nurse visited mothers (1.7% versus 0.2%)<sup>1</sup>

## Additional Maternal and Child Health Outcomes

### Maternal Health Outcomes

**35%** fewer cases of pregnancy-induced hypertension<sup>5</sup>

**18%** fewer preterm births<sup>6</sup>

**79%** reduction in preterm delivery among women who smoke cigarettes<sup>7</sup>

**31%** reduction in very closely spaced (<6 months) subsequent pregnancies<sup>8</sup>

### Child Health Outcomes

**48%** reduction in child abuse and neglect<sup>9</sup>

**39%** fewer health care encounters for injuries or ingestions in the first 2 years of life among children born to mothers with low psychological resources<sup>10</sup>

**67%** less behavioral and intellectual problems in children at age 6<sup>11</sup>

**56%** fewer emergency room visits for accidents and poisonings through age 2<sup>12</sup>



## PREVENTABLE CHILD MORTALITY OVER 20 YEAR FOLLOW-UP

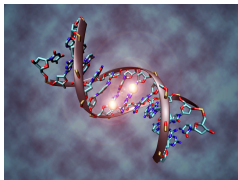
- Among Nurse-Family Partnership participants, there were **lower rates of preventable child mortality** from birth until age 20.<sup>1</sup>
- 1.6% of the children not receiving nurse home visits died from preventable causes – including sudden infant death syndrome, unintentional injuries and homicide – while none of the nurse visited children died from these causes.<sup>1</sup>

# Interactions between Intervention and Genetics

Stanford-Binet Fifth Edition (SB5) classification<sup>[4]</sup>

IQ Range ("deviation IQ")	IQ Classification
145-160	Very gifted or highly advanced
130-144	Gifted or very advanced
120-129	Superior
110-119	High average
90-109	Average
80-89	Low average
70-79	Borderline impaired or delayed
55-69	Mildly impaired or delayed
40-54	Moderately impaired or delayed

~



×



**Phenotype**  
**IQ Score**

**Large Data**  
**Genetic Markers**

**Environment**  
**NFP Intervention**

$$Y = \sum_{j=1}^p X_j \beta_j + \sum_{j=1}^p X_j X_E \tau_j + \varepsilon$$

$$\begin{aligned}
 Y &= \sum_{j=1}^p X_j \beta_j & + & \sum_{j=1}^p X_j X_E \tau_j & + & \varepsilon \\
 &= \begin{array}{c} \text{blue square} \\ \mathbf{X} \\ n \times p \end{array} \begin{array}{c} \text{yellow rectangle} \\ \boldsymbol{\beta} \\ p \times 1 \end{array} & + & \begin{array}{c} \text{red rectangle} \\ X_E \\ n \times 1 \end{array} \circ \begin{array}{c} \text{blue square} \\ \mathbf{X} \\ n \times p \end{array} \begin{array}{c} \text{green rectangle} \\ \boldsymbol{\tau} \\ p \times 1 \end{array} & + & \begin{array}{c} \text{gray rectangle} \\ \boldsymbol{\varepsilon} \\ n \times 1 \end{array}
 \end{aligned}$$



$$\begin{aligned}
 Y &= \sum_{j=1}^p X_j \beta_j & + & \sum_{j=1}^p X_j X_E \tau_j & + & \varepsilon \\
 &= & & & & \\
 &= \begin{array}{c} \text{Main effects} \\ \begin{array}{c} \text{Blue square } \mathbf{X} \\ n \times p \\ \text{Yellow vertical bar } \boldsymbol{\beta} \\ p \times 1 \end{array} & + & \begin{array}{c} \text{Interaction effects} \\ \begin{array}{c} \text{Red vertical bar } X_E \\ n \times 1 \\ \text{Blue square } \mathbf{X} \\ n \times p \\ \text{Green vertical bar } \boldsymbol{\tau} \\ p \times 1 \end{array} & + & \begin{array}{c} \text{Error} \\ \text{Grey vertical bar } \boldsymbol{\varepsilon} \\ n \times 1 \end{array}
 \end{array}
 \end{aligned}$$

Main effects

Interaction effects

Error

$$\begin{aligned}
 Y &= \sum_{j=1}^p X_j \beta_j & + & \sum_{j=1}^p X_j X_E \tau_j & + & \varepsilon \\
 &= \begin{array}{c} \text{Main effects} \\ \begin{array}{c} \mathbf{X} \\ n \times p \end{array} \begin{array}{c} \boldsymbol{\beta} \\ p \times 1 \end{array} & + & \begin{array}{c} \text{Interaction effects} \\ \begin{array}{c} X_E \circ \mathbf{X} \\ n \times 1 \quad n \times p \end{array} \begin{array}{c} \boldsymbol{\tau} \\ p \times 1 \end{array} & + & \begin{array}{c} \boldsymbol{\varepsilon} \\ n \times 1 \end{array} \\
 & & & & & \text{Error}
 \end{array}
 \end{aligned}$$

Let  $Z_{jE} = X_E X_j$

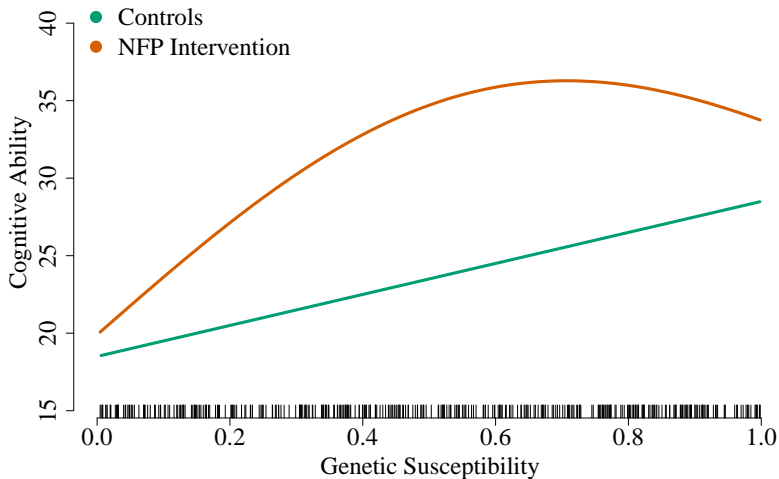
$$Y = \sum_{j=1}^p X_j \beta_j + \sum_{j=1}^p Z_{jE} \tau_j + \epsilon$$

$=$

The diagram illustrates the matrix equation  $Y = X\beta + Z\tau + \epsilon$ . Matrix  $X$  is a blue rectangle labeled  $X$  with dimensions  $n \times 2p$ . Matrix  $Z$  is a blue rectangle labeled  $Z$  with dimensions  $n \times p$ . Vector  $\beta$  is a vertical bar with a yellow-to-orange gradient, labeled  $\beta$  with dimensions  $2p \times 1$ . Vector  $\tau$  is a vertical bar with a green-to-dark-green gradient, labeled  $\tau$  with dimensions  $2p \times 1$ . Vector  $\epsilon$  is a gray vertical bar labeled  $\epsilon$  with dimensions  $n \times 1$ .



## Motivation 1: Non-linear Interactions



## Motivation 2: Heredity Property

$$Y = \beta_0 \cdot \mathbf{1} + \underbrace{\sum_{j=1}^p \beta_j X_j + \beta_E X_E}_{\text{main effects}} + \underbrace{\sum_{j=1}^p \tau_j X_E X_j}_{\text{interactions}} + \varepsilon$$

---

<sup>1</sup>Chipman. Canadian Journal of Statistics (1996)

<sup>2</sup>McCullagh and Nelder. Generalized Linear Models (1983)

<sup>3</sup>Cox. International Statistical Review (1984)

## Motivation 2: Heredity Property

$$Y = \beta_0 \cdot \mathbf{1} + \underbrace{\sum_{j=1}^p \beta_j X_j + \beta_E X_E}_{\text{main effects}} + \underbrace{\sum_{j=1}^p \tau_j X_E X_j}_{\text{interactions}} + \varepsilon$$

### Strong Heredity<sup>1</sup>

$$\hat{\tau}_j \neq 0 \quad \Rightarrow \quad \hat{\beta}_j \neq 0 \quad \text{and} \quad \hat{\beta}_E \neq 0$$

- Heredity property is desired for the purposes of **interpretability**<sup>2</sup>
- Large main effects are more likely to lead to appreciable interactions<sup>3</sup>

---

<sup>1</sup>Chipman. Canadian Journal of Statistics (1996)

<sup>2</sup>McCullagh and Nelder. Generalized Linear Models (1983)

<sup>3</sup>Cox. International Statistical Review (1984)

# Lasso interaction model

- $Y \rightarrow$  response
- $X_E \rightarrow$  environment
- $X_j \rightarrow$  predictors,  $j = 1, \dots, p$

$$Y = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \beta_j X_j + \beta_E X_E + \sum_{j=1}^p \tau_j X_E X_j + \varepsilon$$

$$\underset{\Theta := (\beta_0, \beta, \tau)}{\operatorname{argmin}} \quad \mathcal{L}(\Theta) + \lambda(\|\beta\|_1 + \|\tau\|_1)$$



# Strong Heredity Interactions: Current State of the Art

Type	Model	Software
Linear	CAP (Zhao et al. 2009, <a href="#">Ann. Stat</a> )	$\times$
	SHIM (Choi et al. 2009, <a href="#">JASA</a> )	$\times$
	hiernet (Bien et al. 2013, <a href="#">Ann. Stat</a> )	hierNet(x, y)
	GRESH (She and Jiang 2014, <a href="#">JASA</a> )	$\times$
	FAMILY (Haris et al. 2014, <a href="#">JCGS</a> )	FAMILY(x, z, y)
	glinternet (Lim and Hastie 2015, <a href="#">JCGS</a> )	glinternet(x, y)
	RAMP (Hao et al. 2016, <a href="#">JASA</a> )	RAMP(x, y)
	LassoBacktracking (Shah 2018, <a href="#">JMLR</a> )	LassoBT(x, y)
Non-linear	VANISH (Radchenko and James 2010, <a href="#">JASA</a> )	$\times$
	<a href="#">sail</a> (Bhatnagar et al. 2020+, in revision <a href="#">CSDA</a> )	<code>sail(x, e, y, basis)</code>

# Our Extension to Nonlinear Effects

Consider the basis expansion

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

$$f(X_1) = \underbrace{\begin{bmatrix} \psi_{11}(X_{11}) & \psi_{12}(X_{12}) & \cdots & \psi_{11}(X_{15}) \\ \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \cdots & \vdots \\ \psi_{11}(X_{i1}) & \psi_{12}(X_{i2}) & \cdots & \psi_{11}(X_{i5}) \\ \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \cdots & \vdots \\ \psi_{11}(X_{N1}) & \psi_{12}(X_{N2}) & \cdots & \psi_{11}(X_{N5}) \end{bmatrix}}_{\Psi_1} \quad N \times 5 \quad \times \quad \underbrace{\begin{bmatrix} \beta_{11} \\ \beta_{12} \\ \beta_{13} \\ \beta_{14} \\ \beta_{15} \end{bmatrix}}_{\theta_1} \quad 5 \times 1$$

# sail: Additive Interactions

- $\theta_j = (\beta_{j1}, \dots, \beta_{jm_j}) \in \mathbb{R}^{m_j}$
- $\tau_j = (\tau_{j1}, \dots, \tau_{jm_j}) \in \mathbb{R}^{m_j}$
- $\Psi_j \rightarrow n \times m_j$  matrix of evaluations of the  $\psi_{j\ell}$
- In our implementation, we use cubic bsplines with 5 degrees of freedom

## Model

$$Y = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \Psi_j \theta_j + \beta_E X_E + \sum_{j=1}^p (X_E \circ \Psi_j) \tau_j + \varepsilon$$

# sail: Strong Heredity

## Reparametrization<sup>1</sup>

$$\boldsymbol{\tau}_j = \gamma_j \beta_E \boldsymbol{\theta}_j$$

## Model

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

## Objective Function

$$\underset{\boldsymbol{\Theta} := (\beta_E, \boldsymbol{\theta}, \boldsymbol{\gamma})}{\operatorname{argmin}} \quad \mathcal{L}(\boldsymbol{\Theta}) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda \alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

---

<sup>1</sup>Choi et al. JASA (2010)

# sail: Weak Heredity

## Reparametrization

$$\boldsymbol{\tau}_j = \gamma_j(\beta_E \cdot \mathbf{1}_{m_j} + \boldsymbol{\theta}_j)$$

## Model

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

## Objective Function

$$\operatorname{argmin}_{\beta_E, \boldsymbol{\theta}, \boldsymbol{\gamma}} \mathcal{L}(\boldsymbol{\Theta}) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\boldsymbol{\theta}_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$



# Block Relaxation (De Leeuw, 1994)

---

**Algorithm 1:** Block Relaxation Algorithm

---

Set the iteration counter  $k \leftarrow 0$  and fix  $\alpha \in (0, 1)$ ;

**for** each  $\lambda$  **do**

**repeat**

$$\gamma^{(k+1)} \leftarrow \underset{\gamma}{\operatorname{argmin}} \quad Q_{\lambda} \left( \gamma, \beta_E^{(k)}, \boldsymbol{\theta}^{(k)} \right)$$

$$\boldsymbol{\theta}^{(k+1)} \leftarrow \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \quad Q_{\lambda} \left( \boldsymbol{\theta}, \beta_E^{(k)}, \gamma^{(k+1)} \right)$$

$$\beta_E^{(k+1)} \leftarrow \underset{\beta_E}{\operatorname{argmin}} \quad Q_{\lambda} \left( \boldsymbol{\theta}^{(k+1)}, \beta_E, \gamma^{(k+1)} \right)$$

$$k \leftarrow k + 1$$

**until** convergence criterion is satisfied;

**end**

---

# Implementation

## Objective Function

$$\operatorname{argmin}_{\beta_E, \boldsymbol{\theta}, \boldsymbol{\gamma}} \mathcal{L}(Y; \boldsymbol{\Theta}) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

---

<sup>1</sup><https://cran.r-project.org/package=sail>



# Implementation

## Objective Function

$$\operatorname{argmin}_{\beta_E, \theta, \gamma} \mathcal{L}(Y; \Theta) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

## Lasso problem

$$\operatorname{argmin}_{\gamma} \mathcal{L}(Y; \Theta) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

---

<sup>1</sup><https://cran.r-project.org/package=sail>

# Implementation

## Objective Function

$$\operatorname{argmin}_{\beta_E, \boldsymbol{\theta}, \boldsymbol{\gamma}} \mathcal{L}(Y; \boldsymbol{\Theta}) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

---

<sup>1</sup><https://cran.r-project.org/package=sail>

# Implementation

## Objective Function

$$\operatorname{argmin}_{\beta_E, \boldsymbol{\theta}, \gamma} \mathcal{L}(Y; \boldsymbol{\Theta}) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

## Group Lasso problem

$$\operatorname{argmin}_{\beta_E, \boldsymbol{\theta}} \mathcal{L}(Y; \boldsymbol{\Theta}) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

---

<sup>1</sup><https://cran.r-project.org/package=sail>



# Sparsity

## Theorem 1

$$\widehat{\Theta}_n = \operatorname{argmin}_{\beta_E, \theta, \gamma} \mathcal{L}(\Theta) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

$$\mathcal{A}_1 = \{j : \theta_j \neq 0, \beta_j \neq 0\}$$

$$\mathcal{A}_2 = \{k : \gamma_k \neq 0\}, \quad \mathcal{A} = \mathcal{A}_1 \cup \mathcal{A}_2$$

Under certain regularity conditions and the existence of a local minimizer  $\widehat{\Theta}_n$  that is  $\sqrt{n}$ -consistent

$$P\left(\widehat{\Theta}_{\mathcal{A}^c} = 0\right) \rightarrow 1$$

## Theorem 1

$$\widehat{\Theta}_n = \underset{\beta_E, \theta, \gamma}{\operatorname{argmin}} \quad \mathcal{L}(\Theta) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

$$\mathcal{A}_1 = \{j : \theta_j \neq 0, \beta_j \neq 0\}$$

$$\mathcal{A}_2 = \{k : \gamma_k \neq 0\}, \quad \mathcal{A} = \mathcal{A}_1 \cup \mathcal{A}_2$$

Under certain regularity conditions and the existence of a local minimizer  $\widehat{\Theta}_n$  that is  $\sqrt{n}$ -consistent

$$P\left(\widehat{\Theta}_{\mathcal{A}^c} = 0\right) \rightarrow 1$$

Theorem 1 shows that when the tuning parameters for the nonzero coefficients converge to 0 faster than  $n^{-1/2}$  we can consistently remove the noise terms with probability tending to 1.

# Asymptotic normality

## Theorem 2

$$\widehat{\Theta}_n = \operatorname{argmin}_{\beta_E, \theta, \gamma} \mathcal{L}(\Theta) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

Under certain regularity conditions, the component  $\widehat{\Theta}_{\mathcal{A}}$  of the local minimizer  $\widehat{\Theta}_n$  satisfies

$$\sqrt{n} \left( \widehat{\Theta}_{\mathcal{A}} - \Theta_{\mathcal{A}} \right) \rightarrow_d \mathcal{N} \left( 0, \mathbf{I}^{-1} \left( \Theta_{\mathcal{A}} \right) \right)$$

Theorem 2 shows that the `sail` estimates for nonzero coefficients in the true model have the same asymptotic distribution as they would have if the zero coefficients were known in advance.

# Asymptotic normality

## Theorem 2

$$\widehat{\Theta}_n = \underset{\beta_E, \theta, \gamma}{\operatorname{argmin}} \quad \mathcal{L}(\Theta) + \lambda(1 - \alpha) \left( w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right) + \lambda\alpha \sum_{j=1}^p w_{jE} |\gamma_j|$$

Under certain regularity conditions, the component  $\widehat{\Theta}_{\mathcal{A}}$  of the local minimizer  $\widehat{\Theta}_n$  satisfies

$$\sqrt{n} \left( \widehat{\Theta}_{\mathcal{A}} - \Theta_{\mathcal{A}} \right) \rightarrow_d \mathcal{N} \left( 0, \mathbf{I}^{-1} \left( \Theta_{\mathcal{A}} \right) \right)$$

Theorem 2 shows that the `sail` estimates for nonzero coefficients in the true model have the same asymptotic distribution as they would have if the zero coefficients were known in advance.

Theorem 1 + 2  $\rightarrow$  Oracle property (Fan and Li, 2001)





# Simulation Scenarios

1. **Truth obeys strong hierarchy (right in our wheel house):**

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

# Simulation Scenarios

1. **Truth obeys strong hierarchy (right in our wheel house):**

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

2. **Truth obeys weak hierarchy**

# Simulation Scenarios

1. **Truth obeys strong hierarchy (right in our wheel house):**

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

2. **Truth obeys weak hierarchy**
3. **Truth only has interactions**

# Simulation Scenarios

1. **Truth obeys strong hierarchy (right in our wheel house):**

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

2. **Truth obeys weak hierarchy**
3. **Truth only has interactions**
4. **Truth is linear**

# Simulation Scenarios

1. **Truth obeys strong hierarchy (right in our wheel house):**

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

2. **Truth obeys weak hierarchy**
3. **Truth only has interactions**
4. **Truth is linear**
5. **Truth only has main effects**

# Simulation Scenarios

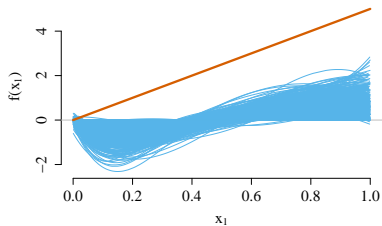
1. **Truth obeys strong hierarchy (right in our wheel house):**

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

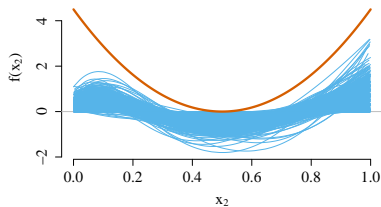
2. **Truth obeys weak hierarchy**
  3. **Truth only has interactions**
  4. **Truth is linear**
  5. **Truth only has main effects**
- $n_{train} = n_{tuning} = 200, n_{test} = 800, p = 1000, \beta_E = 1, SNR = 2$
  - $X_j \sim \text{truncnorm}(0, 1), j = 1, \dots, 1000, E \sim \text{truncnorm}(-1, 1)$
  - sail needs to estimate  $1000 \times 5 \times 2 = 10\text{k}$  parameters

# Scenario 1: Main Effects for 500 Simulations

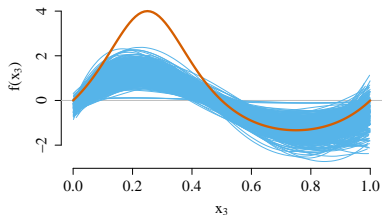
$$f(x_1) = 5x_1$$



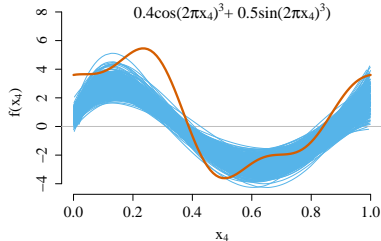
$$f(x_2) = 4.5(2x_2 - 1)^2$$



$$f(x_3) = \frac{4\sin(2\pi x_3)}{2 - \sin(2\pi x_3)}$$



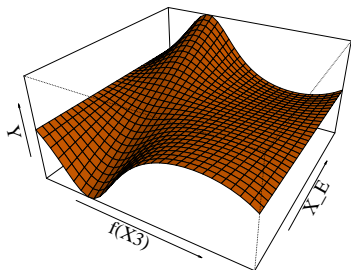
$$f(x_4) = 6(0.1\sin(2\pi x_4) + 0.2\cos(2\pi x_4) + 0.3\sin(2\pi x_4)^2 + 0.4\cos(2\pi x_4)^3 + 0.5\sin(2\pi x_4)^3)$$



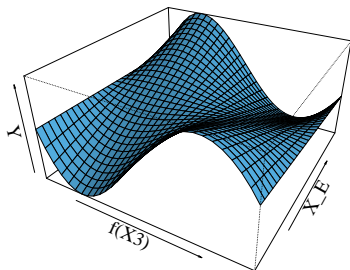


# Scenario 1: Estimated Interaction Effects for $E \cdot f(X_3)$

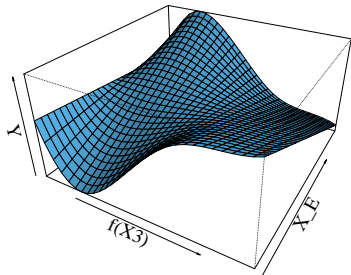
**Truth**



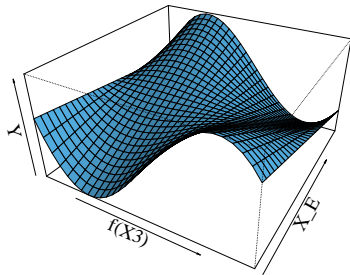
**Estimated: 25th Percentile**



**Estimated: 50th Percentile**

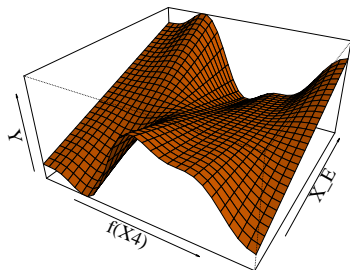


**Estimated: 75th Percentile**

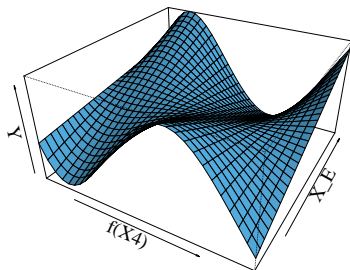


# Scenario 1: Estimated Interaction Effects for $E \cdot f(X_4)$

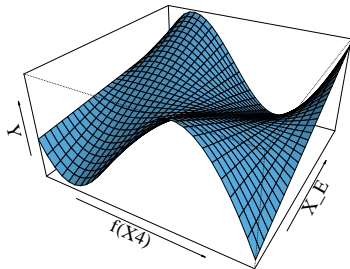
**Truth**



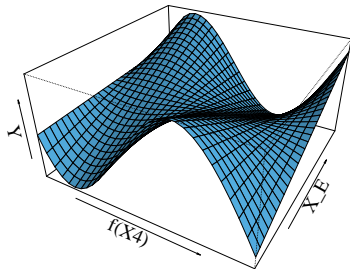
**Estimated: 25th Percentile**



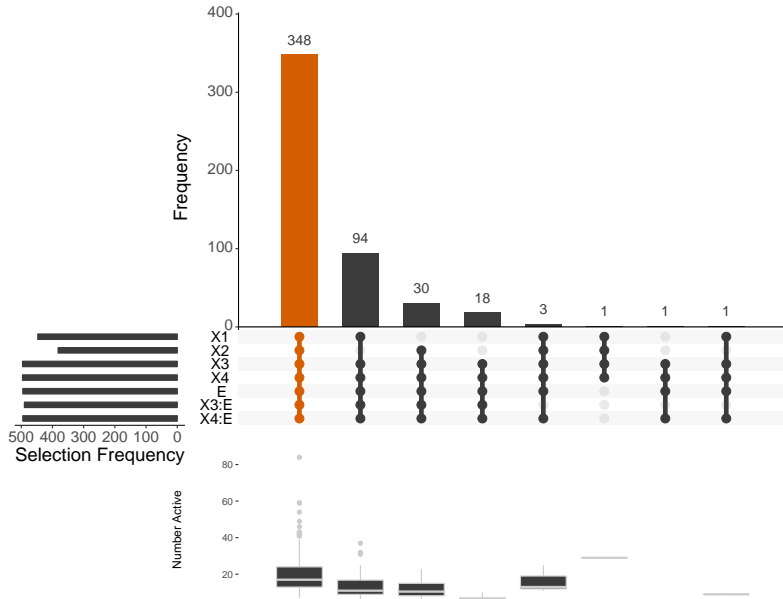
**Estimated: 50th Percentile**



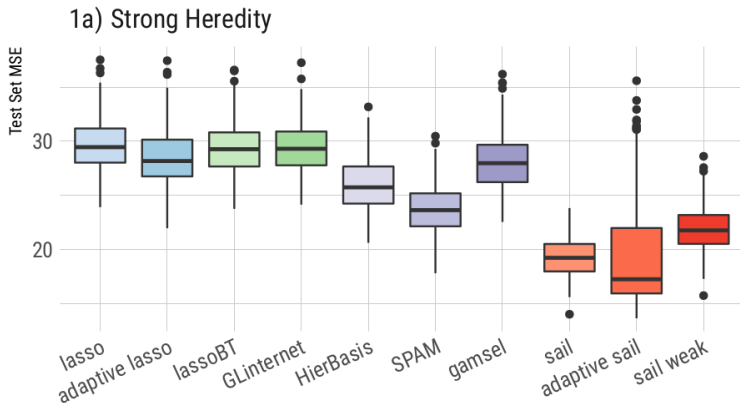
**Estimated: 75th Percentile**



# Right in Our Wheel House Simulation Results

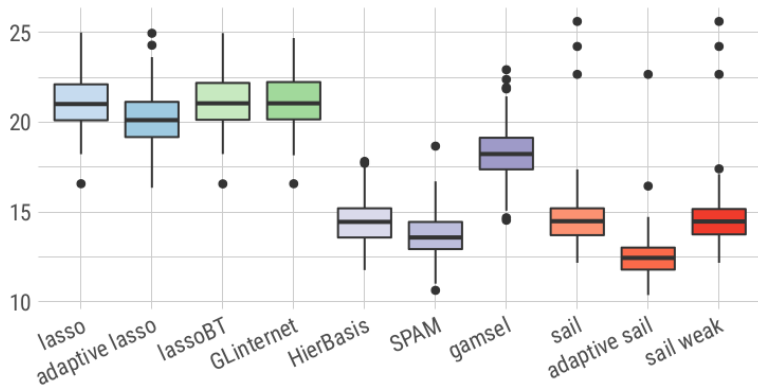


# Strong Heredity



# Main Effects Only

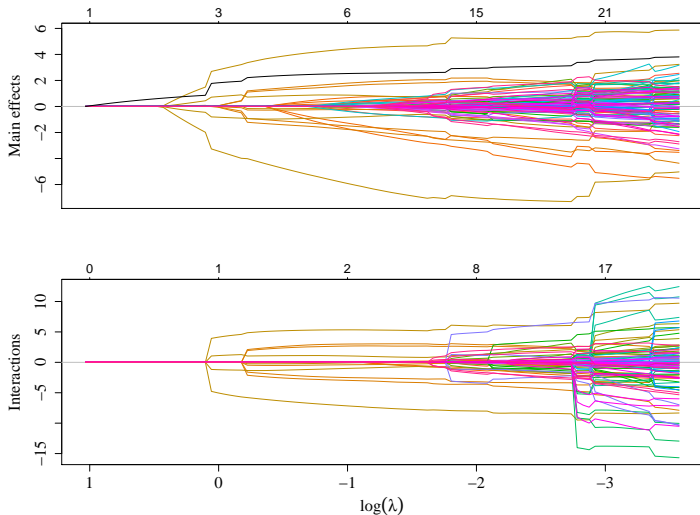
## 3) Main Effects Only





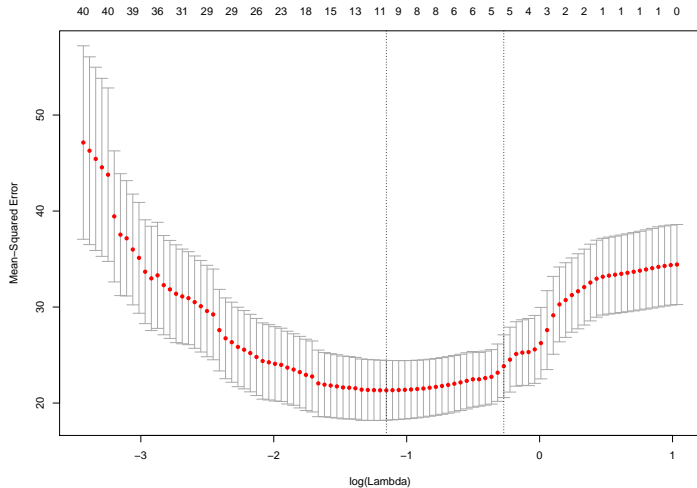
# sail R package: Solution Path results

```
f.basis <- function(x) splines::bs(x, degree = 5)
fit <- sail(x, y, e, basis = f.basis)
plot(fit)
```



# sail R package: Cross-validation results

```
sail::plot(cvfit)
```







# Nurse Family Partnership Program

- Early intervention in young children has been shown to positively impact intellectual abilities.
- Genome-wide association studies (GWAS) suggest that 20% of the variance in educational attainment (years of education) may be accounted for by common genetic variation.
- An interesting query that arises is how the environment interacts with these genetics variants to predict measures of cognitive function.

# Nurse Family Partnership Program

- Early intervention in young children has been shown to positively impact intellectual abilities.
- Genome-wide association studies (GWAS) suggest that 20% of the variance in educational attainment (years of education) may be accounted for by common genetic variation.
- An interesting query that arises is how the environment interacts with these genetics variants to predict measures of cognitive function.

# Nurse Family Partnership Program

- The Stanford Binet IQ scores at 4 years of age were collected for 189 subjects born to women randomly assigned to control ( $n = 100$ ) or nurse-visited intervention groups ( $n = 89$ ).

# Nurse Family Partnership Program

- The Stanford Binet IQ scores at 4 years of age were collected for 189 subjects born to women randomly assigned to control ( $n = 100$ ) or nurse-visited intervention groups ( $n = 89$ ).
- For each subject, we calculated a polygenic risk score (PRS) for educational attainment at different p-value thresholds using weights from a previous GWAS.

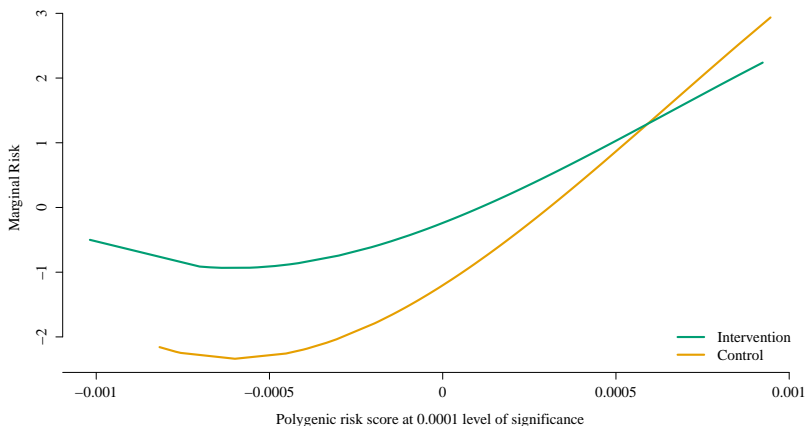
# Nurse Family Partnership Program

- The Stanford Binet IQ scores at 4 years of age were collected for 189 subjects born to women randomly assigned to control ( $n = 100$ ) or nurse-visited intervention groups ( $n = 89$ ).
- For each subject, we calculated a polygenic risk score (PRS) for educational attainment at different p-value thresholds using weights from a previous GWAS.
- In this context, individuals with a higher PRS have a propensity for higher educational attainment.

# Nurse Family Partnership Program

- The Stanford Binet IQ scores at 4 years of age were collected for 189 subjects born to women randomly assigned to control ( $n = 100$ ) or nurse-visited intervention groups ( $n = 89$ ).
- For each subject, we calculated a polygenic risk score (PRS) for educational attainment at different p-value thresholds using weights from a previous GWAS.
- In this context, individuals with a higher PRS have a propensity for higher educational attainment.
- The goal of this analysis was to determine if there was an interaction between genetic predisposition to educational attainment ( $X$ ) and maternal participation in the NFP program ( $E$ ) on child IQ at 4 years of age ( $Y$ ).

## Application of sail to NFP data



**Fig.:** The selected model, chosen via 10-fold cross-validation, contained three variables: the main effects for the intervention and the PRS for educational attainment using genetic variants significant at the 0.0001 level, as well as their interaction.





# Strengths and Limitations

## Strengths

- Non-linear environment interactions with strong heredity property in  $p \gg N$
- `sail` allows for flexible modeling of input variables

# Strengths and Limitations

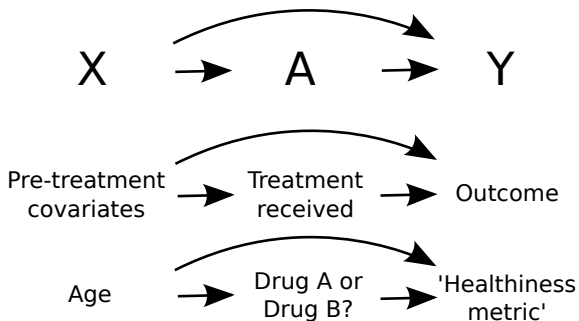
## Strengths

- Non-linear environment interactions with strong heredity property in  $p \gg N$
- `sail` allows for flexible modeling of input variables

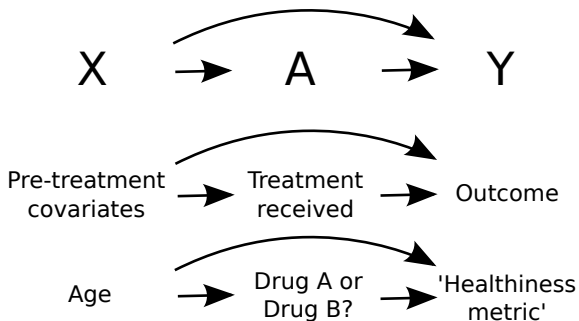
## Limitations

- `sail` can currently only handle  $E \cdot f(X)$  or  $f(E) \cdot X$
- Does not allow for  $f(X_1, E)$  or  $f(X_1, X_2)$
- Memory footprint is an issue

# Dynamic Treatment Regimes (DTRs)



# Dynamic Treatment Regimes (DTRs)



$$\mathbb{E}[Y \mid \mathbf{X}, A; \boldsymbol{\psi}, \boldsymbol{\beta}] = \underbrace{\mathbf{X}\boldsymbol{\beta}}_{\text{Impact of patient history in the absence of treatment}} + \underbrace{\psi_0 A + \boldsymbol{\psi} \mathbf{A} \mathbf{X}}_{\text{Impact of treatment on outcome}}$$

# Extension of sail to DTRs



Cornell University

arXiv.org > stat > arXiv:2101.07359

Statistics > Methodology

[Submitted on 18 Jan 2021]

## Variable Selection in Regression-based Estimation of Dynamic Treatment Regimes

Zeyu Bian, Erica EM Moodie, Susan M Shortreed, Sahir Bhatnagar

Dynamic treatment regimes (DTRs) consist of a sequence of decision rules, one per stage of intervention, that finds effective treatments for individual patients between treatment and a small number of covariates which are often chosen a priori. However, with increasingly large and complex data being collected, a driven approach of selecting these covariates might improve the estimated decision rules and simplify models to make them easier to interpret. We propose a method that has the strong heredity property, that is, an interaction term can be included in the model only if the corresponding main terms have also been selected. The newly proposed methods compare favorably with other variable selection approaches.

Subjects: **Methodology (stat.ME)**; Computation (stat.CO)

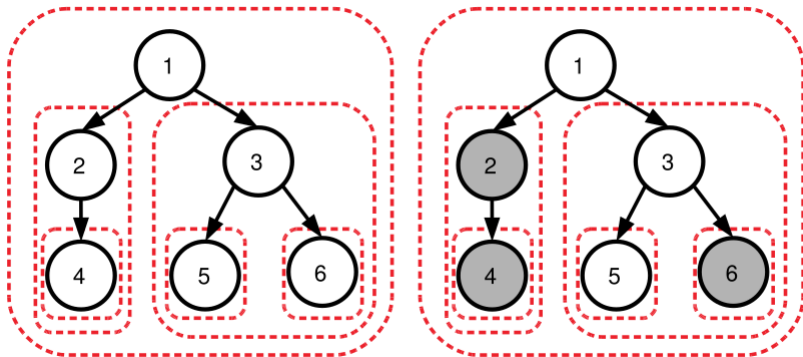
Cite as: [arXiv:2101.07359](https://arxiv.org/abs/2101.07359) [**stat.ME**]

(or [arXiv:2101.07359v1](https://arxiv.org/abs/2101.07359v1) [**stat.ME**] for this version)

---

<sup>1</sup>*In revision at Biometrics.* <https://arxiv.org/abs/2101.07359>

# Hierarchical Penalty Structure



<sup>1</sup>Bach, Jenatton, Mairal and Obozinski (2011). Optimization with Sparsity-Inducing Penalties.

# Bi-level selection

- Bi-level selection:

$$f(X_1) = \underbrace{\begin{bmatrix} X_{11} & \psi_{11}(X_{11}) & \psi_{12}(X_{12}) & \cdots & \psi_{11}(X_{15}) \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ X_{i1} & \psi_{11}(X_{i1}) & \psi_{12}(X_{i2}) & \cdots & \psi_{11}(X_{i5}) \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ X_{N1} & \psi_{11}(X_{N1}) & \psi_{12}(X_{N2}) & \cdots & \psi_{11}(X_{N5}) \end{bmatrix}}_{\Psi_1} \quad N \times 5 \quad \times \quad \underbrace{\begin{bmatrix} \beta_{\text{linear}} \\ \beta_{11} \\ \beta_{12} \\ \beta_{13} \\ \beta_{14} \\ \beta_{15} \end{bmatrix}}_{\theta_1} \quad 6 \times 1$$





# Acknowledgements



Zeyu Bian, PhD (c)

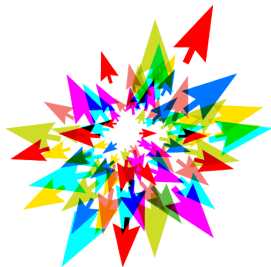


# Acknowledgements

- Tianyuan Lu (McGill)
- Yi Yang (McGill)
- Celia Greenwood (Lady Davis Institute)
- Erica Moodie (McGill)
- Kieran O'Donnell (Yale)



**compute** | **calcul**  
canada | canada



# References

1. Bhatnagar, SR, Lu, T, Lovato, A, Olds, DL, Kobor, MS, Meaney, MJ, O'Donnell, K, Yang, Y, and Greenwood, CMT (2021+). A Sparse Additive Model for High-Dimensional Interactions with an Exposure Variable. bioRxiv. DOI [10.1101/445304](https://doi.org/10.1101/445304). *In revision at Computational Statistics and Data Analysis*.
2. **Bian Z**, Moodie EEM, Shortreed S, Bhatnagar SR (2021+). Variable Selection in Regression-based Estimation of Dynamic Treatment Regimes. <https://arxiv.org/abs/2101.07359>. *In revision at Biometrics*.
3. De Leeuw, J. (1994). Block-relaxation algorithms in statistics. In Information systems and data analysis (pp. 308-324). Springer Berlin Heidelberg.
4. Choi, N. H., Li, W., & Zhu, J. (2010). Variable selection with the strong heredity constraint and its oracle property. *Journal of the American Statistical Association*, 105(489), 354-364.
5. Chipman, H. (1996). Bayesian variable selection with related predictors. *Canadian Journal of Statistics*, 24(1), 17-36.

[sahirbhatnagar.com](http://sahirbhatnagar.com)

# Session Info

```
R version 3.6.0 (2019-04-26)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Pop!_OS 18.10

Matrix products: default
BLAS: /usr/lib/x86_64-linux-gnu/blas/libblas.so.3.8.0
LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.8.0

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods   base

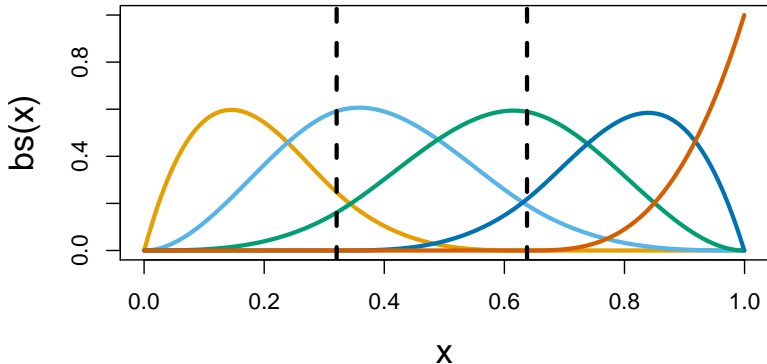
other attached packages:
[1] xtable_1.8-4      rpart.plot_3.0.6  rpart_4.1-15     data.table_1.12.0
[5] ISLR_1.2          ggplot2_3.1.0    knitr_1.23

loaded via a namespace (and not attached):
 [1] Rcpp_1.0.1          magrittr_1.5      tidyselect_0.2.5
 [4] munsell_0.5.0      colorspace_1.4-0  R6_2.4.0
 [7] rlang_0.3.4        highr_0.8         stringr_1.4.0
[10] plyr_1.8.4         dplyr_0.8.0.1    tools_3.6.0
[13] grid_3.6.0         gtable_0.2.0     xfun_0.7
[16] pacman_0.5.0       withr_2.1.2      digest_0.6.19
[19] lazyeval_0.2.1    assertthat_0.2.1 tibble_2.1.1
[22] crayon_1.3.4      RSkittleBrewer_1.1 purrr_0.3.2
[25] glue_1.3.1        evaluate_0.14     stringi_1.4.3
[28] compiler_3.6.0    pillar_1.4.0     scales_1.0.0
[31] pkgconfig_2.0.2
```

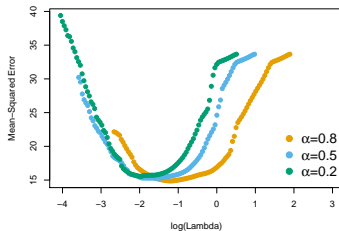
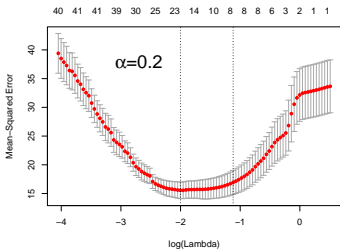
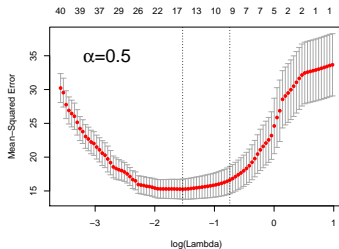
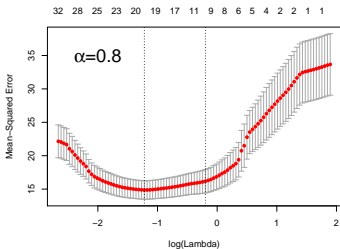
# B-Spline Expansion

```
x <- truncnorm::rtruncnorm(1000, a = 0, b = 1)
B <- splines::bs(x, df = 5, degree=3, intercept = FALSE)
```

**df=5, degree=3, inner.knots at c(33.33%, 66.66%) percentile**



# sail A Note on the Second Tuning Parameter results

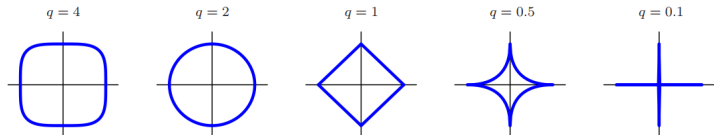


# Why the L1 norm ?

- For a fixed real number  $q \geq 0$  consider the criterion

$$\tilde{\beta} = \arg \min_{\beta} \left\{ \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|^q \right\}$$

- Why do we use the  $\ell_1$  norm? Why not use the  $q = 2$  (Ridge) or any  $\ell_q$  norm?



- $q = 1$  is the smallest value that yields a sparse solution **and** yields a **convex** problem  $\rightarrow$  scalable to high-dimensional data
- For  $q < 1$  the constrained region is **nonconvex**



# Linear Effects Simulation - Comparison

