

Betting on Sparsity

Sahir Bhatnagar, PhD Candidate, McGill Biostatistics

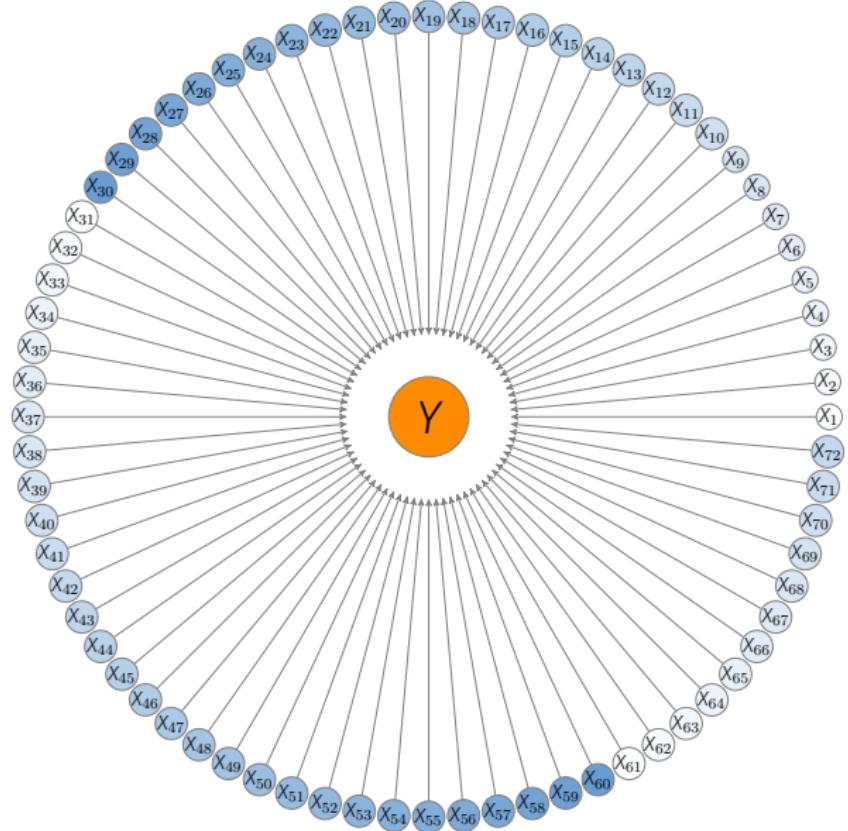
Joint work with Yi Yang and Celia Greenwood (McGill)

March 1, 2018

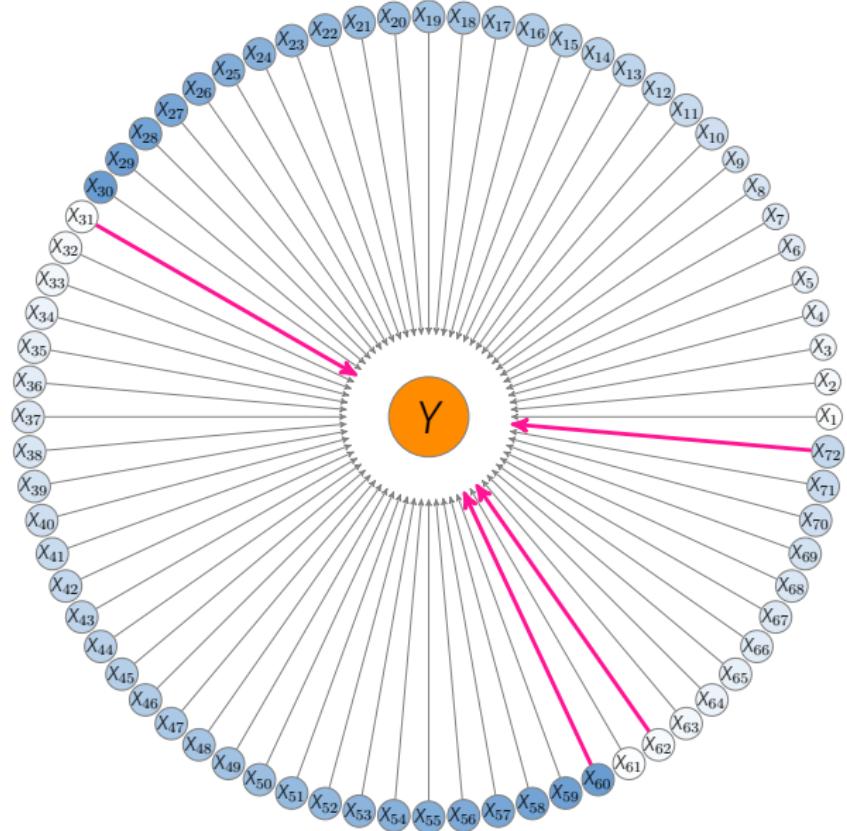


Betting on Sparsity

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.¹

¹The elements of statistical learning. Springer series in statistics, 2001.

Bet on Sparsity Principle

Use a procedure that does well in sparse problems,
since no procedure does well in dense problems.¹

- 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.

¹The elements of statistical learning. Springer series in statistics, 2001.

A Thought Experiment

How would you schedule a meeting of 20 people?

How would you schedule a meeting of 20 people?

March 2017												
11 participants	Thu 9	Fri 10	Sat 11	Sun 12	Mon 13	Tue 14	Wed 15	Thu 16	Fri 17	Sat 18	Sun 19	
JayZ	✓	✓	✓			✓			✓	✓		
Evan									✓	✓	✓	
Omar	✓	✓		✓		✓			✓	✓	✓	
Caitlin	✓	✓	✓						✓	✓	✓	
Austin	✓	✓	✓									
Ethan			✓	✓					✓		✓	
Max	✓	✓	✓			✓			✓	✓	✓	
Tycho	✓	✓	✓	✓		✓			✓	✓	✓	
Janavi Chadha	✓		✓	✓		✓	✓			✓	✓	
Charlotte										✓		
Darshanye	✓	✓				✓			✓	✓		
Your name	□	□	□	□	□	□	□	□	□	□	□	
	5:00 PM – 9:00 PM	5:00 PM – 9:00 PM	9:00 AM – 3:00 PM	3:00 PM – 9:00 PM	1:00 PM – 9:00 PM							
March 2017												
	7	8	7	4	0	6	1	0	7	8	9	2

Doctors Bet on Sparsity Also

Doctors Bet on Sparsity Also



Motivating Example

Predictors of NHL Salary²



²<https://www.kaggle.com/camnugent/nhl-salary-data-prediction-cleaning-and-modelling>

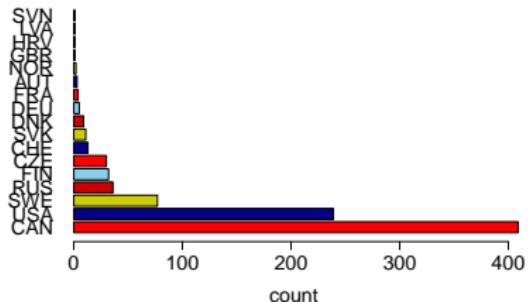
Supervised Learning

- Learn the function f

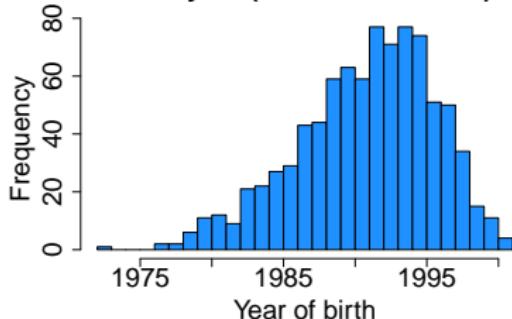


Predictors of NHL Salary

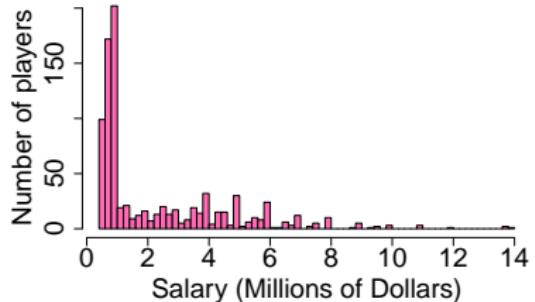
Country



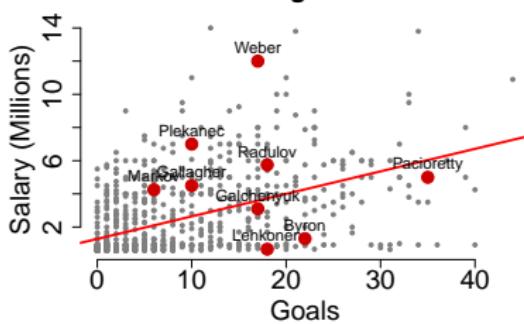
Birth year (2016/2017 season)



NHL Salary Distribution: 2016/2017

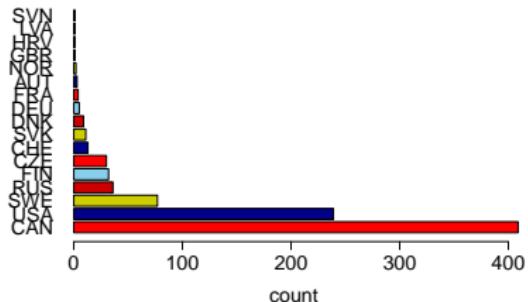


Linear Regression Fit

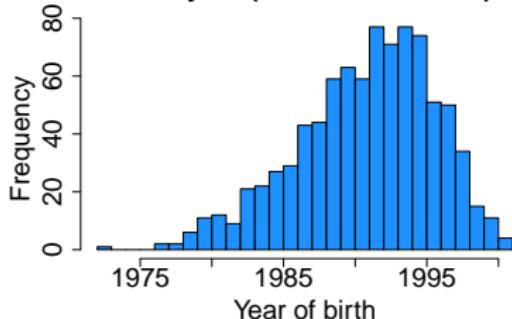


Predictors of NHL Salary

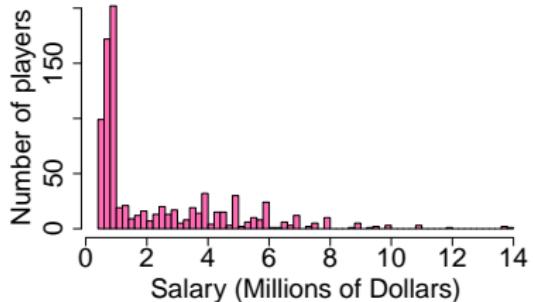
Country



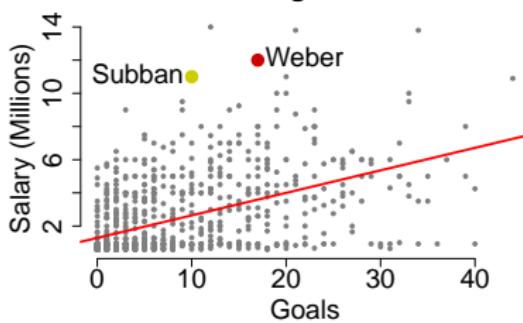
Birth year (2016/2017 season)



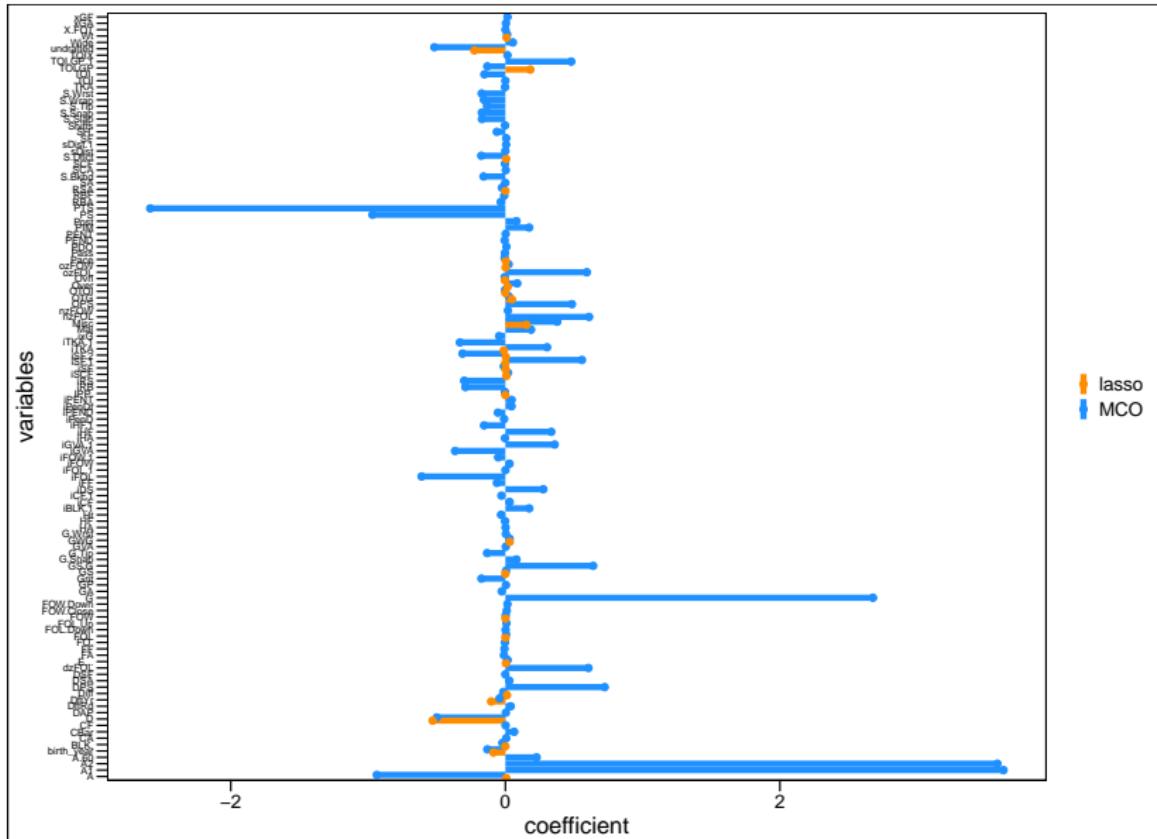
NHL Salary Distribution: 2016/2017



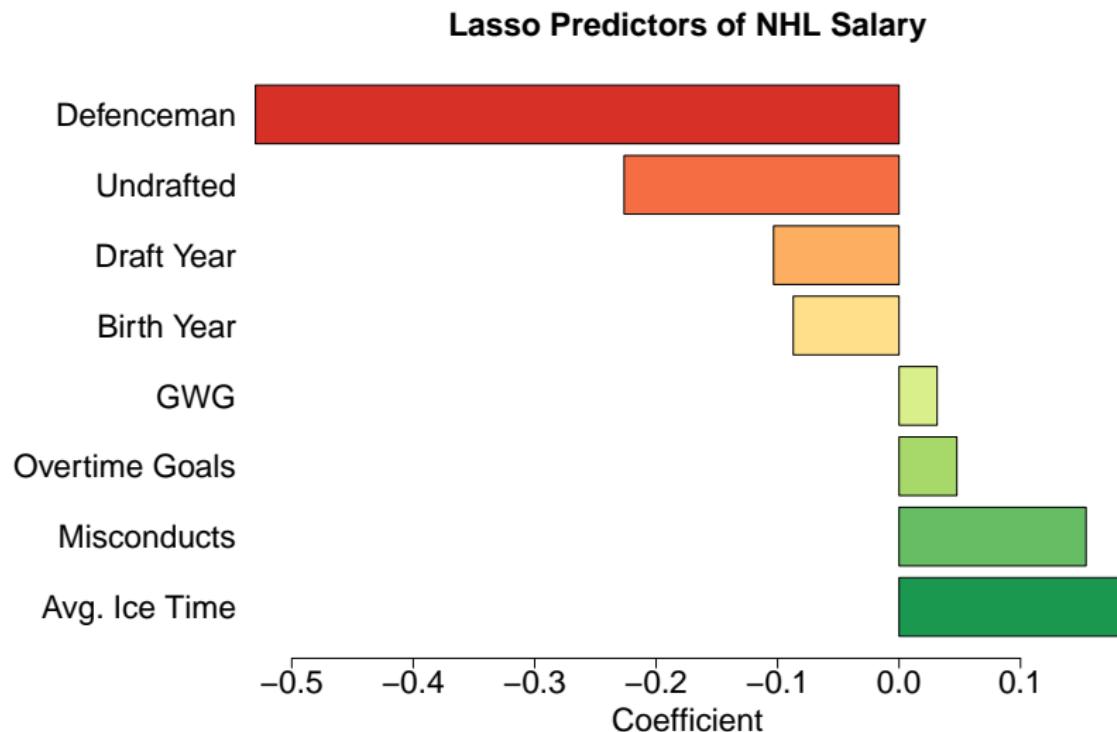
Linear Regression Fit



OLS vs. Lasso Coefficients



Lasso Selected Predictors



Background on the Lasso

- Predictors $x_{ij}, j = 1, \dots, p$ and outcome values y_i for the i th observation, $i = 1, \dots, n$
- Assume x_{ij} are standardized so that $\sum_i x_{ij}/n = 0$ and $\sum_i x_{ij}^2 = 1$.

¹Tibshirani. JRSSB (1996)

Background on the Lasso

- Predictors $x_{ij}, j = 1, \dots, p$ and outcome values y_i for the i th observation, $i = 1, \dots, n$
- Assume x_{ij} are standardized so that $\sum_i x_{ij}/n = 0$ and $\sum_i x_{ij}^2 = 1$. The lasso¹ solves

$$\hat{\boldsymbol{\beta}}^{lasso} = \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

$$\text{subject to } \sum_{j=1}^p |\beta_j| \leq s, \quad s > 0$$

¹Tibshirani. JRSSB (1996)

Background on the Lasso

- Predictors $x_{ij}, j = 1, \dots, p$ and outcome values y_i for the i th observation, $i = 1, \dots, n$
- Assume x_{ij} are standardized so that $\sum_i x_{ij}/n = 0$ and $\sum_i x_{ij}^2 = 1$. The lasso¹ solves

$$\hat{\boldsymbol{\beta}}^{lasso} = \arg \min_{\boldsymbol{\beta}} \frac{1}{2} \sum_{i=1}^n \left(y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

$$\text{subject to } \sum_{j=1}^p |\beta_j| \leq s, \quad s > 0$$

- Equivalently, the Lagrange version of the problem, for $\lambda > 0$

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

¹Tibshirani. JRSSB (1996)

Inspection of the Lasso Solution

- Consider a single predictor setting based on the observed data $\{(x_i, y_i)\}_{i=1}^n$. The problem then is to solve

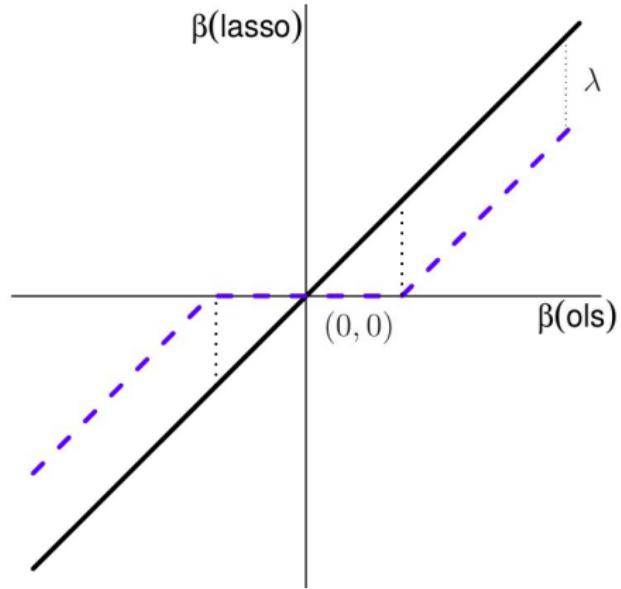
$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \frac{1}{2} \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda |\beta| \quad (1)$$

- With a **standardized** predictor, the lasso solution (1) is a **soft-thresholded** version of the least-squares (LS) estimate $\hat{\beta}^{\text{LS}}$

$$\begin{aligned}\hat{\beta}^{\text{lasso}} &= s_{\lambda}(\hat{\beta}^{\text{LS}}) = \text{sign}(\hat{\beta}^{\text{LS}}) \left(|\hat{\beta}^{\text{LS}}| - \lambda \right)_+ \\ &= \begin{cases} \hat{\beta}^{\text{LS}} - \lambda, & \hat{\beta}^{\text{LS}} > \lambda \\ 0 & |\hat{\beta}^{\text{LS}}| \leq \lambda \\ \hat{\beta}^{\text{LS}} + \lambda & \hat{\beta}^{\text{LS}} \leq -\lambda \end{cases}\end{aligned}$$

Inspection of the Lasso Solution

- When the data are standardized, the lasso solution shrinks the LS estimate toward zero by the amount λ



¹Hastie et al. Statistical learning with sparsity: the lasso and generalizations

Choosing the Model Complexity

Group Lasso Illustration

Extended from the lasso penalty, the group lasso estimator is:

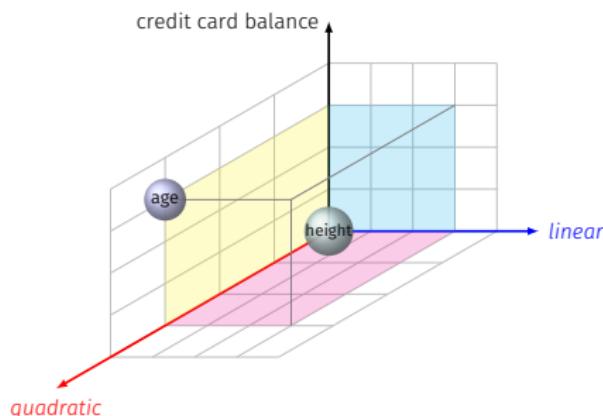
$$\min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\mathbf{y} - \beta_0 - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{k=1}^K \sqrt{p_k} \|\boldsymbol{\beta}^{(k)}\|_2 \quad p_k - \text{group size}$$

Group Lasso Illustration

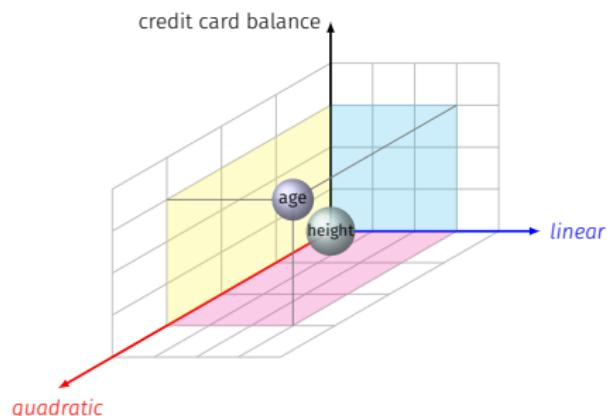
Extended from the lasso penalty, the group lasso estimator is:

$$\min_{(\beta_0, \boldsymbol{\beta})} \frac{1}{2} \|\mathbf{y} - \beta_0 - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \sum_{k=1}^K \sqrt{p_k} \|\boldsymbol{\beta}^{(k)}\|_2 \quad p_k - \text{group size}$$

Credit card balance \sim age + age² + height + height²



(c) Lasso



(d) Group Lasso

Our Software

Overview of Our Software Packages

- **eclust** – Bhatnagar et al. (2017, Genetic Epidemiology)
<https://cran.r-project.org/package=eclust>
- **sail** – Bhatnagar, Yang and Greenwood (2018+, preprint)
<https://github.com/sahirbhatnagar/sail>
- **gmmix** – Bhatnagar, Oualkacha, Yang, Greenwood (2018+, preprint)
<https://github.com/sahirbhatnagar/gmmix>
- **casebase** – Bhatnagar¹, Turgeon¹, Yang, Hanley and Saarela (2018+, preprint)
<https://cran.r-project.org/package=casebase>

¹joint co-authors

Overview of Our Software Packages

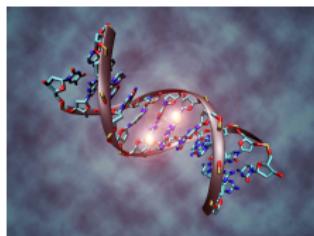
	eclust	sail	gmmix	casebase
Model				
Least-Squares	✓	✓	✓	
Binary Classification	✓			
Survival Analysis				✓
Penalty				
Ridge	✓		✓	✓
Lasso	✓	✓	✓	✓
Elastic Net	✓		✓	✓
Group Lasso		✓	✓	
Feature				
Interactions	✓	✓		✓
Flexible Modeling	✓	✓		✓
Random Effects			✓	
Data	(x, y, e)	(x, y, e)	(x, y, Ψ)	(x, t, δ)

sail: Strong Additive Interaction Learning

Motivation 1: Non-linear Interactions



~



×

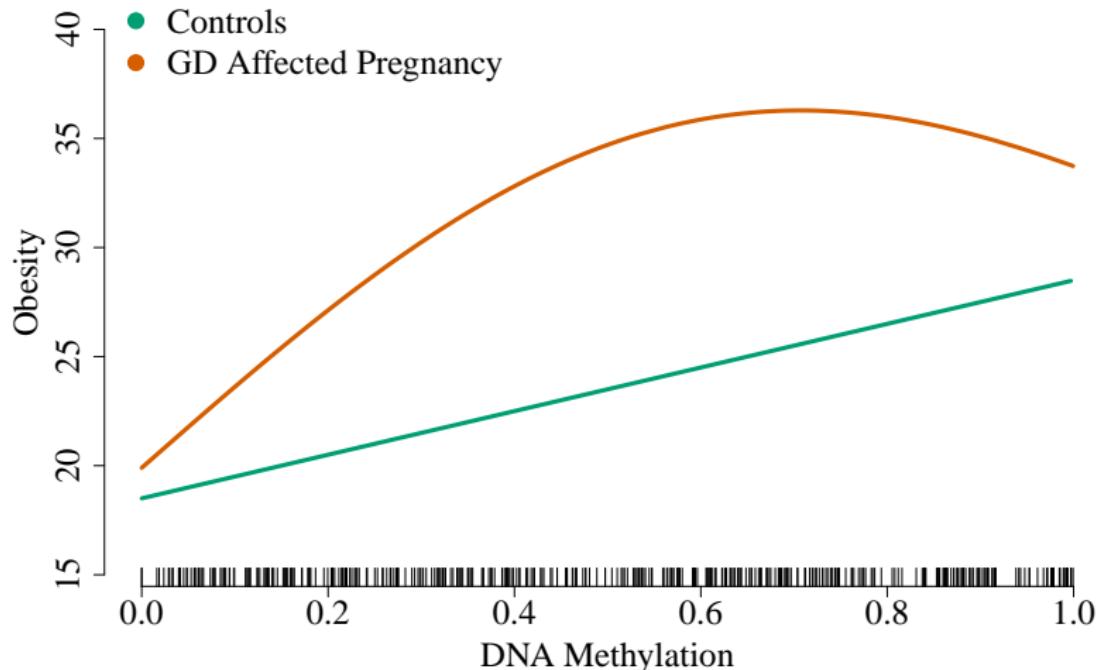


Phenotype
Obesity measures

Large Data
Child's epigenome
($p \approx 450k$)

Environment
Gestational
Diabetes

Motivation 1: Non-linear Interactions



Motivation 2: Heredity Property

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

¹Chipman. Canadian Journal of Statistics (1996)

²McCullagh and Nelder. Generalized Linear Models (1983)

³Cox. International Statistical Review (1984)

Motivation 2: Heredity Property

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

Strong Heredity¹

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

- Heredity property is desired for the purposes of **interpretability**²
- Large main effects are more likely to lead to appreciable interactions³

¹Chipman. Canadian Journal of Statistics (1996)

²McCullagh and Nelder. Generalized Linear Models (1983)

³Cox. International Statistical Review (1984)

Lasso interaction model

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

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

$$\operatorname{argmin}_{\beta_0, \boldsymbol{\beta}, \boldsymbol{\alpha}} \mathcal{L}(Y; \boldsymbol{\Theta}) + \lambda(\|\boldsymbol{\beta}\|_1 + \|\boldsymbol{\alpha}\|_1)$$

Strong Heredity Interactions: Current State of the Art

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

Our Extension to Nonlinear Effects

Consider the basis expansion

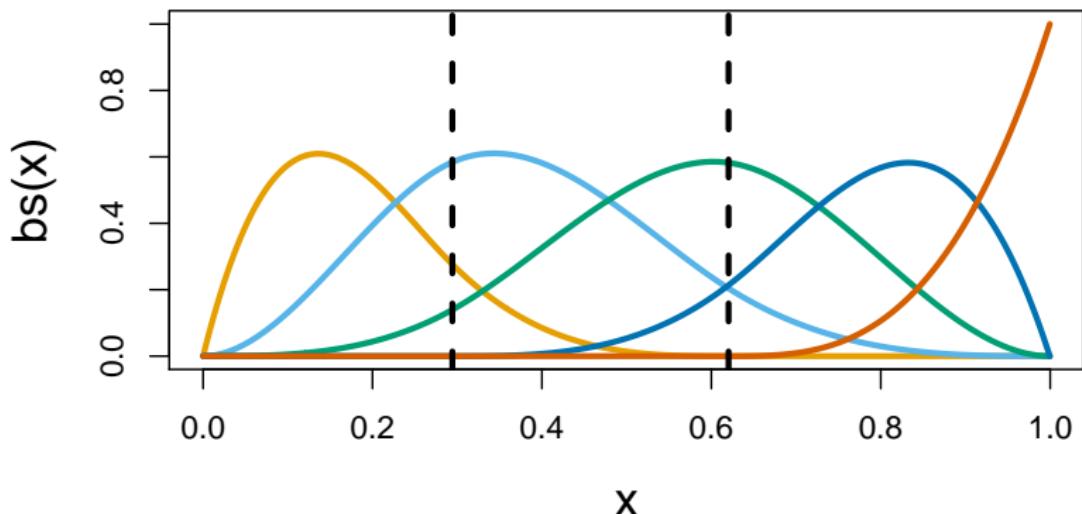
$$f_j(X_j) = \sum_{\ell=1}^{p_j} \psi_{j\ell}(X_j) \beta_{j\ell}$$

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

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: Additive Interactions

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

Model

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

sail: Strong Heredity

Reparametrization¹

$$\alpha_j = \gamma_j \beta_E \theta_j$$

Model

$$Y = \beta_0 \cdot \mathbf{1} + \sum_{j=1}^p \boldsymbol{\Psi}_j \theta_j + \beta_E X_E + \sum_{j=1}^p \color{red}{\gamma_j \beta_E X_E \boldsymbol{\Psi}_j \theta_j} + \varepsilon$$

Objective Function

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

¹Choi et al. JASA (2010)

Algorithm

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 λ **do**

repeat

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

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

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

$$k \leftarrow k + 1$$

until convergence criterion is satisfied;

end

Implementation

Objective Function

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

¹<https://github.com/sahirbhatnagar/sail>

Implementation

Objective Function

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

Lasso problem

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

¹<https://github.com/sahirbhatnagar/sail>

Implementation

Objective Function

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

¹<https://github.com/sahirbhatnagar/sail>

Implementation

Objective Function

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

Group Lasso problem

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

¹<https://github.com/sahirbhatnagar/sail>

Simulations

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 only has main effects:

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

Simulation Scenarios

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

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

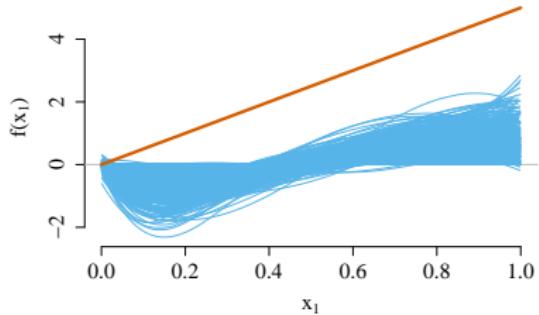
2. Truth only has main effects:

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

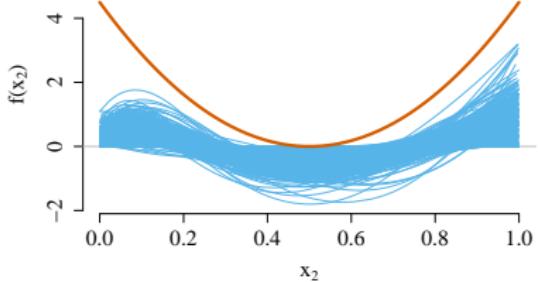
- $n = 200, p = 1000, \beta_E = 1, SNR = 2$
- $X_j \sim \text{truncnorm}(0, 1), j = 1, \dots, 1000,$
 $E \sim \text{truncnorm}(-1, 1)$
- sail needs to estimate $1000 \times 5 \times 2 = 10k$ parameters

Scenario 1: Main Effects for 500 Simulations

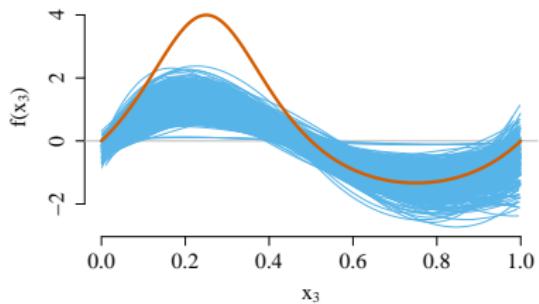
$$f(x_1) = 5x_1$$



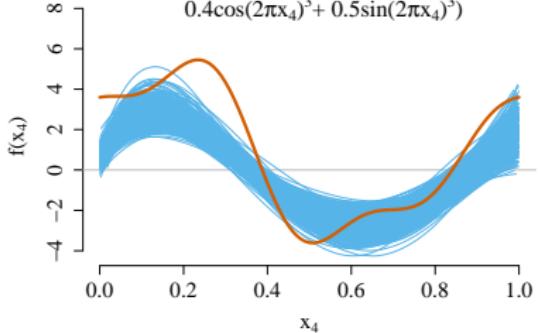
$$f(x_2) = 4.5(2x_2 - 1)^2$$



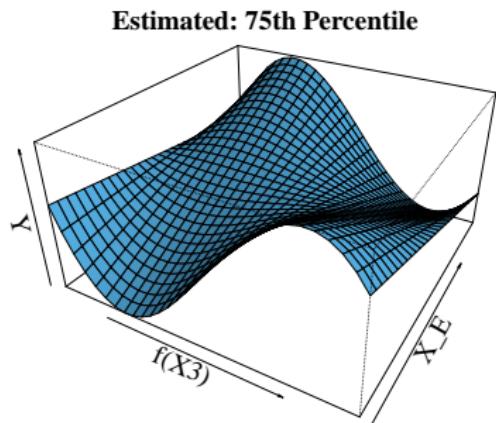
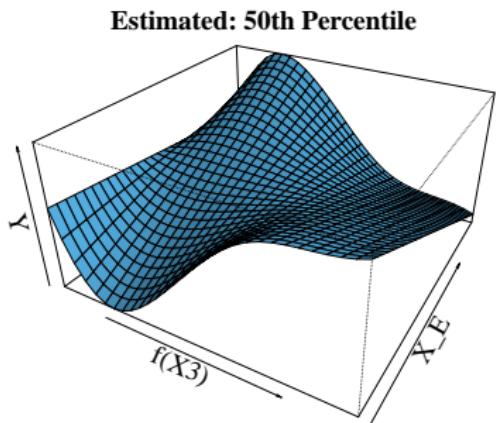
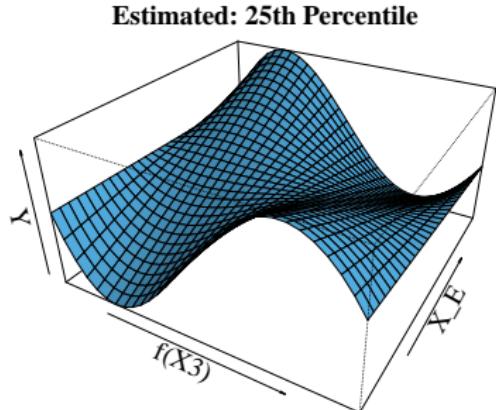
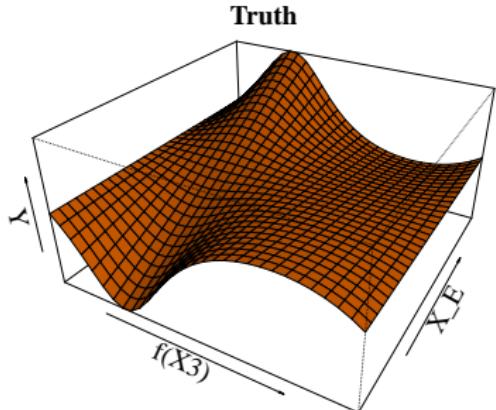
$$f(x_3) = \frac{4\sin(2\pi x_3)}{2 - \sin(2\pi x_3)}$$



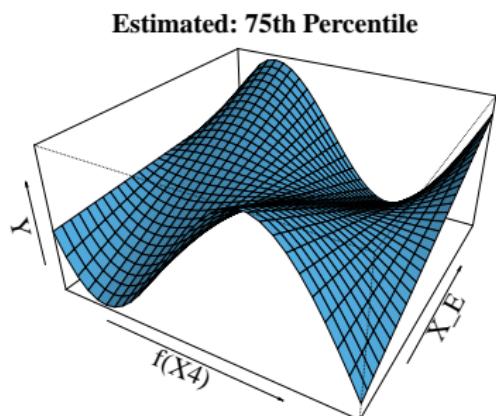
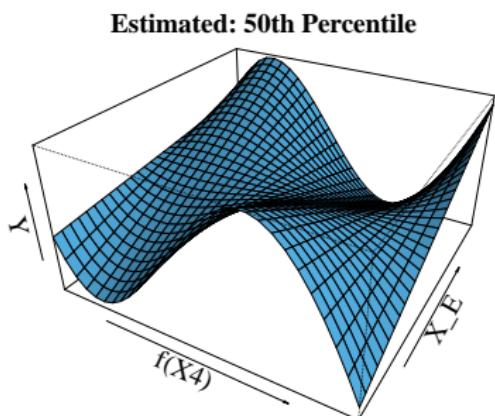
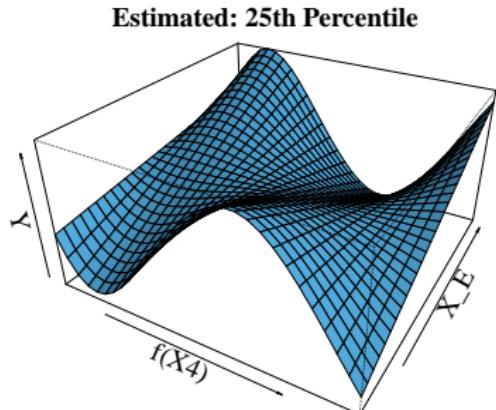
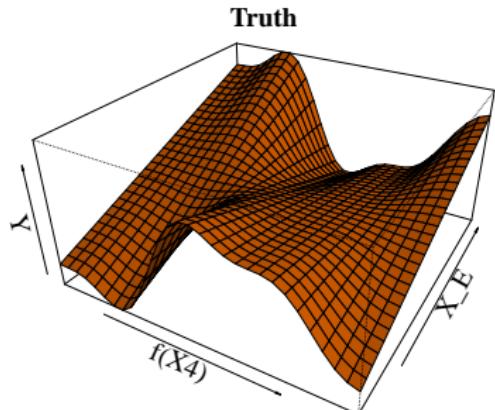
$$f(x_4) = 6(0.1\sin(2\pi x_4) + 0.2\cos(2\pi x_4) + 0.3\sin(2\pi x_4)^2 + 0.4\cos(2\pi x_4)^3 + 0.5\sin(2\pi x_4)^3)$$



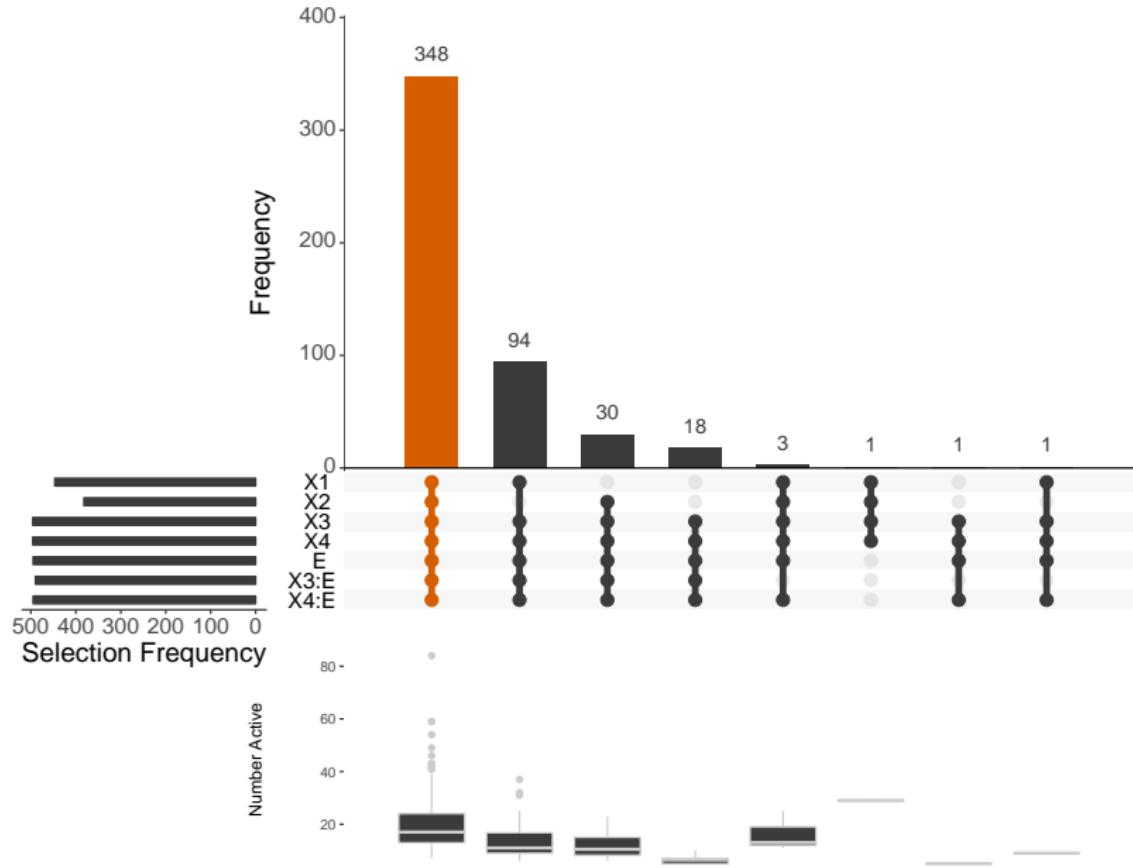
Scenario 1: Estimated Interaction Effects for $E \cdot f(X_3)$



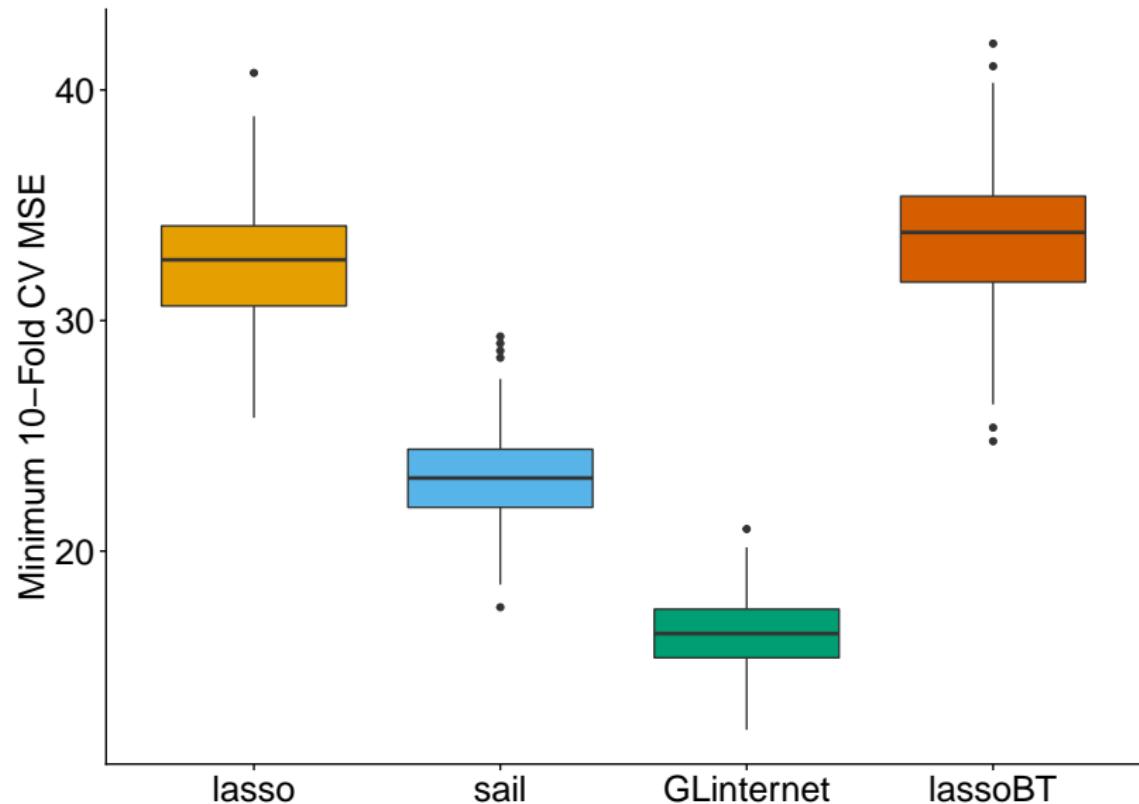
Scenario 1: Estimated Interaction Effects for $E \cdot f(X_4)$



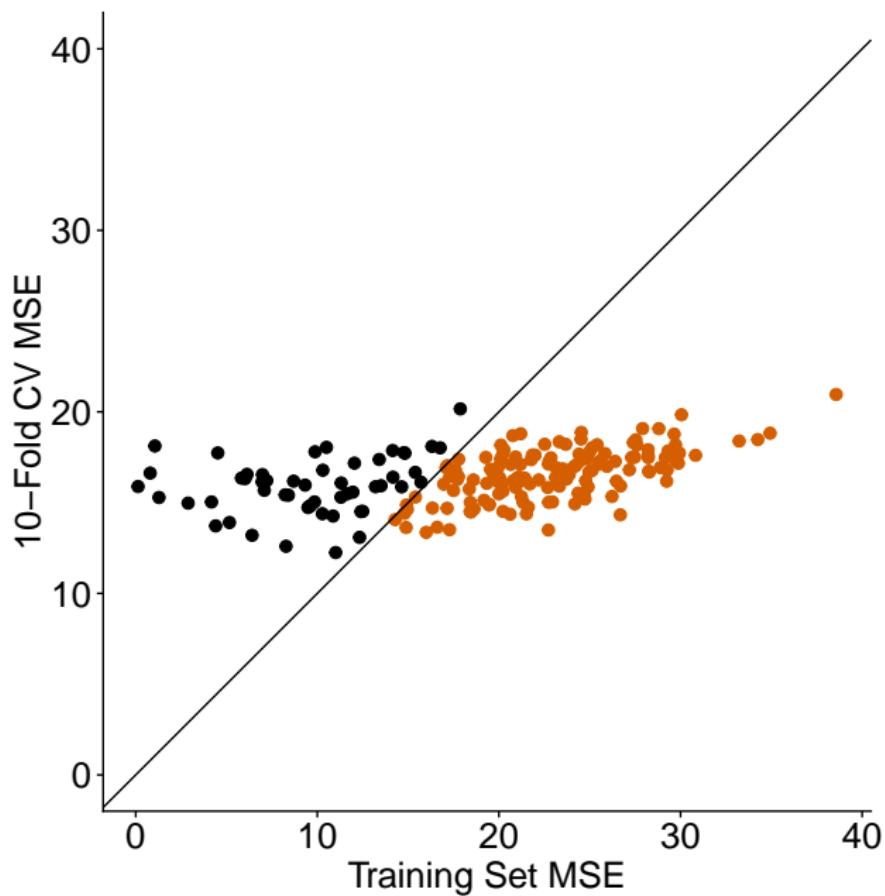
Right in Our Wheel House Simulation Results



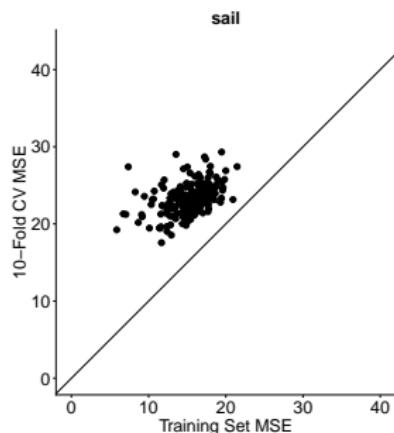
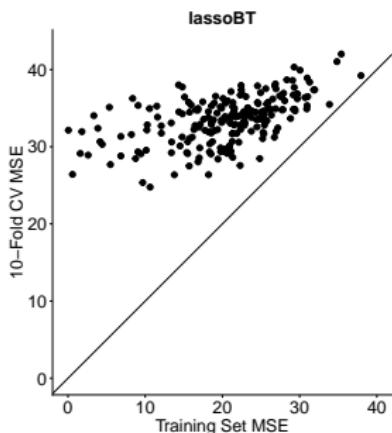
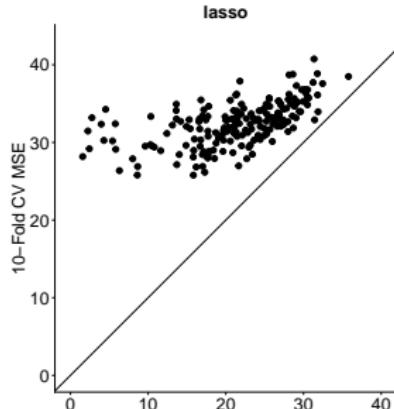
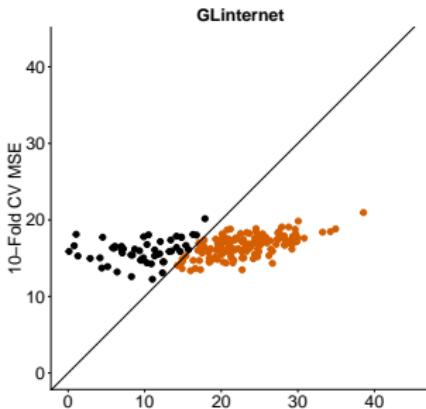
Right in Our Wheel House Simulation - Comparison



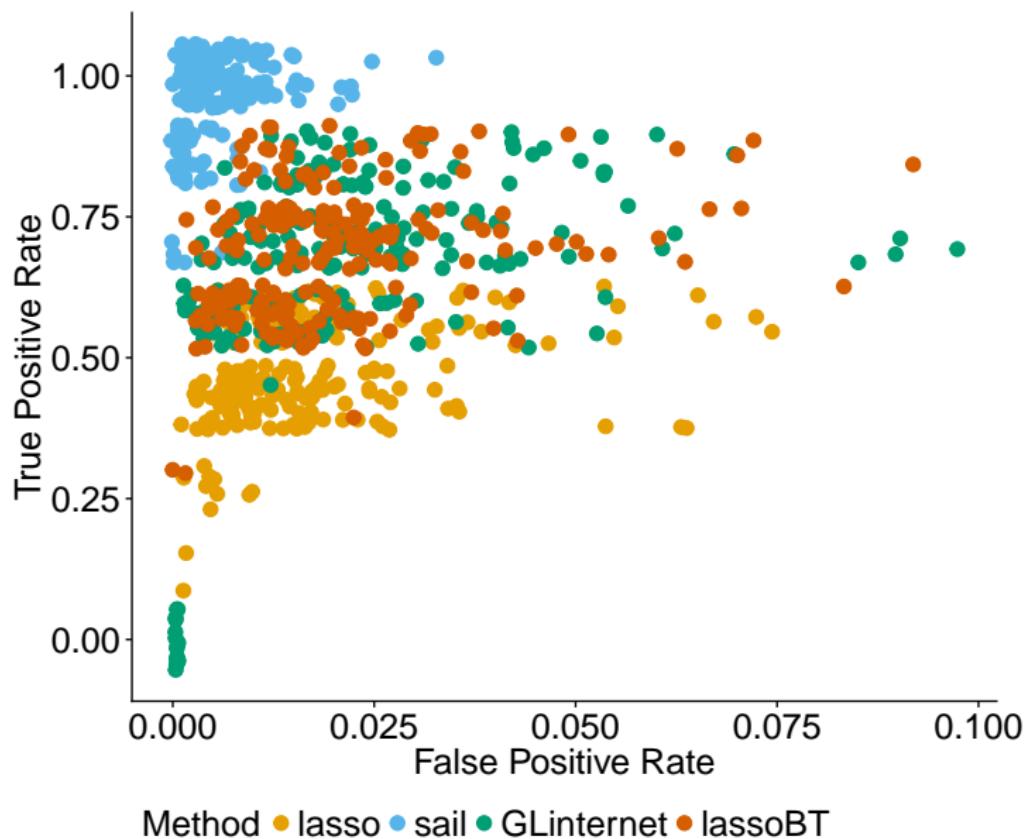
GLinternet: 70% of points below the line



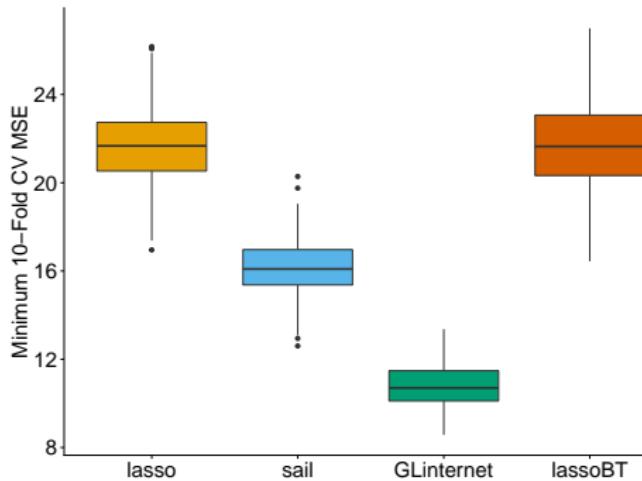
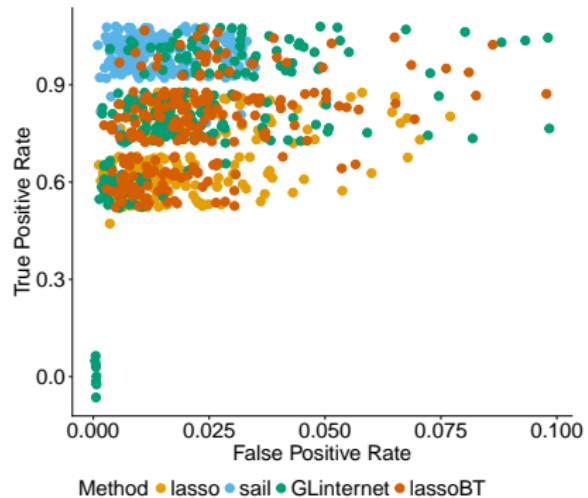
10-Fold CV MSE vs. Training MSE Comparison



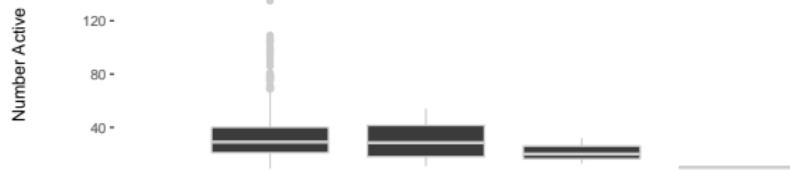
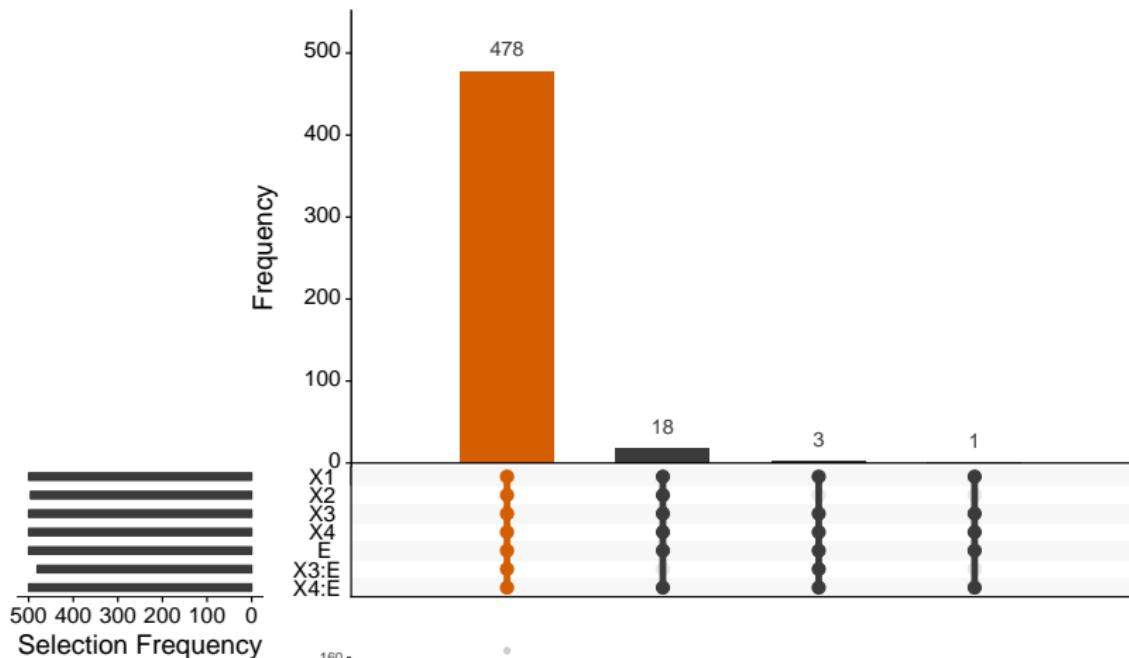
Right in Our Wheel House Simulation - Comparison



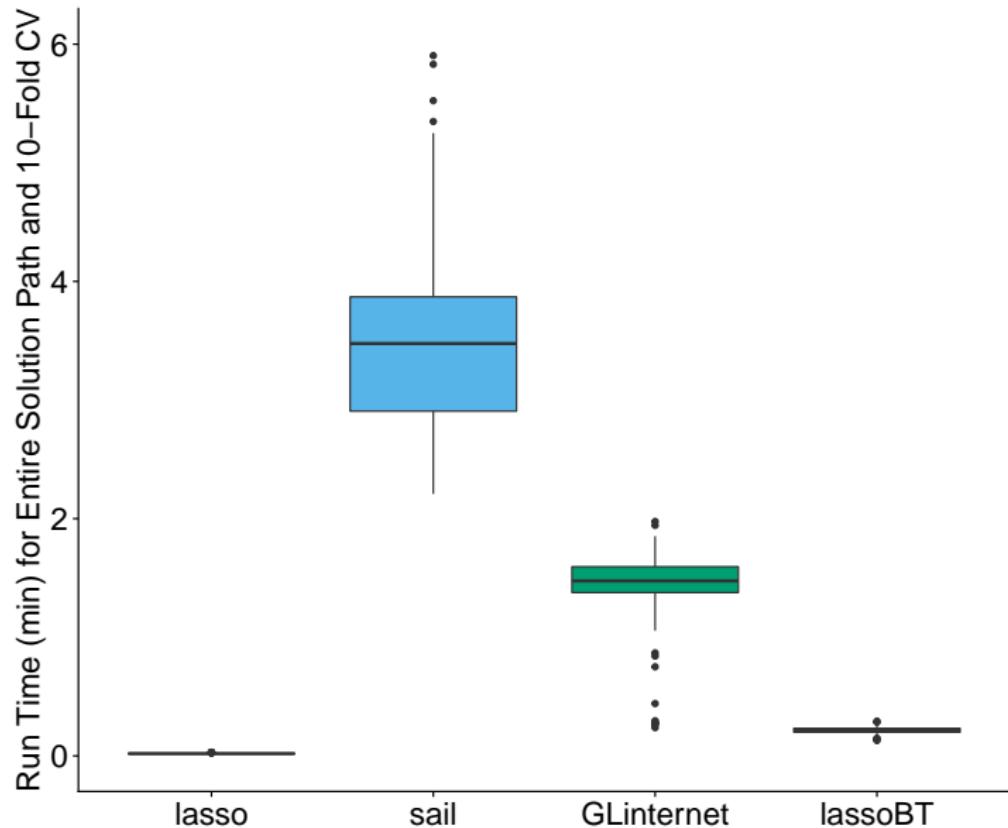
No Interactions Simulation - Comparison



sail with degree=1 when Truth is Linear



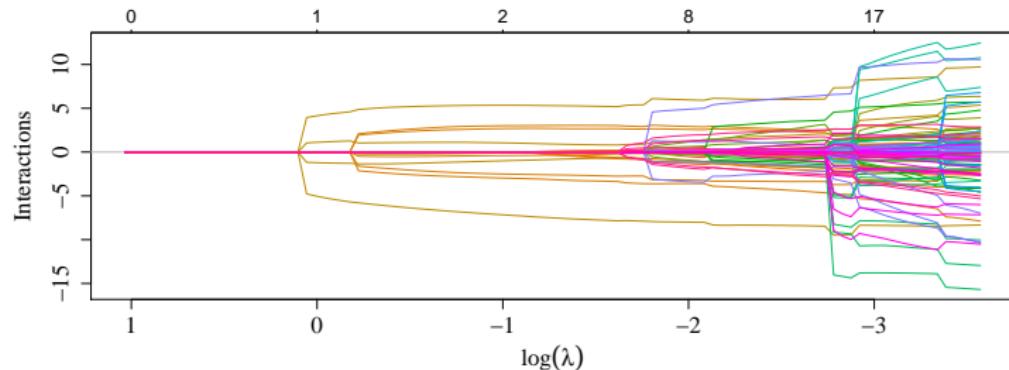
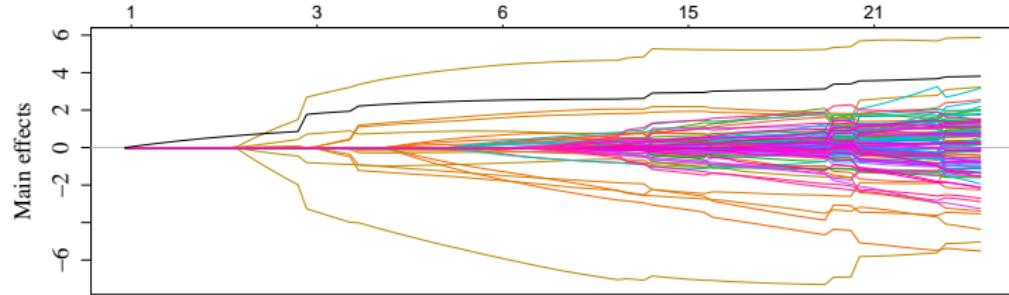
Computing time



sail R package

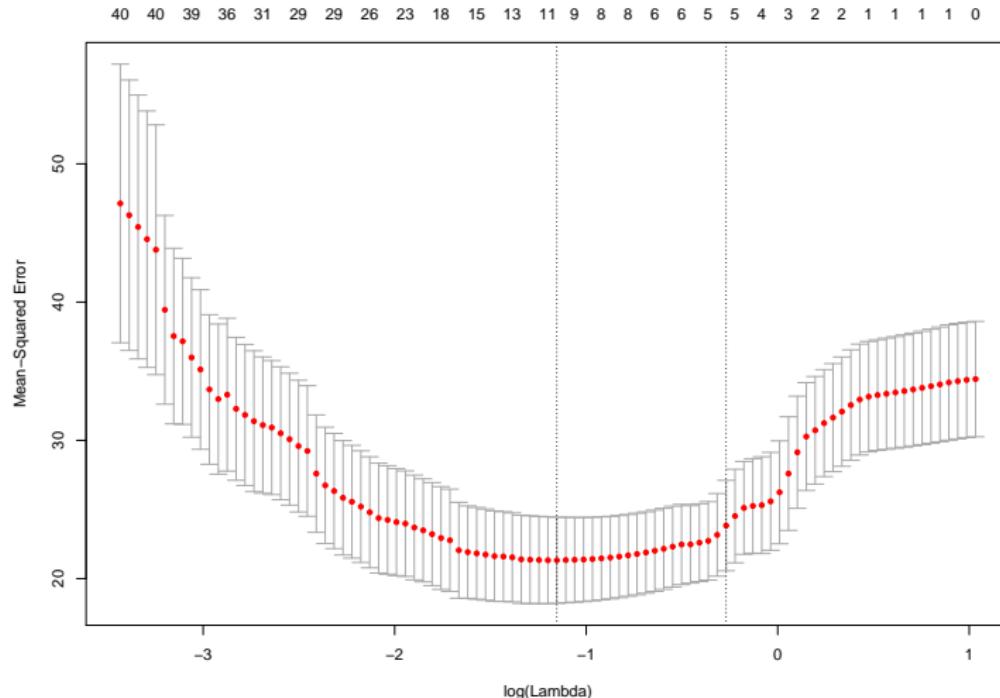
sail R package: Solution Path results

```
sail:::plot(fit)
```

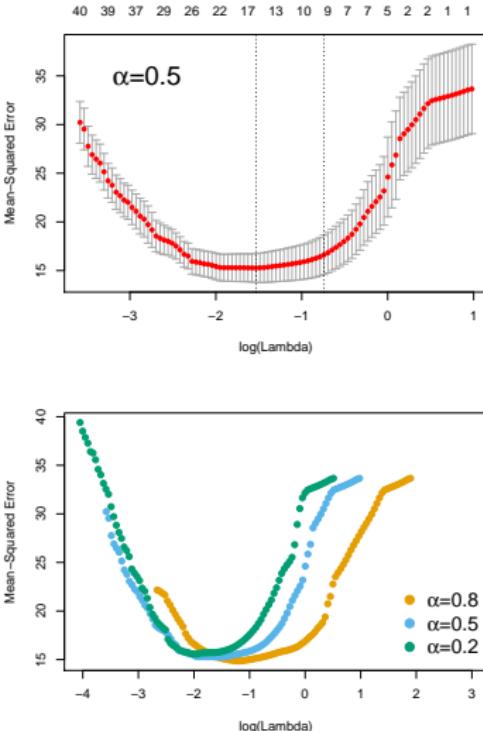
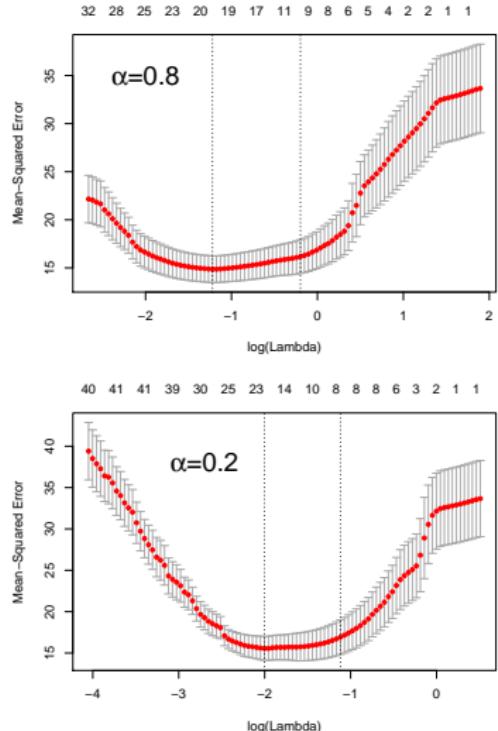


sail R package: Cross-validation results

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



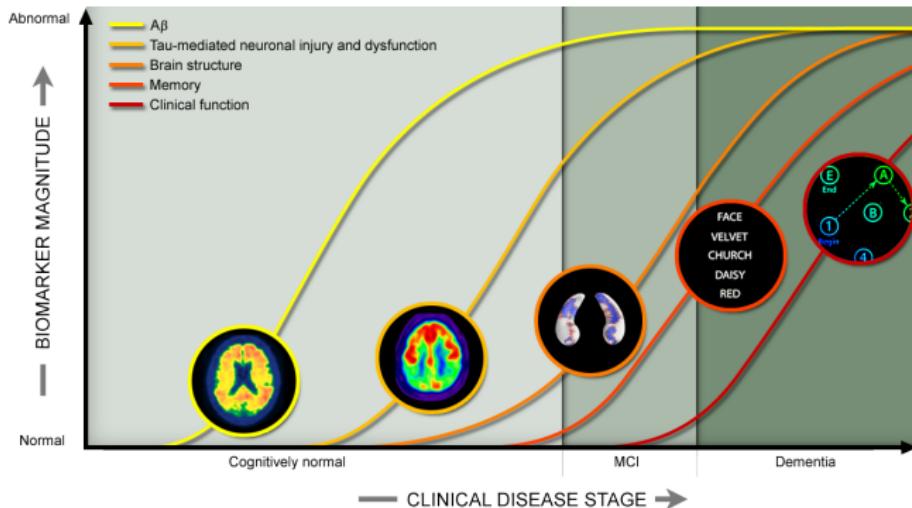
sail A Note on the Second Tuning Parameter results



Real Data Application

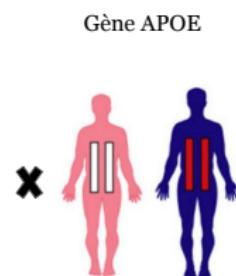
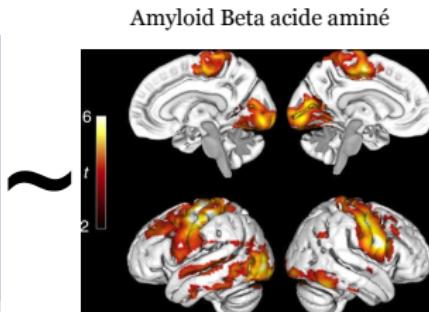
Alzheimer's Disease Neuroimaging Initiative (ADNI)

- 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 ADNI is to **validate biomarkers** for use in Alzheimer's disease clinical treatment trials



Interaction between A β Protein and APOE gene

- **E:** APOE4 allele increases the risk for Alzheimer's and lowers the age of onset
- **X:** PET amyloid imaging to assess A β protein load in 96 brain regions
- **Y:** General cognitive decline measured by mini-mental state examination
- $3 \times 96 \times 2 + 1 = 577$ parameters to estimate



Y
 343×1

X
 343×96

E
 343×1

Variable Selection Results: sail vs. lasso

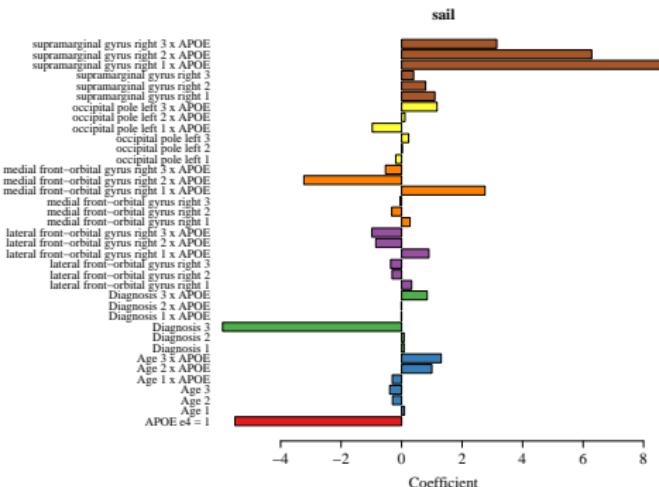


Fig.: sail: 7 variables

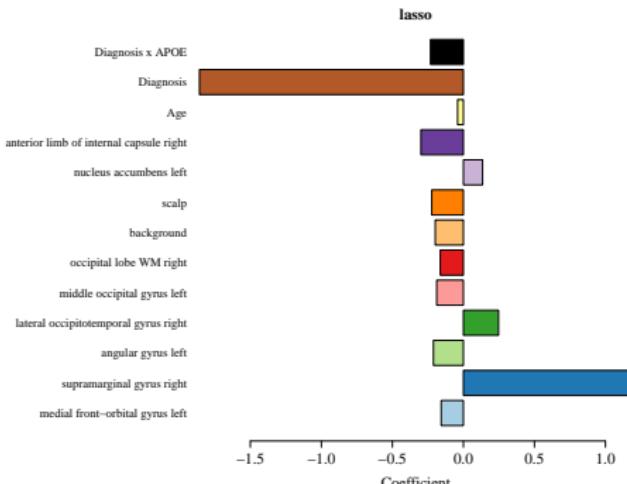
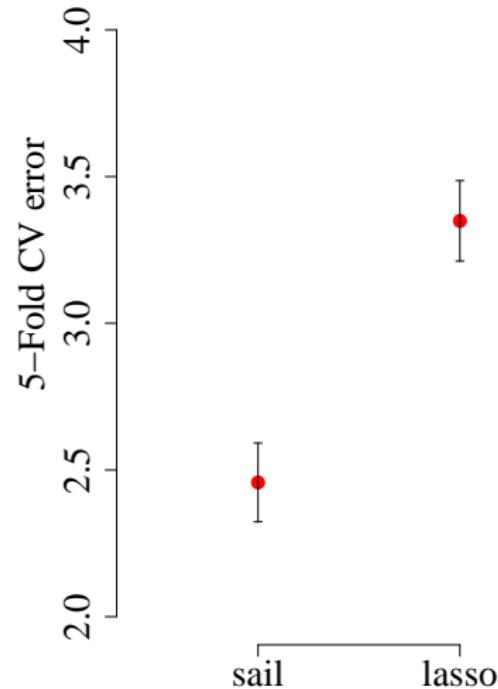


Fig.: lasso: 13 variables

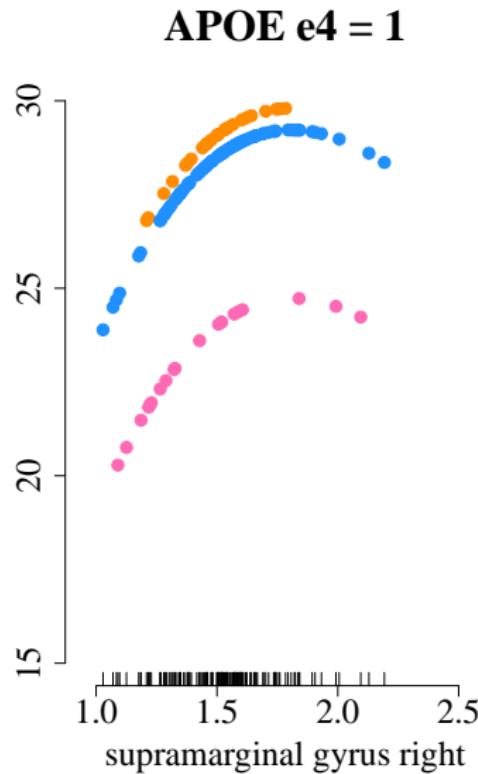
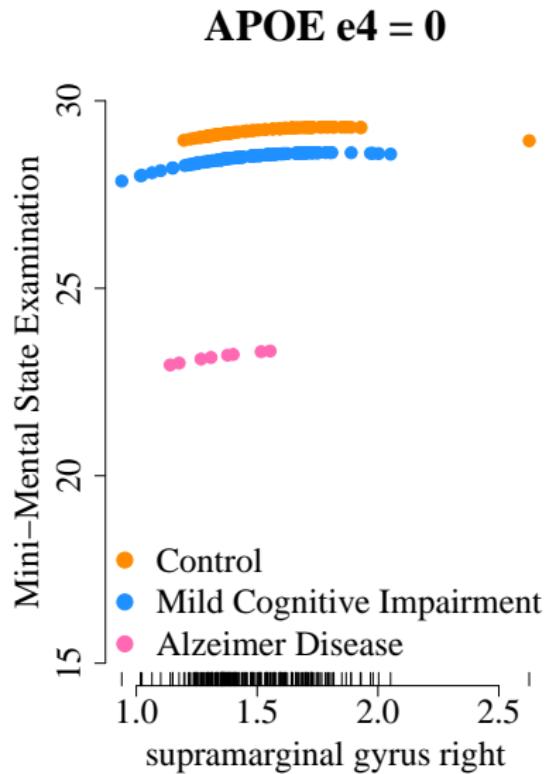
5-Fold Cross-Validated MSE

$$R^2 = 0.61$$

$$R^2 = 0.44$$



sail: Interactions with the supramarginal gyrus region



Discussion

Strengths and Limitations

Strengths

- Non-linear environment interactions with strong heredity property in $p >> N$
- `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

Future Directions

- Weak heredity property $\rightarrow \alpha_j = \gamma_j(|\beta_j| + |\beta_E|)$
- Implement ADMM algorithm for scalability. Distributed computing (GPU)
- Binary Outcomes
- bi-level selection:

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

Acknowledgements



References

- Radchenko, P., & James, G. M. (2010). Variable selection using adaptive nonlinear interaction structures in high dimensions. *Journal of the American Statistical Association*, 105(492), 1541-1553.
- 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.
- Chipman, H. (1996). Bayesian variable selection with related predictors. *Canadian Journal of Statistics*, 24(1), 17-36.
- Friedman, J., Hastie, T., & Tibshirani, R. (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software*, 33(1)
- Yang, Y., & Zou, H. (2015). A fast unified algorithm for solving group-lasso penalize learning problems. *Statistics and Computing*, 25(6), 1129-1141
- De Leeuw, J. (1994). Block-relaxation algorithms in statistics. In *Information systems and data analysis* (pp. 308-324). Springer Berlin Heidelberg.

sahirbhatnagar.com

Session Info

```
R version 3.4.1 (2017-06-30)
Platform: x86_64-pc-linux-gnu (64-bit)
Running under: Ubuntu 16.04.3 LTS
```

```
Matrix products: default
BLAS: /usr/lib/openblas-base/libblas.so.3
LAPACK: /usr/lib/libopenblas-p0.2.18.so
```

```
attached base packages:
```

```
[1] stats      graphics   grDevices utils      datasets  base
```

```
other attached packages:
```

```
[1] xtable_1.8-2          rpart.plot_2.1.2       rpart_4.1-11
[4] data.table_1.10.4-3    ISLR_1.2             ggplot2_2.2.1.9000
[7] knitr_1.19
```

```
loaded via a namespace (and not attached):
```

```
[1] Rcpp_0.12.15        magrittr_1.5       splines_3.4.1
[4] munsell_0.4.3       colorspace_1.3-2    rlang_0.1.6
[7] stringr_1.2.0       highr_0.6          plyr_1.8.4
[10] tools_3.4.1         grid_3.4.1         gtable_0.2.0
[13] pacman_0.4.6        lazyeval_0.2.1      digest_0.6.15
[16] tibble_1.4.2         RSkittleBrewer_1.1  codetools_0.2-15
[19] evaluate_0.10.1     stringi_1.1.5      compiler_3.4.1
[22] pillar_1.1.0         methods_3.4.1      scales_0.5.0.9000
[25] truncnorm_1.0-7
```

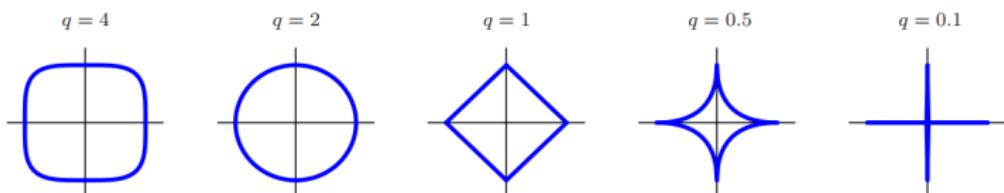
Appendix

Why the L1 norm ?

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

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

- Why do we use the ℓ_1 norm? Why not use the $q = 2$ (Ridge) or any ℓ_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

