

BETTING ON SPARSITY

Sahir Bhatnagar, PhD Candidate, McGill Biostatistics

Joint work with Karim Oualkacha, Yi Yang and Celia Greenwood

November 16, 2017

INTRODUCTION

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)
- 2009: Aon Hewitt - Régimes de retraite

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)
- 2009: Aon Hewitt - Régimes de retraite
- 2010: Associé de la Société des Actuaires (ASA)

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)
- 2009: Aon Hewitt - Régimes de retraite
- 2010: Associé de la Société des Actuaires (ASA)
- 2012: Maîtrise en Biostatistique (Queen's)

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)
- 2009: Aon Hewitt - Régimes de retraite
- 2010: Associé de la Société des Actuaires (ASA)
- 2012: Maîtrise en Biostatistique (Queen's)
- 2013: Doctorat en Biostatistique (McGill, diplôme prévu mai 2018)

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)
- 2009: Aon Hewitt - Régimes de retraite
- 2010: Associé de la Société des Actuaires (ASA)
- 2012: Maîtrise en Biostatistique (Queen's)
- 2013: Doctorat en Biostatistique (McGill, diplôme prévu mai 2018)
- 2016: Wellcome Trust Sanger Institute (Cambridge)

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)
- 2009: Aon Hewitt - Régimes de retraite
- 2010: Associé de la Société des Actuaires (ASA)
- 2012: Maîtrise en Biostatistique (Queen's)
- 2013: Doctorat en Biostatistique (McGill, diplôme prévu mai 2018)
- 2016: Wellcome Trust Sanger Institute (Cambridge)
- 2017: Chargé de cours théorie des probabilités et l'inférence statistique (McGill)

MON CHEMINEMENT

- 2008: Baccalauréat en actuariat (Concordia)
- 2009: Aon Hewitt - Régimes de retraite
- 2010: Associé de la Société des Actuaires (ASA)
- 2012: Maîtrise en Biostatistique (Queen's)
- 2013: Doctorat en Biostatistique (McGill, diplôme prévu mai 2018)
- 2016: Wellcome Trust Sanger Institute (Cambridge)
- 2017: Chargé de cours théorie des probabilités et l'inférence statistique (McGill)
- sahirbhatnagar.com

OUTLINE

1. A motivating example

OUTLINE

1. A motivating example
2. Background on penalization methods lasso and group lasso

OUTLINE

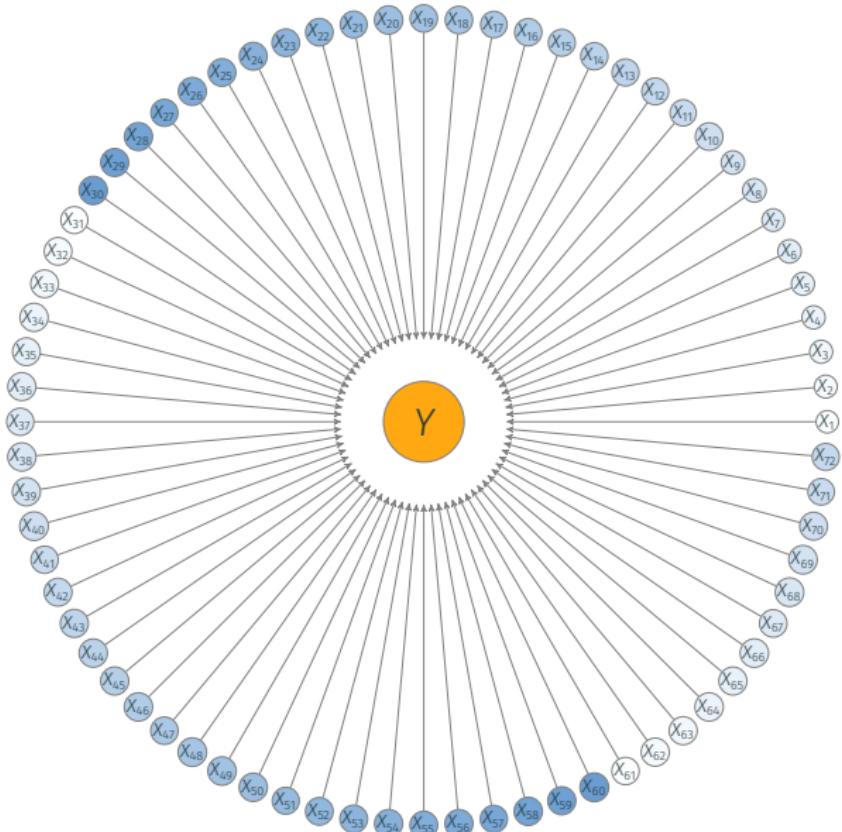
1. A motivating example
2. Background on penalization methods lasso and group lasso
3. Overview of our software packages

OUTLINE

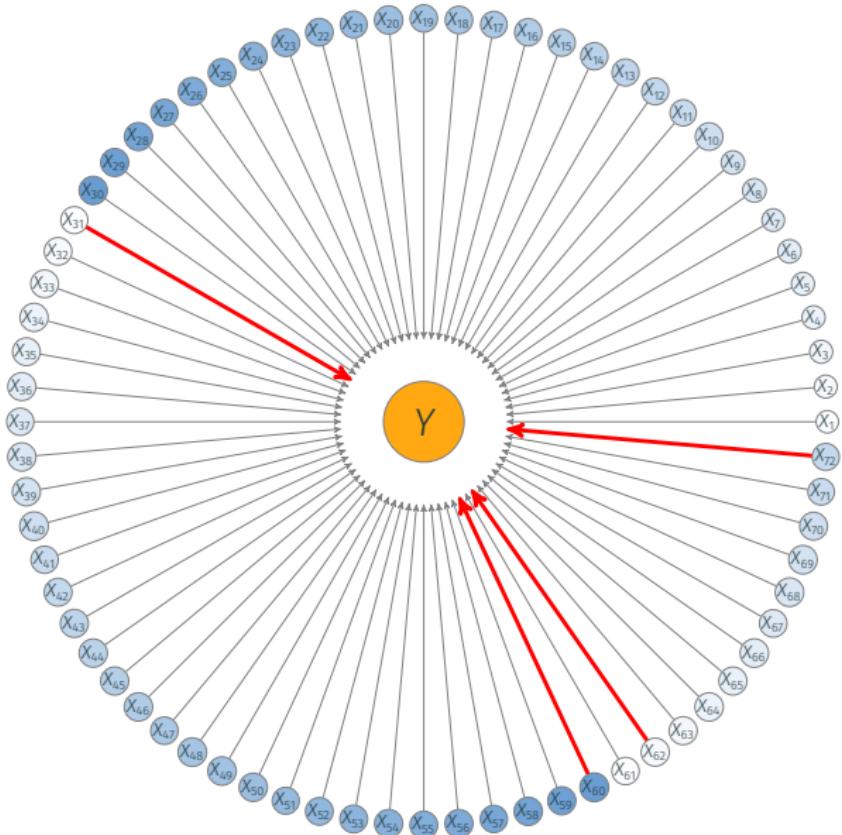
1. A motivating example
2. Background on penalization methods lasso and group lasso
3. Overview of our software packages
4. Present two of our penalization methods: **sail** and **ggmix**

BET 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

- An underlying assumption of **simplicity** in high-dimensional data ($N \ll p$)
- A sparse statistical model is one in which only a **relatively small number of predictors** $k < N$ play an important role
- Sparse models can be faster to compute, easier to understand, and yield more stable predictions.

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

MOTIVATING EXAMPLE

PREDICTORS OF NHL SALARY²



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

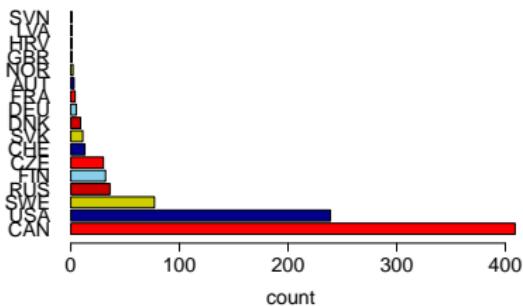
SUPERVISED LEARNING

- Learn the function f

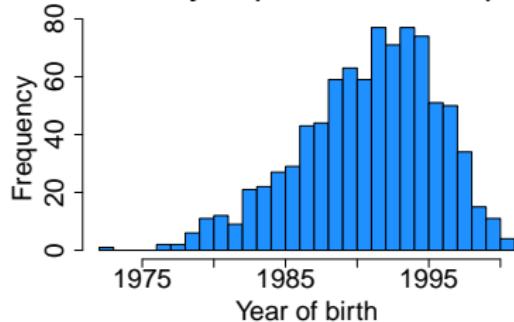


PREDICTORS OF NHL SALARY

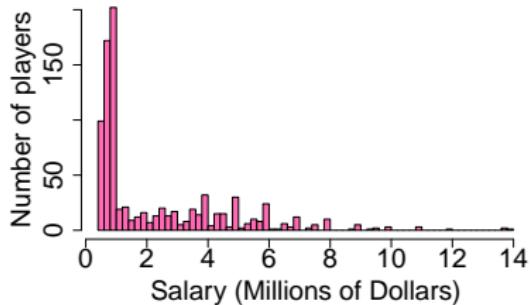
Country



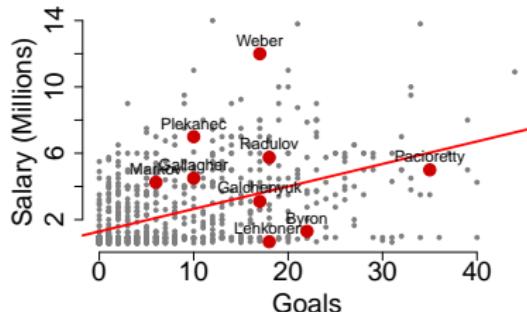
Birth year (2016/2017 season)



NHL Salary Distribution: 2016/2017

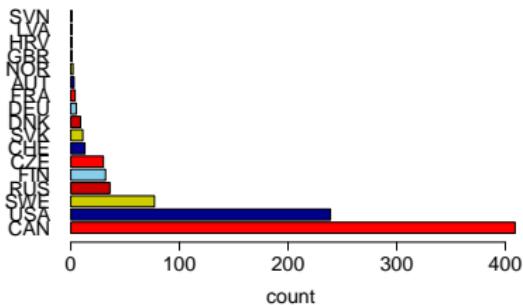


Linear Regression Fit

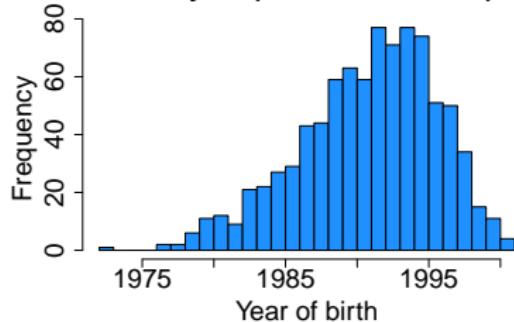


PREDICTORS OF NHL SALARY

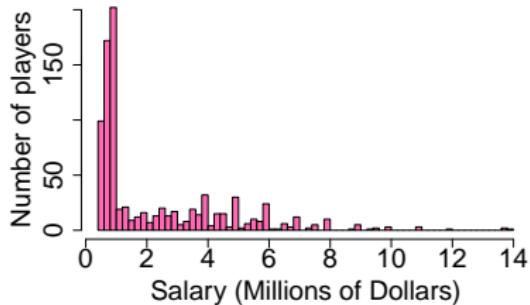
Country



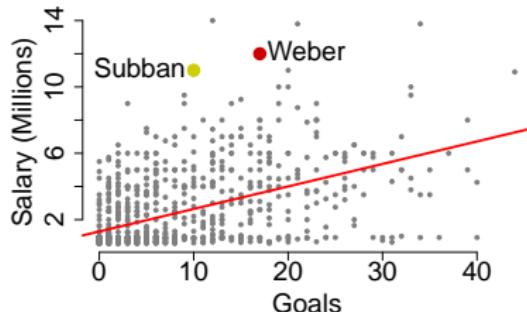
Birth year (2016/2017 season)



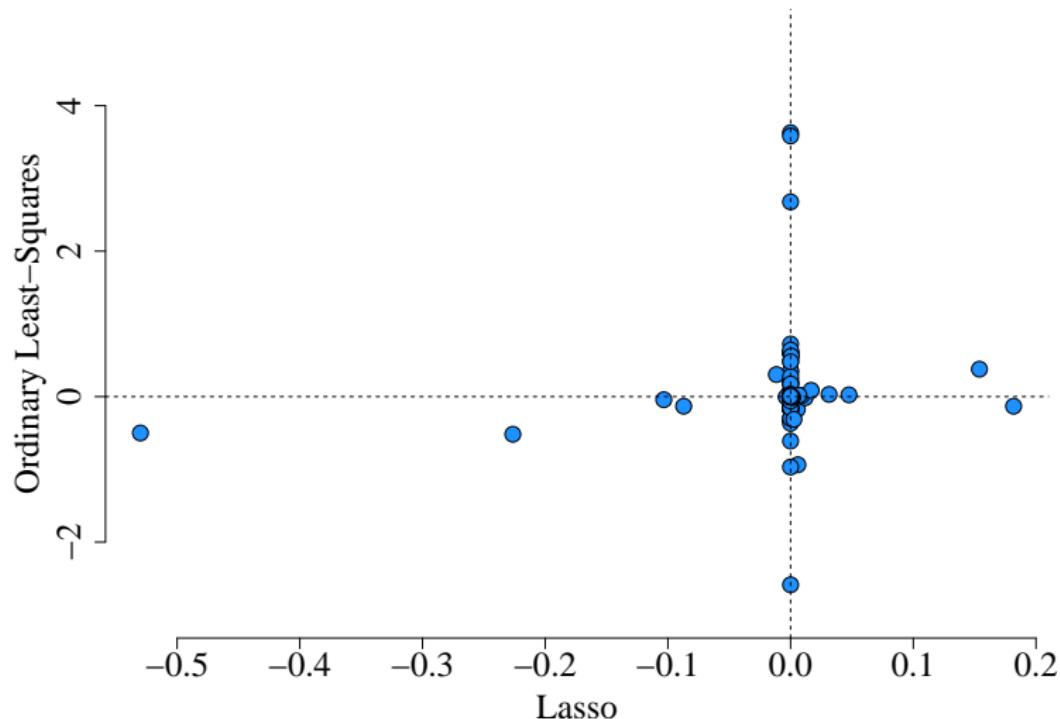
NHL Salary Distribution: 2016/2017



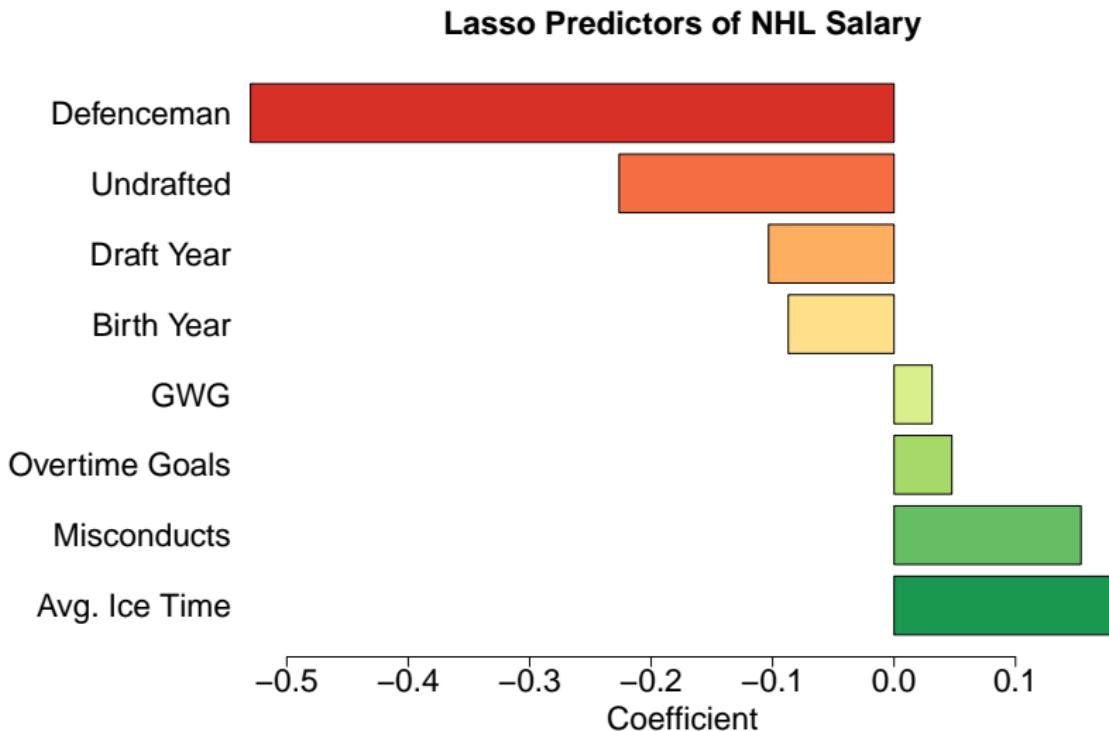
Linear Regression Fit



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

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

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

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

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

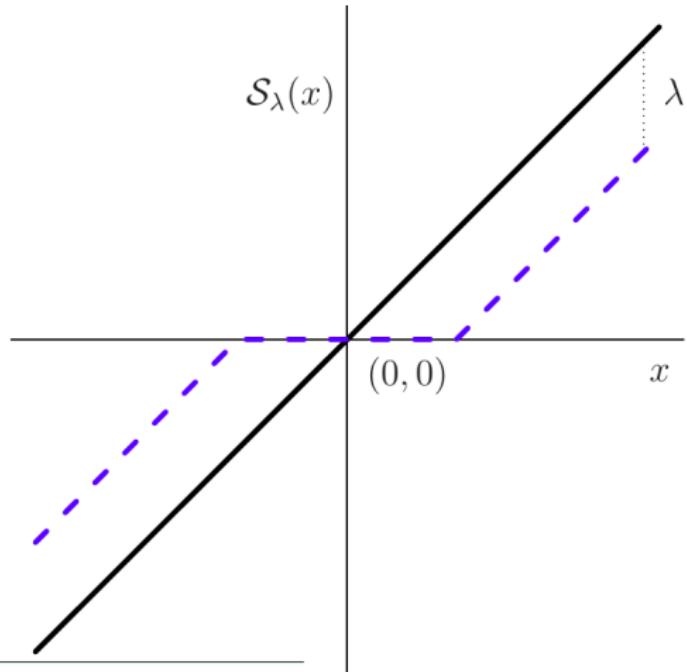
$$\hat{\beta}^{lasso} = \arg \min_{\beta} \frac{1}{2} \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda |\beta| \quad (1)$$

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

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

INSPECTION OF THE LASSO SOLUTION

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



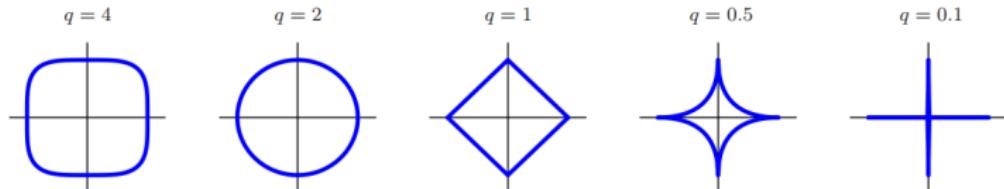
¹Hastie et al. Statistical learning with sparsity: the lasso and generalizations. CRC press, (2015).

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 ℓ_1 norm? Why not use the $q = 2$ (Ridge) or any ℓ_q norm?



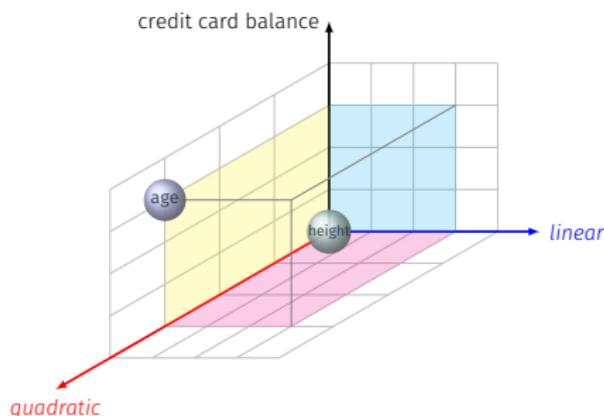
- $q = 1$ is the smallest value that yields a sparse solution **and** yields a **convex** problem → scalable to high-dimensional data
- For $q < 1$ the constrained region is **nonconvex**

CHOOSING THE MODEL COMPLEXITY

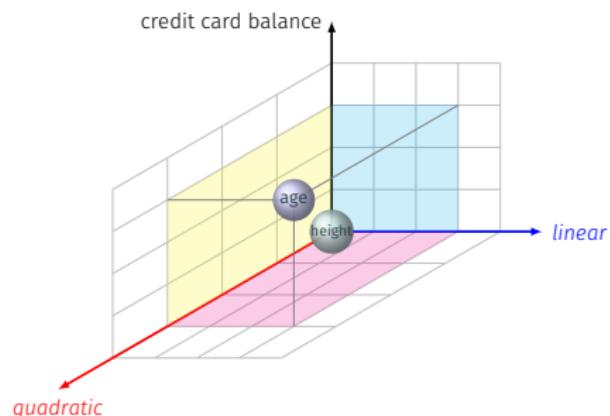
GROUP LASSO

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

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



(a) Lasso



(b) 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 (2017+, preprint)
<https://github.com/sahirbhatnagar/sail>
- **gmmix** – Bhatnagar, Oualkacha, Yang, Greenwood (2017+, preprint)
<https://github.com/sahirbhatnagar/gmmix>
- **casebase** – Bhatnagar¹, Turgeon¹, Yang, Hanley and Saarela (2017+, 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: HEREDITY PROPERTY

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

¹Chipman. Canadian Journal of Statistics (1996)

²McCullagh and Nelder. Generalized Linear Models (1983)

³Cox. International Statistical Review (1984)

MOTIVATION 1: HEREDITY PROPERTY

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

Strong Heredity¹

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

Weak Heredity¹

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

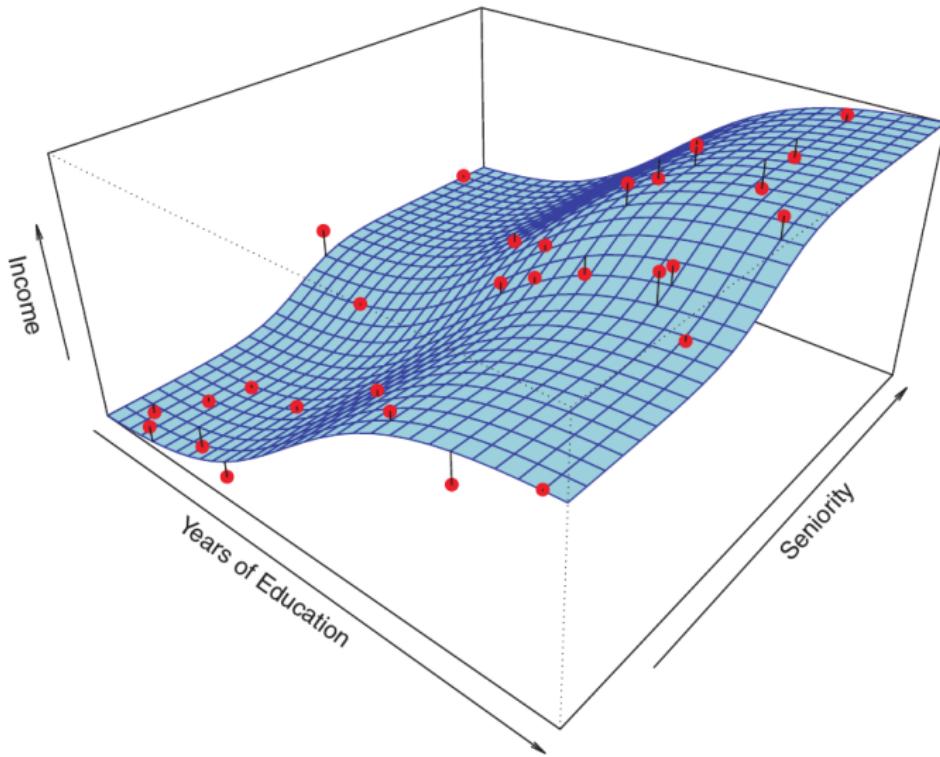
- Heredity property is desired for the purposes of **interpretability**²
- Large main effects are more likely to lead to appreciable interactions³

¹Chipman. Canadian Journal of Statistics (1996)

²McCullagh and Nelder. Generalized Linear Models (1983)

³Cox. International Statistical Review (1984)

MOTIVATION 2: NON-LINEAR INTERACTIONS



¹James et al. Introduction to Statistical Learning (2013)

LASSO INTERACTION MODEL

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

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

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

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>)	hierNet(x, y)
	GRESH (She and Jiang 2014, <i>JASA</i>)	X
	FAMILY (Haris et al. 2014, <i>JCGS</i>)	FAMILY(x, z, y)
	glinternet (Lim and Hastie 2015, <i>JCGS</i>)	glinternet(x, y)
	RAMP (Hao et al. 2016, <i>JASA</i>)	RAMP(x, y)
Non-linear	VANISH (Radchenko and James 2010, <i>JASA</i>)	X
	sail (Bhatnagar et al. 2017+)	sail(x, e, y)

OUR EXTENSION TO NONLINEAR EFFECTS

Consider the basis expansion

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

$$f(x_1) = \underbrace{\begin{bmatrix} \psi_{11}(x_{11}) & \psi_{12}(x_{12}) & \cdots & \psi_{11}(x_{15}) \\ \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \cdots & \vdots \\ \psi_{11}(x_{i1}) & \psi_{12}(x_{i2}) & \cdots & \psi_{11}(x_{i5}) \\ \vdots & \vdots & \cdots & \vdots \\ \vdots & \vdots & \cdots & \vdots \\ \psi_{11}(x_{N1}) & \psi_{12}(x_{N2}) & \cdots & \psi_{11}(x_{N5}) \end{bmatrix}}_{\Psi_1}_{N \times 5} \times \underbrace{\begin{bmatrix} \beta_{11} \\ \beta_{12} \\ \beta_{13} \\ \beta_{14} \\ \beta_{15} \end{bmatrix}}_{\theta_1}_{5 \times 1}$$

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 **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

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

Objective Function

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

¹Choi et al. JASA (2010)

ALGORITHM

BLOCK RELAXATION (DE LEEUW, 1994)

Algorithm 1: Block Relaxation Algorithm

Set the iteration counter $k \leftarrow 0$, initial values for the parameter vector $\Theta^{(0)}$;

for each pair $(\lambda_\beta, \lambda_\gamma)$ **do**

repeat

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

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

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

$$k \leftarrow k + 1$$

until convergence criterion is satisfied;

end

Objective Function

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

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

IMPLEMENTATION

Objective Function

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

Lasso problem

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

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

Objective Function

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

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

Objective Function

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

Group Lasso problem

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

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

SIMULATIONS

SIMULATIONS SCENARIOS

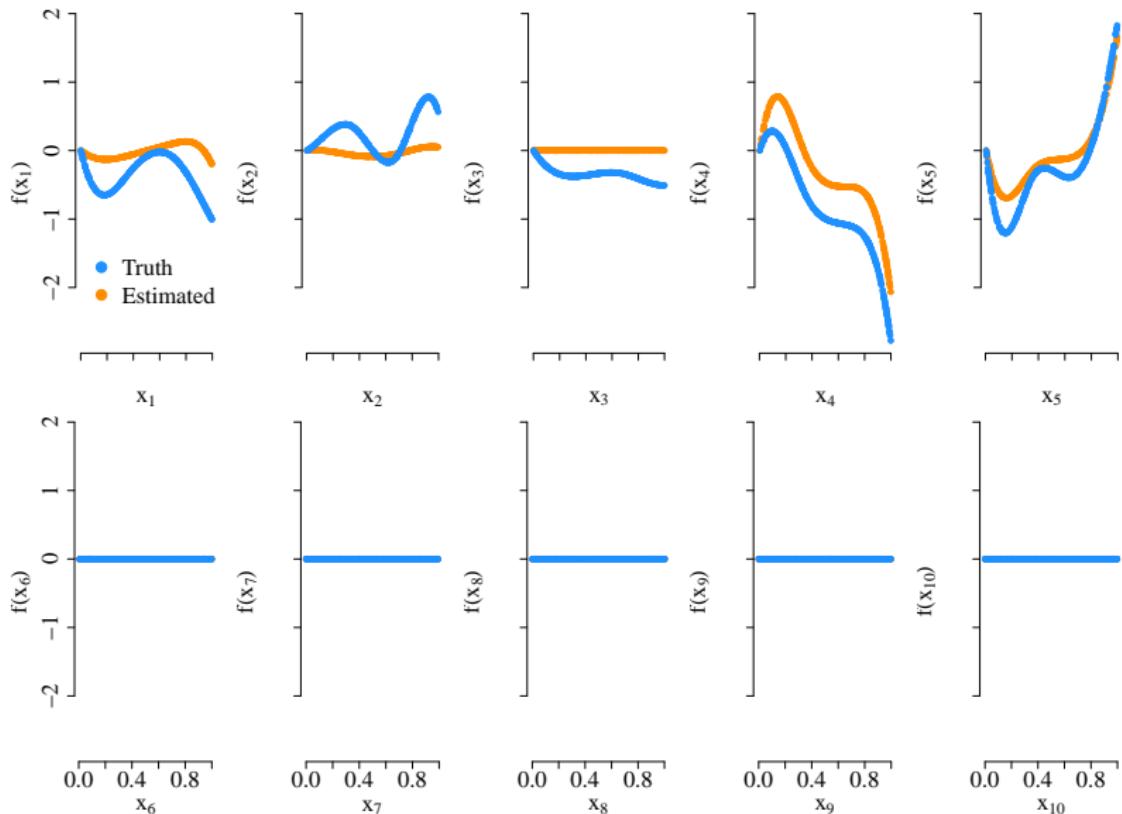
Scenario 1: Easy

- $Y = \sum_{j=1}^5 f(X_j) + X_E + E \times (f(X_1) + f(X_2))$
- $f(\cdot) \rightarrow$ B-splines with 5 df
- $\theta_j \sim \mathcal{N}(0, 1)$
- $N = 400, p = 50$
- $50 \times 5 \times 2 + 1 = 501$ parameters to estimate

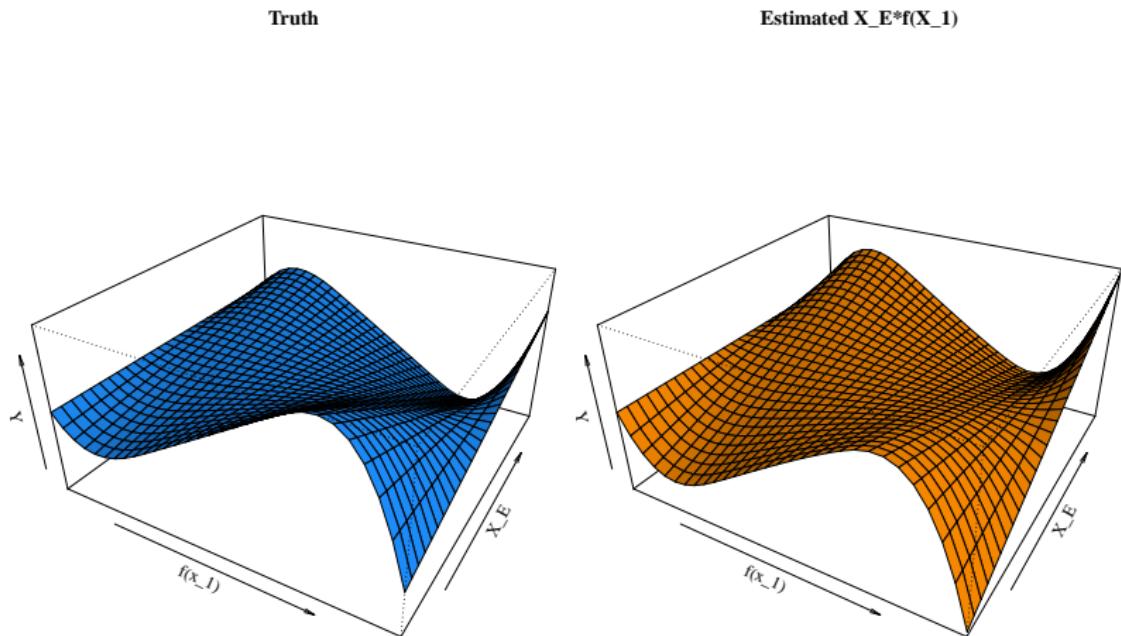
Scenario 2: Hard

- $Y = \sum_{j=1}^4 f(X_j) + X_E + E \times (f(X_3) + f(X_4))$
- $f(X_1) \rightarrow$ linear
- $f(X_2) \rightarrow$ quadratic
- $f(X_3) \rightarrow$ sinusoidal
- $f(X_4) \rightarrow$ complicated sinusoidal

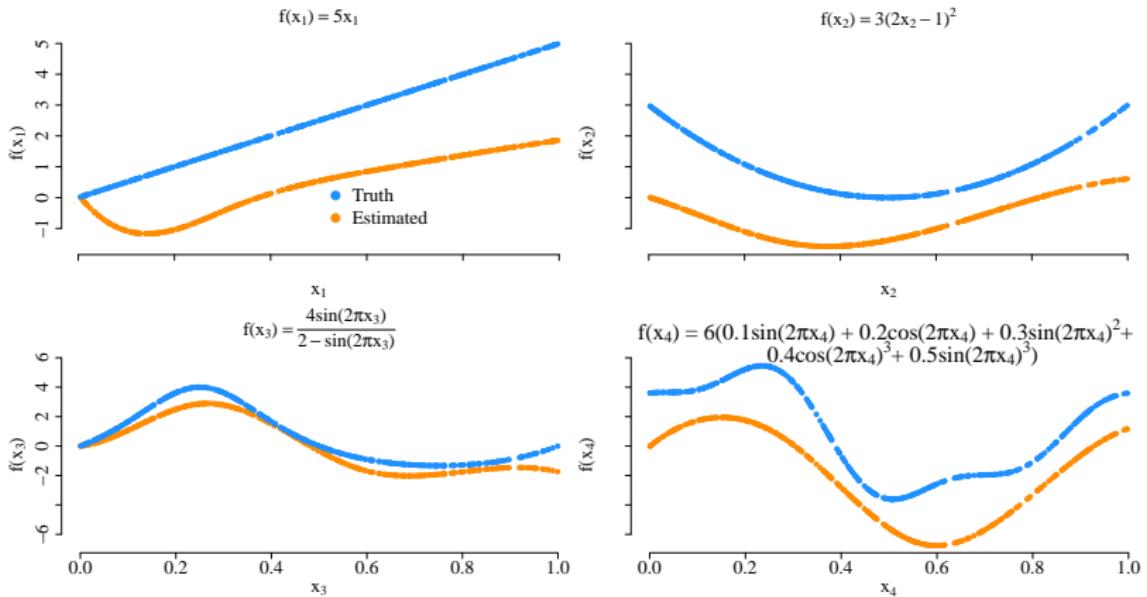
SCENARIO 1: MAIN EFFECTS



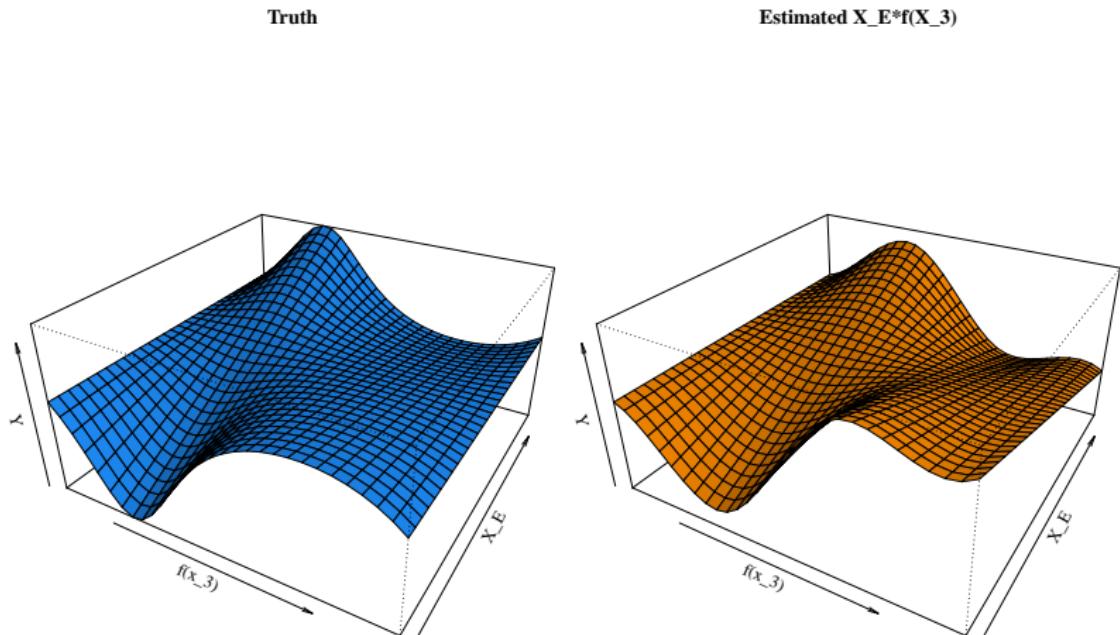
SCENARIO 1: INTERACTION EFFECTS



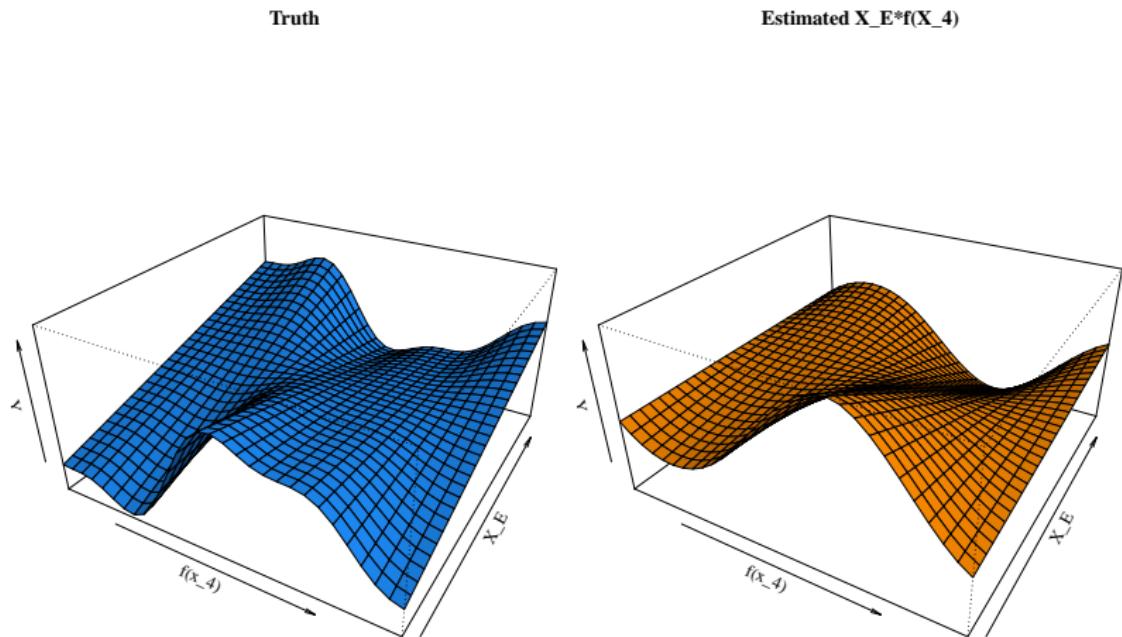
SCENARIO 2: MAIN EFFECTS



SCENARIO 2: INTERACTION EFFECTS

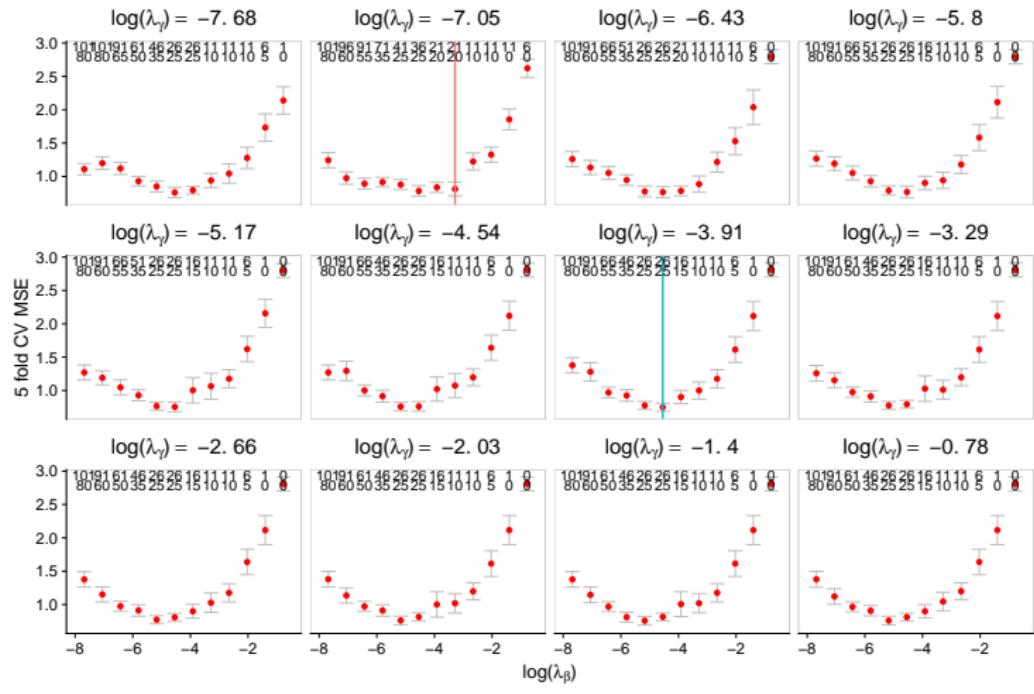


SCENARIO 2: INTERACTION EFFECTS



sail R PACKAGE: CROSS-VALIDATION RESULTS

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



| lambda.1se | lambda.min

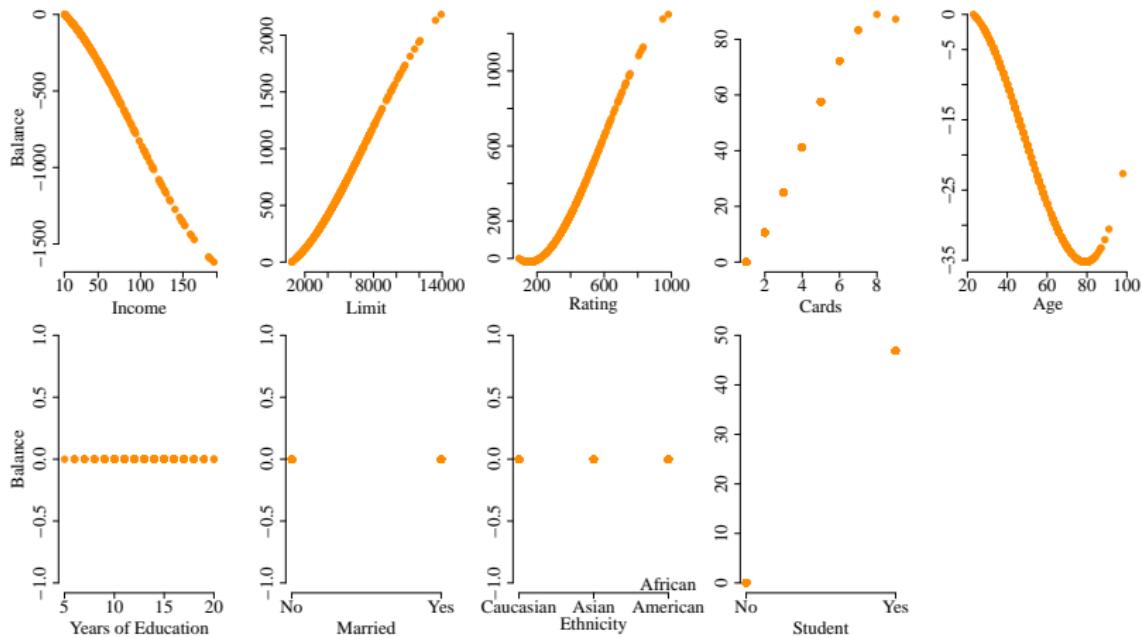
sail: RESULTS ON A REAL DATA SET

- Credit dataset¹: $Y_{400 \times 1} \rightarrow$ Credit card balance, $X_{400 \times 9} \rightarrow$ predictors

¹ISLR package on CRAN

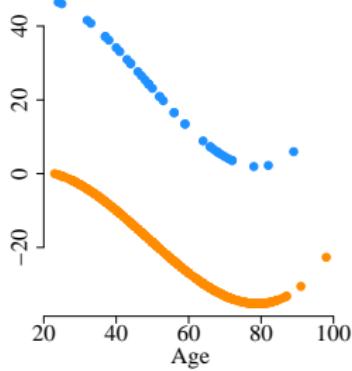
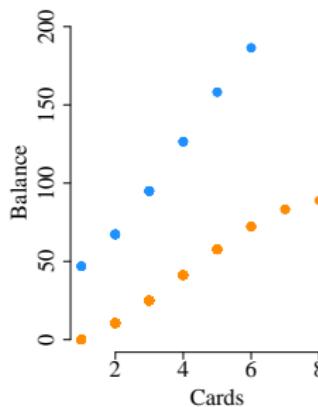
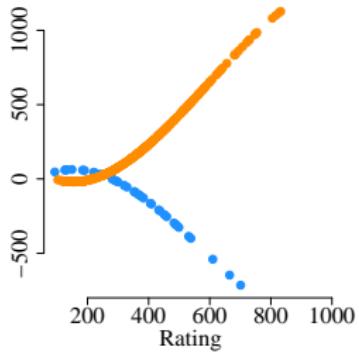
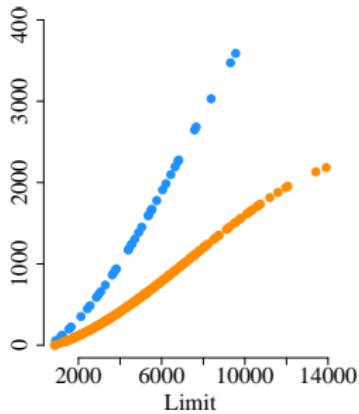
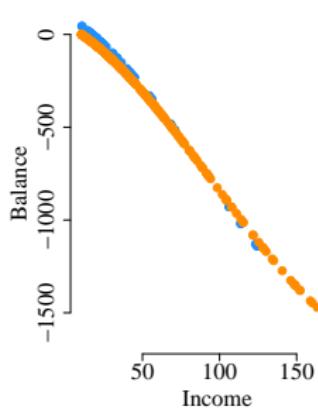
sail: RESULTS ON A REAL DATA SET

■ Credit dataset¹: $\mathbf{Y}_{400 \times 1} \rightarrow$ Credit card balance, $\mathbf{X}_{400 \times 9} \rightarrow$ predictors



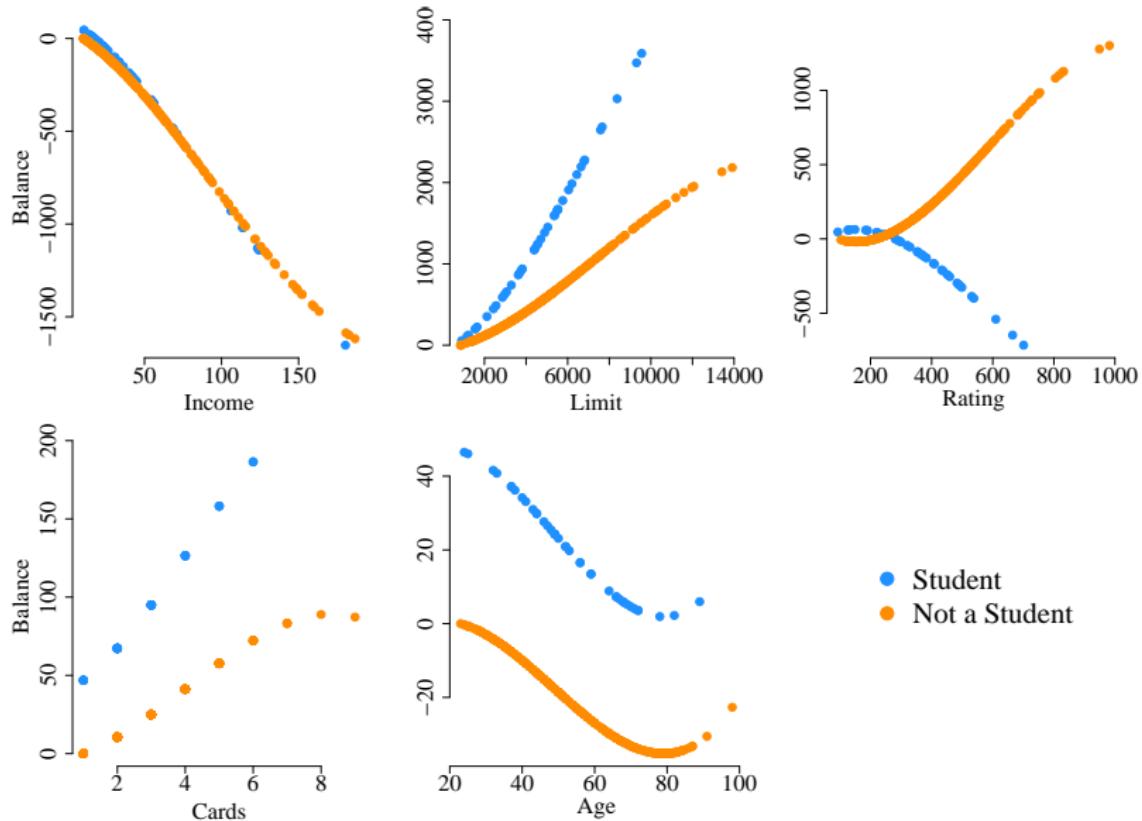
¹ISLR package on CRAN

sail: INTERACTIONS WITH Student



Which color
represents
a Student?

sail: INTERACTIONS WITH Student



ggmix: MIXED MODELS WITH THE GROUP LASSO

MOTIVATING DATASET: UK BIOBANK

- 500,000 individuals with over 40 million predictors
- 1000's of responses (e.g. disease status, bone mineral density)
- **Problem:** Which predictors are associated with the response?



MOTIVATING DATASET: TWO PROBLEMS

	ID	Response	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
1	2610781	-1.255	1	2	0	0	0	1
2	4114347	-0.339	1	2	0	2	0	1
3	4399930	-0.6	1	2	1	1	0	1
4	2081319	0.809	1	2	0	1	0	2
5	1347380	0.279	2	2	0	0	0	0
6	3262449	-0.421	2	2	0	1	0	1
7	4870063	-0.454	2	2	0	0	0	2
8	1141212	1.383	2	2	1	1	1	0
9	2997954	-2.29	1	2	0	0	0	1
10	5805218	2.289	1	2	0	1	1	1

PROBLEM 1: GROUPS OF PREDICTORS AFFECT THE RESPONSE

	ID	Response	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
1	2610781	-1.255	1	2	0	0	0	1
2	4114347	-0.339	1	2	0	2	0	1
3	4399930	-0.6	1	2	1	1	0	1
4	2081319	0.809	1	2	0	1	0	2
5	1347380	0.279	2	2	0	0	0	0
6	3262449	-0.421	2	2	0	1	0	1
7	4870063	-0.454	2	2	0	0	0	2
8	1141212	1.383	2	2	1	1	1	0
9	2997954	-2.29	1	2	0	0	0	1
10	5805218	2.289	1	2	0	1	1	1

PROBLEM 2: OBSERVATIONS ARE NOT INDEPENDENT

- Observations are correlated, but this information is unknown
- However it can be estimated from the data

ID	Response	Gene1	Gene2	Gene3	Gene4	Gene5	Gene6
1	2610781	-1.255	1	2	0	0	0
2	4114347	-0.339	1	2	0	2	0
3	4399930	-0.6	1	2	1	1	0
4	2081319	0.809	1	2	0	1	0
5	1347380	0.279	2	2	0	0	0
6	3262449	-0.421	2	2	0	1	0
7	4870063	-0.454	2	2	0	0	2
8	1141212	1.383	2	2	1	1	1
9	2997954	-2.29	1	2	0	0	1
10	5805218	2.289	1	2	0	1	1

ggmix: SIMULATION RESULTS

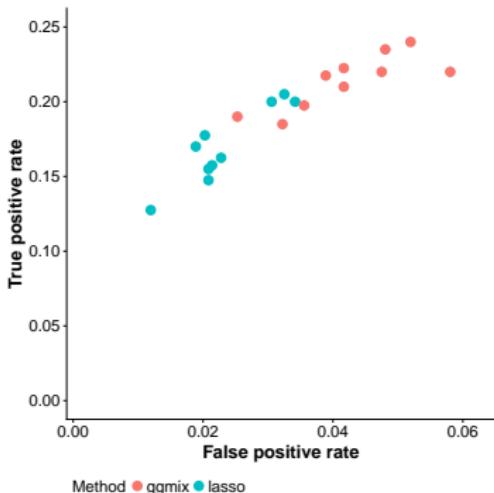


Figure 1: True Positive vs. False Positive (10 simulations). $N = 1000$, $p = 5000$, $p_{active} = 500$

Lasso	ggmix
32.8 (0.87)	26.7 (1.06)

Table 1: Mean root mean squared error (sd)

DATA

- Response: $\mathbf{Y} = (y_1, \dots, y_n) \in \mathbb{R}^n$
- Predictors: $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times p}$, where $p \gg n$
- Similarity Matrix: $\Phi \in \mathbb{R}^{n \times n}$
- Regression Coefficients: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$
- Random effect: $\mathbf{b} = (b_1, \dots, b_n) \in \mathbb{R}^n$
- Error: $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$

DATA

- Response: $\mathbf{Y} = (y_1, \dots, y_n) \in \mathbb{R}^n$
- Predictors: $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times p}$, where $p \gg n$
- Similarity Matrix: $\Phi \in \mathbb{R}^{n \times n}$
- Regression Coefficients: $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T \in \mathbb{R}^p$
- Random effect: $\mathbf{b} = (b_1, \dots, b_n) \in \mathbb{R}^n$
- Error: $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n) \in \mathbb{R}^n$
- We consider the following LMM with a single random effect:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{b} + \boldsymbol{\varepsilon} \quad (2)$$

$$\mathbf{b} \sim \mathcal{N}(0, \eta\sigma^2\Phi) \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, (1 - \eta)\sigma^2\mathbf{I})$$

- σ^2 and $\eta \in [0, 1]$ determine how the variance is divided between \mathbf{b} and $\boldsymbol{\varepsilon}$
- $\mathbf{Y}|(\boldsymbol{\beta}, \eta, \sigma^2) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \eta\sigma^2\Phi + (1 - \eta)\sigma^2\mathbf{I})$

LIKELIHOOD

- The negative log-likelihood is given by

$$-\ell(\Theta) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \log(\det(V)) + \frac{1}{2\sigma^2} (Y - X\beta)^T V^{-1} (Y - X\beta)$$

where $V = \eta\Phi + (1 - \eta)\mathcal{I}$ ¹

¹Pirinen et al. Annals of Applied Statistics (2013)

- The negative log-likelihood is given by

$$-\ell(\Theta) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \log(\det(V)) + \frac{1}{2\sigma^2} (Y - X\beta)^T V^{-1} (Y - X\beta)$$

where $V = \eta\Phi + (1 - \eta)\mathcal{I}$ ¹

- Assume we have a low-rank matrix $K \in \mathbb{R}^{n \times k}$ ($k < n$), to compute the factored similarity matrix $\Phi = KK^T$

¹Pirinen et al. Annals of Applied Statistics (2013)

- The negative log-likelihood is given by

$$-\ell(\Theta) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \log(\det(V)) + \frac{1}{2\sigma^2} (Y - X\beta)^T V^{-1} (Y - X\beta)$$

where $V = \eta\Phi + (1 - \eta)\mathbf{I}$ ¹

- Assume we have a low-rank matrix $K \in \mathbb{R}^{n \times k}$ ($k < n$), to compute the factored similarity matrix $\Phi = KK^T$
- Let $K = U\Lambda V^T$ be the SVD of K , then

$$\Phi = U_1\Lambda\Lambda U_1^T = U_1\Sigma U_1^T$$

where $U_1 \in \mathbb{R}^{n \times k}$ is the matrix of singular vectors corresponding to the k non-zero eigenvalues.

¹Pirinen et al. Annals of Applied Statistics (2013)

PENALIZED MAXIMUM LIKELIHOOD ESTIMATOR

- The negative log-likelihood can then be expressed as

$$-\ell(\Theta) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \left(\sum_{i=1}^k \log(1 + \eta(\Sigma_i - 1)) + (n - k) \log(1 - \eta) \right) + \frac{1}{2} \left\{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \left[\frac{1}{\sigma^2(1 - \eta)} \left(\mathbf{I}_n - \mathbf{U}_1 \left(\frac{1 - \eta}{\eta} \Sigma_1^{-1} + \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right) \right] (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

PENALIZED MAXIMUM LIKELIHOOD ESTIMATOR

- The negative log-likelihood can then be expressed as

$$-\ell(\Theta) \propto \frac{n}{2} \log(\sigma^2) + \frac{1}{2} \left(\sum_{i=1}^k \log(1 + \eta(\Sigma_i - 1)) + (n - k) \log(1 - \eta) \right) + \frac{1}{2} \left\{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T \left[\frac{1}{\sigma^2(1 - \eta)} \left(\mathbf{I}_n - \mathbf{U}_1 \left(\frac{1 - \eta}{\eta} \Sigma_1^{-1} + \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right) \right] (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right\}$$

- Define the objective function:

$$Q_\lambda(\Theta) = -\ell(\Theta) + \lambda \sum_j P_j(\beta_j)$$

- $P_j(\cdot)$ is a penalty term on $\beta_1, \dots, \beta_{p+1}$
- An estimate of the regression parameters $\hat{\Theta}_\lambda$ is obtained by

$$\hat{\Theta}_\lambda = \arg \min_{\Theta} Q_\lambda(\Theta)$$

GROUP LASSO FOR MIXED MODELS

- Assume the predictors in $X \in \mathbb{R}^{n \times p}$ belong to K **non-overlapping groups** with pre-defined group membership and cardinality p_k
- Let $\beta_{(k)}$ to denote the segment of β corresponding to group k

GROUP LASSO FOR MIXED MODELS

- Assume the predictors in $\mathbf{X} \in \mathbb{R}^{n \times p}$ belong to K non-overlapping groups with pre-defined group membership and cardinality p_k
- Let $\boldsymbol{\beta}_{(k)}$ to denote the segment of $\boldsymbol{\beta}$ corresponding to group k
- We consider the group lasso penalized estimator

$$\min_{\boldsymbol{\beta}} L(\boldsymbol{\beta} | \mathsf{D}) + \lambda \sum_{k=1}^K w_k \|\boldsymbol{\beta}_{(k)}\|_2, \quad (3)$$

- where

$$L(\boldsymbol{\beta} | \mathsf{D}) = \frac{1}{2} [\mathbf{Y} - \widehat{\mathbf{Y}}]^T \mathbf{W} [\mathbf{Y} - \widehat{\mathbf{Y}}] \quad (4)$$

$\widehat{\mathbf{Y}} = \sum_{j=1}^p \beta_j \mathbf{x}_j$, D is the working data $\{\mathbf{Y}, \mathbf{X}\}$, and

$$\mathbf{W}_{n \times n} = \frac{1}{\sigma^2(1-\eta)} \left(\mathbf{I}_n - \mathbf{U}_1 \left(\frac{1-\eta}{\eta} \Sigma_1^{-1} + \mathbf{I}_k \right)^{-1} \mathbf{U}_1^T \right) \quad (5)$$

GROUPWISE DESCENT: EXPLOITING SPARSITY STRUCTURE

Minimize the objective function

$$\frac{1}{2} [\mathbf{Y} - \widehat{\mathbf{Y}}]^T \mathbf{W} [\mathbf{Y} - \widehat{\mathbf{Y}}] + \lambda \sum_{k=1}^K w_k \|\boldsymbol{\beta}^{(k)}\|_2$$

During each sub-iteration only optimize $\boldsymbol{\beta}^{(k)}$. Set $\boldsymbol{\beta}^{(k')} = \widetilde{\boldsymbol{\beta}}^{(k')}$ for $k' \neq k$ at their current value.

1. Initialization: $\widetilde{\boldsymbol{\beta}}$
2. Cyclic groupwise descent: for $k = 1, 2, \dots, K$, update $\boldsymbol{\beta}^{(k)}$ by minimizing the objective function

$$\widetilde{\boldsymbol{\beta}}^{(k)}(\text{new}) \leftarrow \arg \min_{\boldsymbol{\beta}^{(k)}} L(\boldsymbol{\beta} \mid \mathbf{D}) + \lambda w_k \|\boldsymbol{\beta}^{(k)}\|_2$$

3. Repeat (2) till convergence.

QUADRATIC MAJORIZATION CONDITION

$$\arg \min_{\beta^{(k)}} \frac{1}{2} [\mathbf{Y} - \hat{\mathbf{Y}}]^T \mathbf{W} [\mathbf{Y} - \hat{\mathbf{Y}}] + \lambda \sum_{k=1}^K w_k \|\beta^{(k)}\|_2 \quad (6)$$

- Unfortunately, there is no closed form solution to (6)

¹Yang and Zou. Statistical Computing (2014)

QUADRATIC MAJORIZATION CONDITION

$$\arg \min_{\beta^{(k)}} \frac{1}{2} [\mathbf{Y} - \hat{\mathbf{Y}}]^T \mathbf{W} [\mathbf{Y} - \hat{\mathbf{Y}}] + \lambda \sum_{k=1}^K w_k \|\beta^{(k)}\|_2 \quad (6)$$

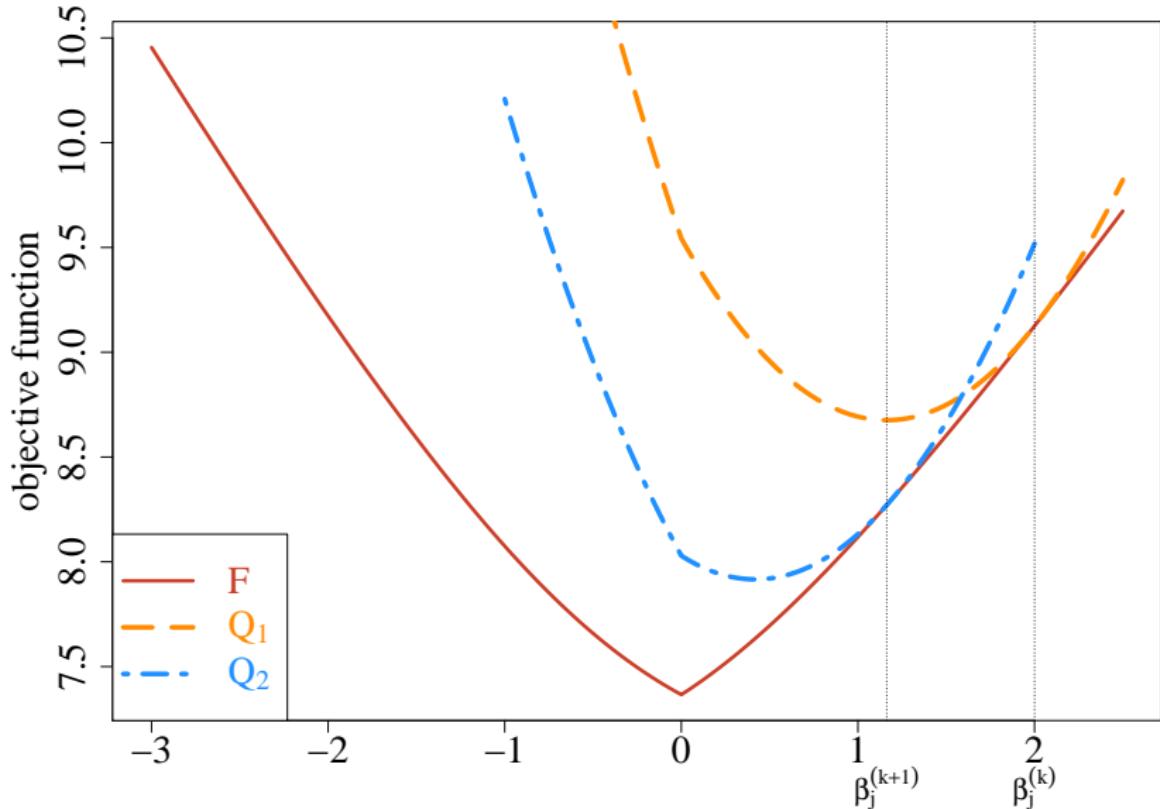
- Unfortunately, there is no closed form solution to (6)
- However, the loss function $L(\beta | \mathbf{D})$ satisfies the quadratic majorization (QM) condition¹, since there exists
 - a $p \times p$ matrix $\mathbf{H} = \mathbf{X}^T \mathbf{W} \mathbf{X}$, and
 - $\nabla L(\beta | \mathbf{D}) = -(\mathbf{Y} - \hat{\mathbf{Y}})^T \mathbf{W} \mathbf{X}$

which may only depend on the data \mathbf{D} , such that for all β, β^* ,

$$L(\beta | \mathbf{D}) \leq L(\beta^* | \mathbf{D}) + (\beta - \beta^*)^T \nabla L(\beta^* | \mathbf{D}) + \frac{1}{2} (\beta - \beta^*)^T \mathbf{H} (\beta - \beta^*)$$

¹Yang and Zou. Statistical Computing (2014)

GENERALIZED COORDINATE DESCENT (GCD)



GROUPWISE MAJORIZATION DESCENT

- Update β in a groupwise fashion

$$\beta - \tilde{\beta} = (\underbrace{0, \dots, 0}_{k-1}, \beta^{(k)} - \tilde{\beta}^{(k)}, \underbrace{0, \dots, 0}_{K-k})$$

GROUPWISE MAJORIZATION DESCENT

- Update β in a groupwise fashion

$$\beta - \tilde{\beta} = (\underbrace{0, \dots, 0}_{k-1}, \beta^{(k)} - \tilde{\beta}^{(k)}, \underbrace{0, \dots, 0}_{K-k})$$

- Only need to compute the majorization function on group level

$$L(\beta | D) \leq L(\tilde{\beta} | D) - (\beta^{(k)} - \tilde{\beta}^{(k)})^\top U^{(k)} + \frac{1}{2} \gamma_k (\beta^{(k)} - \tilde{\beta}^{(k)})^\top (\beta^{(k)} - \tilde{\beta}^{(k)})$$

$$U^{(k)} = \frac{\partial}{\partial \beta^{(k)}} L(\beta | D) = -(\gamma - \hat{\gamma})^\top W X_{(k)}$$

$$H^{(k)} = \frac{\partial^2}{\partial \beta^{(k)} \partial \beta_{(k)}^\top} L(\beta | D) = X_{(k)}^\top W X_{(k)}$$

- $\gamma_k = \text{eigen}_{\max}(H^{(k)})$

GROUPWISE MAJORIZATION DESCENT

- Update β in a groupwise fashion

$$\beta - \tilde{\beta} = (\underbrace{0, \dots, 0}_{k-1}, \beta^{(k)} - \tilde{\beta}^{(k)}, \underbrace{0, \dots, 0}_{K-k})$$

- Only need to compute the majorization function on group level

$$L(\beta | D) \leq L(\tilde{\beta} | D) - (\beta^{(k)} - \tilde{\beta}^{(k)})^\top U^{(k)} + \frac{1}{2} \gamma_k (\beta^{(k)} - \tilde{\beta}^{(k)})^\top (\beta^{(k)} - \tilde{\beta}^{(k)})$$

$$U^{(k)} = \frac{\partial}{\partial \beta^{(k)}} L(\beta | D) = -(\gamma - \hat{\gamma})^\top W X_{(k)}$$

$$H^{(k)} = \frac{\partial^2}{\partial \beta^{(k)} \partial \beta_{(k)}^\top} L(\beta | D) = X_{(k)}^\top W X_{(k)}$$

- $\gamma_k = \text{eigen}_{\max}(H^{(k)})$

- Update $\tilde{\beta}^{(k)}$ with a fast operation:

$$\tilde{\beta}^{(k)} (\text{new}) = \frac{1}{\gamma_k} \left(U^{(k)} + \gamma_k \tilde{\beta}^{(k)} \right) \left(1 - \frac{\lambda w_k}{\| U^{(k)} + \gamma_k \tilde{\beta}^{(k)} \|_2} \right)_+$$

DISCUSSION

STRENGTHS AND LIMITATIONS

Strengths

- Environment interactions with strong heredity property in $p \gg N$
- `sail` allows for flexible modeling of input variables
- `gmmix` package provides first implementation of group lasso in mixed models

STRENGTHS AND LIMITATIONS

Strengths

- Environment interactions with strong heredity property in $p \gg N$
- `sail` allows for flexible modeling of input variables
- `gmmix` package provides first implementation of group lasso in mixed models

Limitations

- `sail` can currently only handle $E \cdot f(X)$ or $f(E) \cdot X$
- Does not allow for $f(X_1, E)$ or $f(X_1, X_2)$
- Current implementation of `sail` is slow due to cross validation for 2 tuning parameters
- Memory footprint is an issue for both `sail` and `gmmix`

FUTURE DIRECTIONS

- Are two tuning parameters really necessary ?

$$\lambda \left\{ (1 - \alpha) \left[w_E |\beta_E| + \sum_{j=1}^p w_j \|\theta_j\|_2 \right] + \alpha \sum_{j=1}^p w_{jE} |\gamma_j| \right\}$$

- Weak heredity property $\rightarrow \alpha_j = \gamma_j(|\beta_j| + |\beta_E|)$
- Implement ADMM algorithm for scalability. Distributed computing (GPU)
- Extension to nonconvex penalties (SCAD, MCP)
- Binary Outcomes

ACKNOWLEDGEMENTS

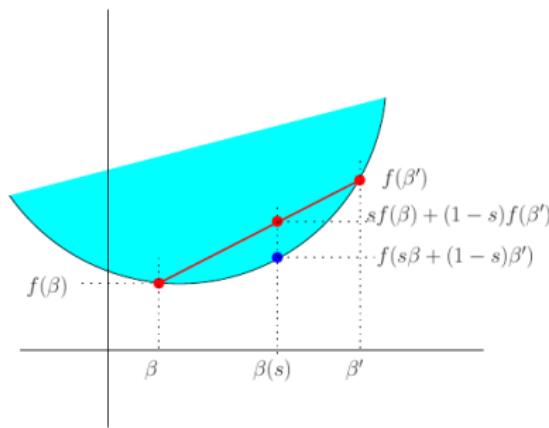


REFERENCES

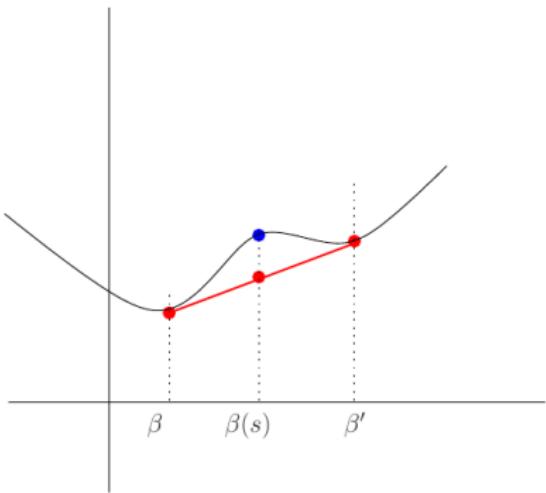
- Radchenko, P., & James, G. M. (2010). Variable selection using adaptive nonlinear interaction structures in high dimensions. *Journal of the American Statistical Association*, 105(492), 1541-1553.
- Choi, N. H., Li, W., & Zhu, J. (2010). Variable selection with the strong heredity constraint and its 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.

APPENDIX

CONVEX FUNCTION



(a)

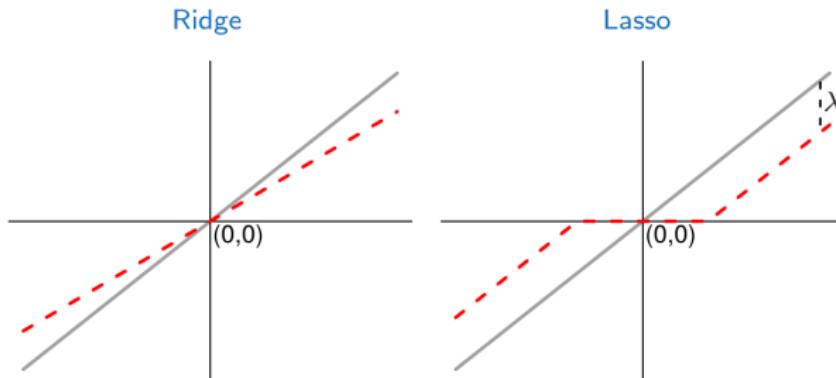


(b)

RIDGE VS. LASSO

- Estimators of β_j in terms of the least-squares estimate $\hat{\beta}_j^{LS}$ for an orthonormal model matrix X

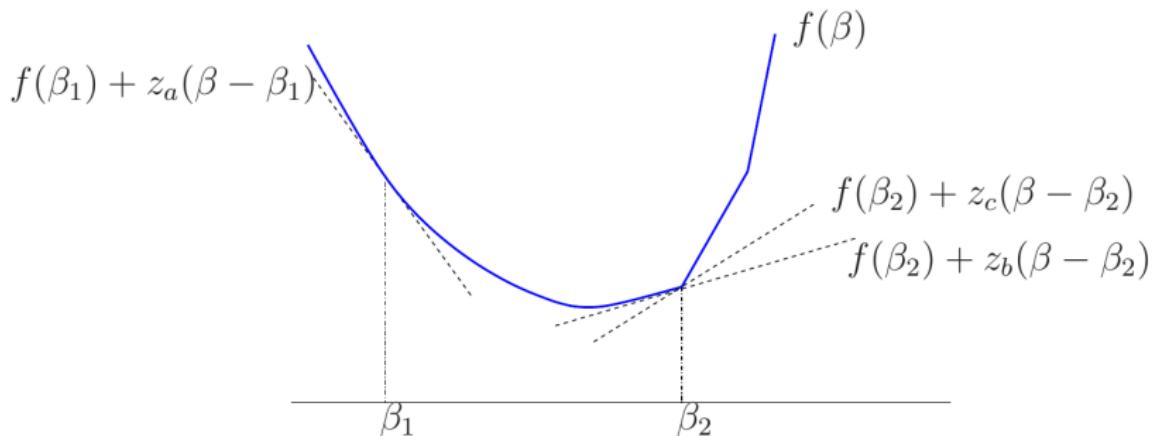
q	Estimator	Formula
1	Lasso	$\text{sign}(\hat{\beta}_j^{LS})(\hat{\beta}_j^{LS} - \lambda)_+$
2	Ridge	$\hat{\beta}_j^{LS} / (1 + \lambda)$



SUBGRADIENT

- A basic property of differentiable convex functions is that the first-order tangent approximation always provides a lower bound.
- The notion of subgradient is based on a natural generalization of this idea. In particular, given a convex function $f : \mathbb{R}^p \rightarrow \mathbb{R}$, a vector $z \in \mathbb{R}^p$ is said to be a *subgradient* of f at β if

$$f(\beta') \geq f(\beta) + \langle z, \beta' - \beta \rangle, \quad \text{for all } \beta' \in \mathbb{R}^p$$



ALGORITHM FOR ggmix

Algorithm 2: Block Relaxation Algorithm

Set the iteration counter $k \leftarrow 0$, initial values for the parameter vector $\Theta^{(0)}$ and convergence threshold ϵ ;

for $\lambda \in \{\lambda_{\max}, \dots, \lambda_{\min}\}$ **do**

repeat

$$\boldsymbol{\beta}^{(k+1)} \leftarrow \arg \min_{\boldsymbol{\beta}} Q_{\lambda} \left(\boldsymbol{\beta}, \eta^{(k)}, \sigma^2^{(k)} \right)$$

$$\eta^{(k+1)} \leftarrow \arg \min_{\eta} Q_{\lambda} \left(\boldsymbol{\beta}^{(k+1)}, \eta, \sigma^2^{(k)} \right)$$

$$\sigma^2^{(k+1)} \leftarrow \arg \min_{\sigma^2} Q_{\lambda} \left(\boldsymbol{\beta}^{(k+1)}, \eta^{(k+1)}, \sigma^2 \right)$$

$k \leftarrow k + 1$

until convergence criterion is satisfied: $\|\Theta^{(k+1)} - \Theta^{(k)}\|_2 < \epsilon$;

end

COORDINATE DESCENT

- Minimize the objective function

$$\min_{\beta} F(\beta) \equiv L(\beta_1, \dots, \beta_p) + p_\lambda(\beta)$$

- $L : \mathbb{R}^p \rightarrow \mathbb{R}$ is differentiable convex, $p_\lambda : \mathbb{R}^p \rightarrow \mathbb{R}$ bounded from below but not necessarily smooth. e.g. $p_\lambda(\beta) = \lambda \sum_{j=1}^p |\beta_j|$

Coordinate Descent (Fu (1998), Friedman et al. (2007), Wu and Lange (2008))

1. Initialization: $\tilde{\beta}$
2. Cyclic coordinate descent: for $j = 1, 2, \dots, p$, update β_j by minimizing the objective function

$$\tilde{\beta}_j^{new} \leftarrow \arg \min_{\beta_j} F(\beta_j | \beta_k = \tilde{\beta}_k, k \neq j)$$

3. Repeat (2) till convergence.

QUADRATIC MAJORIZATION CONDITION

Empirical loss

$$L(\beta \mid D) = \frac{1}{n} \sum_{i=1}^n \Phi(y_i, \beta^\top x_i)$$

Φ satisfies the [QM condition](#), if and only if:

1. $L(\beta \mid D)$ is [differentiable](#) as a function of β .
2. Can find H (p by p), may only depend on the data D , such that for all β, β^*

$$L(\beta \mid D) \leq L(\beta^* \mid D) + (\beta - \beta^*)^\top \nabla L(\beta^* \mid D) + \frac{1}{2} (\beta - \beta^*)^\top H (\beta - \beta^*)$$

Loss	$-\nabla L(\beta \mid D)$	H
Least squares	$\frac{1}{n} \sum_{i=1}^n (y_i - x_i^\top \beta) x_i$	$X^\top X / n$
Logistic regression	$\frac{1}{n} \sum_{i=1}^n y_i x_i \frac{1}{1 + \exp(y_i x_i^\top \beta)}$	$\frac{1}{4} X^\top X / n$
Squared hinge loss	$\frac{1}{n} \sum_{i=1}^n 2y_i x_i (1 - y_i x_i^\top \beta)_+$	$4X^\top X / n$
Huberized hinge loss	$\frac{1}{n} \sum_{i=1}^n y_i x_i \text{hsvm}'(y_i x_i^\top \beta)$	$\frac{2}{\delta} X^\top X / n$

VERIFYING QUADRATIC MAJORIZATION CONDITION

Lemma 1 (Yang and Zou, 2014)

Assume $\Phi(y, f)$ is differentiable with respect to f and write

$$\Phi'_f = \frac{\partial \Phi(y, f)}{\partial f}, \quad \nabla L(\boldsymbol{\beta} | D) = \frac{1}{n} \sum_{i=1}^n \Phi'_f(y_i, \mathbf{x}_i^\top \boldsymbol{\beta}) \mathbf{x}_i.$$

1. If Φ'_f is Lipschitz continuous with constant C such that

$$|\Phi'_f(y, f_1) - \Phi'_f(y, f_2)| \leq C|f_1 - f_2| \quad \forall y, f_1, f_2,$$

then the QM condition holds for Φ and $H = \frac{2C}{n} \mathbf{X}^\top \mathbf{X}$.

2. If $\Phi''_f = \frac{\partial \Phi^2(y, f)}{\partial f^2}$ exists and

$$\Phi''_f \leq C_2 \quad \forall y, f,$$

then the QM condition holds for Φ and $H = \frac{C_2}{n} \mathbf{X}^\top \mathbf{X}$.

STRICT DESCENT PROPERTY OF GMD

Proposition (Yang and Zou, 2014)

- If $\tilde{\beta}^{(k)}$ (new) $\neq \tilde{\beta}^{(k)}$ then

$$L(\tilde{\beta}^{(k)} \text{ (new)} | D) + \lambda w_k \|\tilde{\beta}^{(k)} \text{ (new)}\|_2 < L(\tilde{\beta} | D) + \lambda w_k \|\tilde{\beta}^{(k)}\|_2$$

the objective function is strictly decreased after updating all k in a cycle.

- If $\tilde{\beta}^{(k)}$ (new) $= \tilde{\beta}^{(k)}$ for all k , then the solution must satisfy the KKT conditions:

$$-U^{(k)} + \lambda w_k \cdot \frac{\tilde{\beta}^{(k)}}{\|\tilde{\beta}^{(k)}\|_2} = 0 \quad \text{if } \tilde{\beta}^{(k)} \neq 0,$$

$$\left\| U^{(k)} \right\|_2 \leq \lambda w_k \quad \text{if } \tilde{\beta}^{(k)} = 0,$$

which means that the algorithm converges and finds the right answer.

SEPARABILITY AND COORDINATE DESCENT

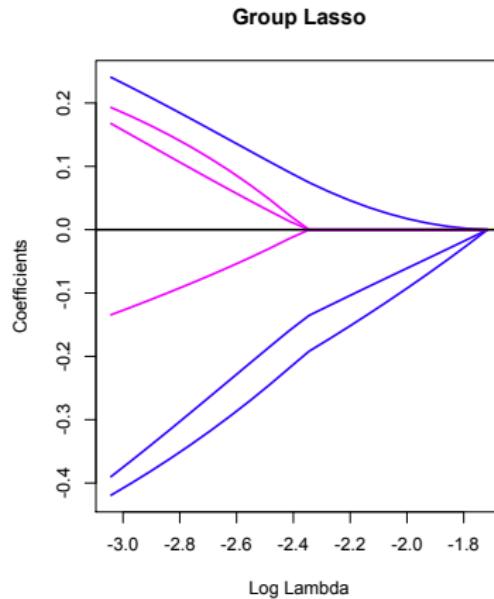
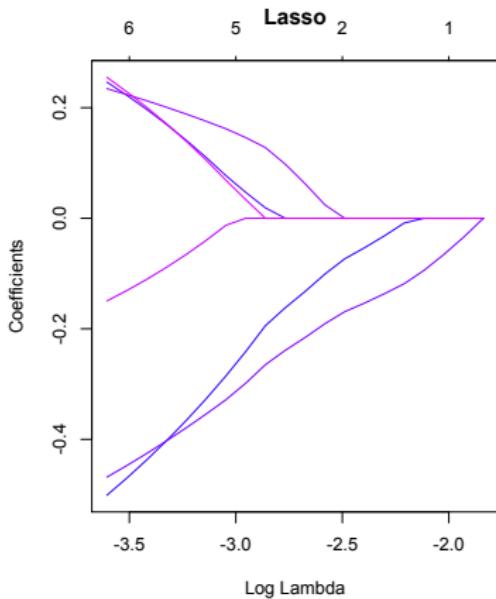
$$f(\beta_1, \dots, \beta_p) = g(\beta_1, \dots, \beta_p) + \sum_{j=1}^p h_j(\beta_j)$$

- $g : \mathbb{R}^p \rightarrow \mathbb{R}$ is differentiable and convex, and the univariate functions $h_j : \mathbb{R} \rightarrow \mathbb{R}$ are convex **but not necessarily differentiable**.
- Tseng (2001) shows that for any convex cost function f with separable structure, the coordinate descent algorithm is guaranteed to converge to the global minimizer
- The key property underlying this result is the separability of the nondifferentiable component $h(\beta) = \sum_{j=1}^p h_j(\beta_j)$, as a sum of functions of each individual parameter.

¹Tseng. Journal of optimization theory and applications (2001)

LASSO VS. GROUP LASSO

- Logistic regression with group lasso: $n = 50, p = 6$.
- Group lasso: specify $(\beta_1, \beta_2, \beta_3), (\beta_4, \beta_5, \beta_6)$. Variable selection at the group level.
- Solution path: view β as function of λ .



GROUP LASSO MOTIVATION

- Categorical predictors (factors): dummy variables
- Additive Model: $\sum_{k=1}^K f_k(x^{(k)}) \approx \sum_{k=1}^K \sum_{m=1}^M \beta_{km} h_m(x^{(k)})$
 - ex. birth weight predicted by the mother's age and weight, Age, Age², Age³ and Weight, Weight², Weight³

Group lasso partitions the variable coefficients into K groups

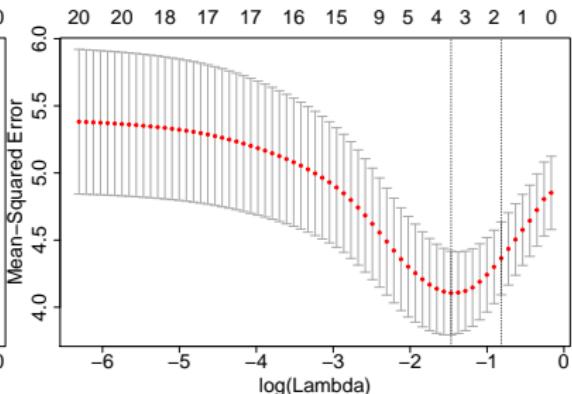
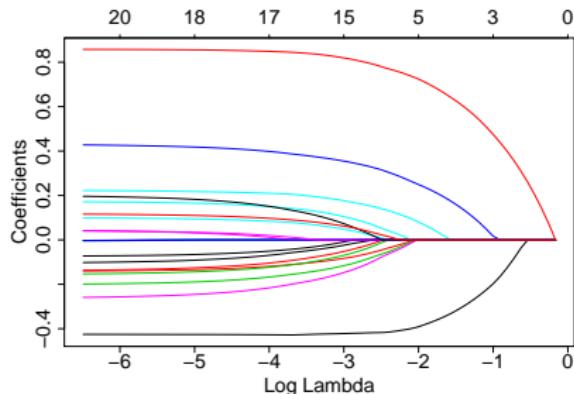
$$\boldsymbol{\beta} = ([\boldsymbol{\beta}^{(1)}]^\top, [\boldsymbol{\beta}^{(2)}]^\top, \dots, [\boldsymbol{\beta}^{(K)}]^\top)^\top$$

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

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

CHOOSING MODEL COMPLEXITY

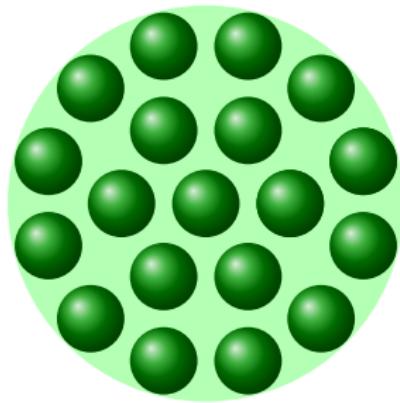
- The tuning parameter λ controls the **complexity** of the model
- **Generalization ability of the model:** we select the λ that gives the most accurate model for predicting independent test data from the same population



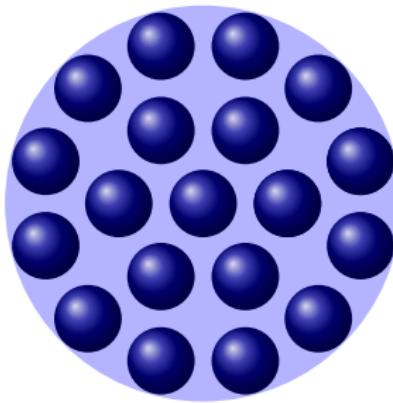
SIMULATION STUDY

- $\mathbf{X}^{(test)}$: 4k SNPs from UK Biobank genotyped data from 1k randomly sampled individuals
- $\mathbf{X}^{(causal)}$: random sample of 400 SNPs from $\mathbf{X}^{(test)}$
- $\mathbf{X}^{(kinship)}$: random sample of 4k SNPs used to construct the kinship matrix (Φ) including the causal SNPs
- $\beta_j \sim Unif(0.1, 0.3)$ for $j = 1, \dots, 400$
- $Y = \beta_0 + \sum_{j=1}^{400} \beta_j \mathbf{x}_j^{(causal)} + P + E$
- $P \sim \mathcal{N}(0, 0.6 \cdot \Phi_{n \times n})$, $E \sim \mathcal{N}(0, 0.4 \cdot \mathbf{I}_{n \times n})$

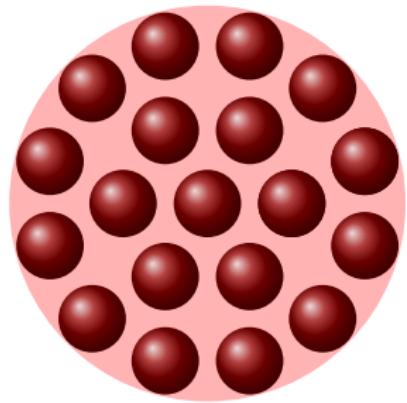
MOTIVATION



(a) Retired

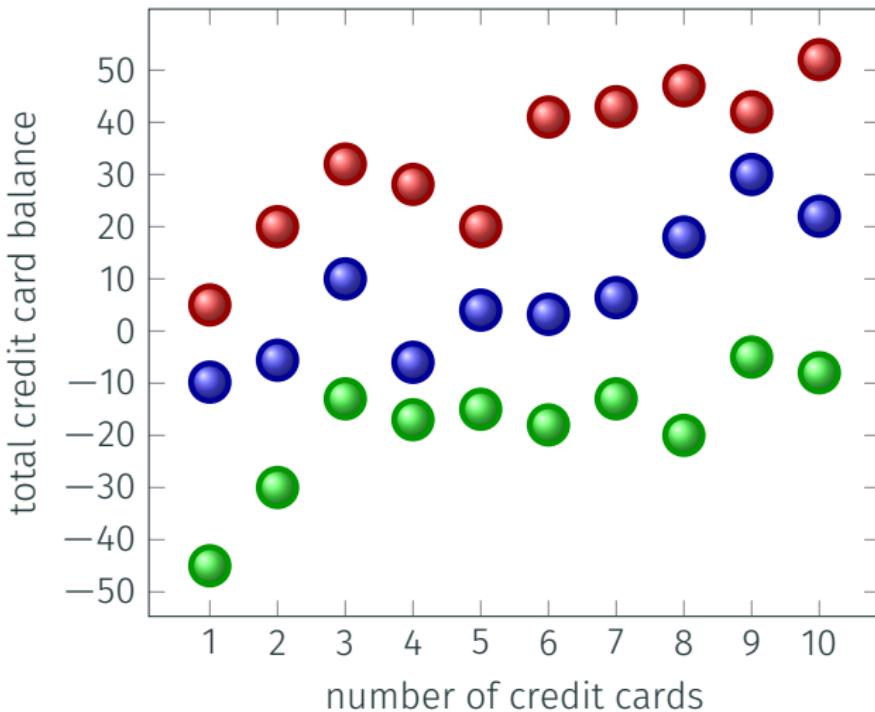


(b) Employed

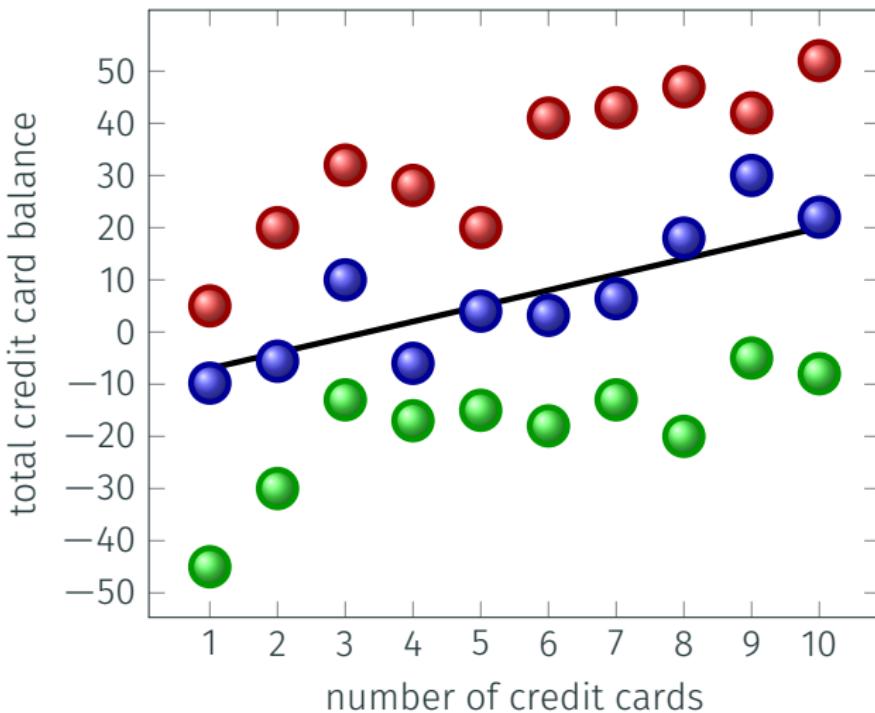


(c) Students

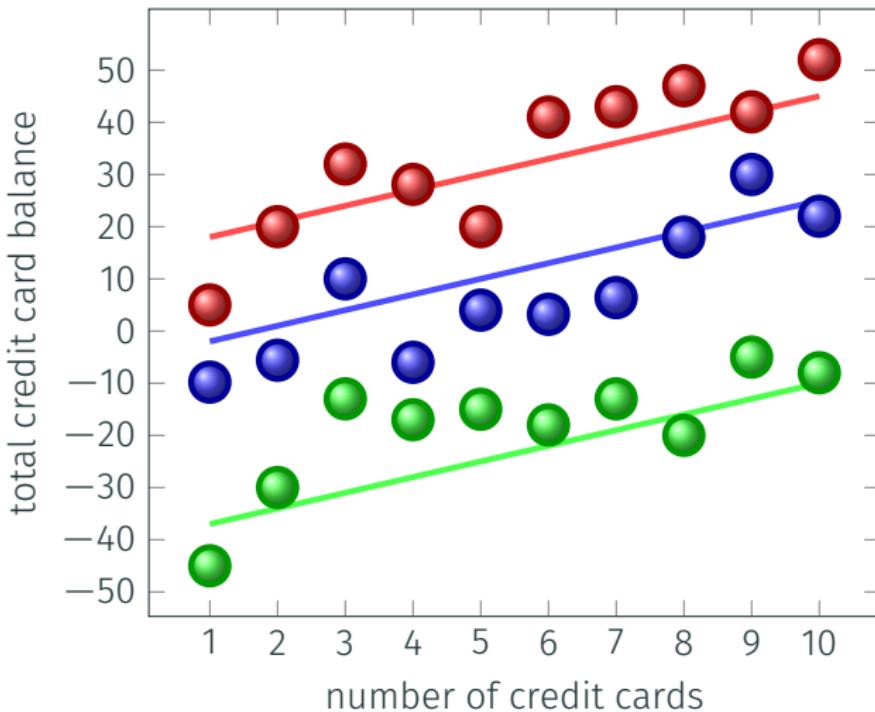
LINEAR MODEL



LINEAR MODEL

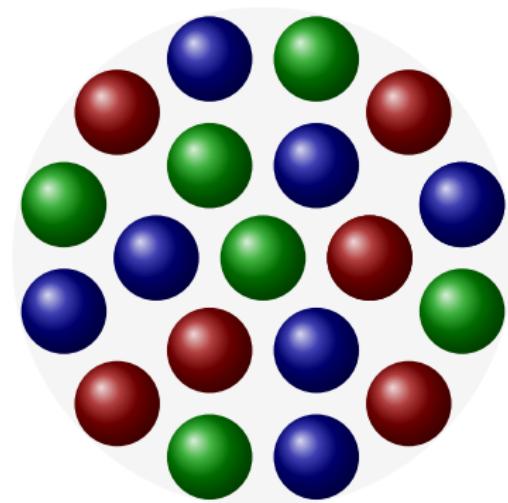


RANDOM INTERCEPT

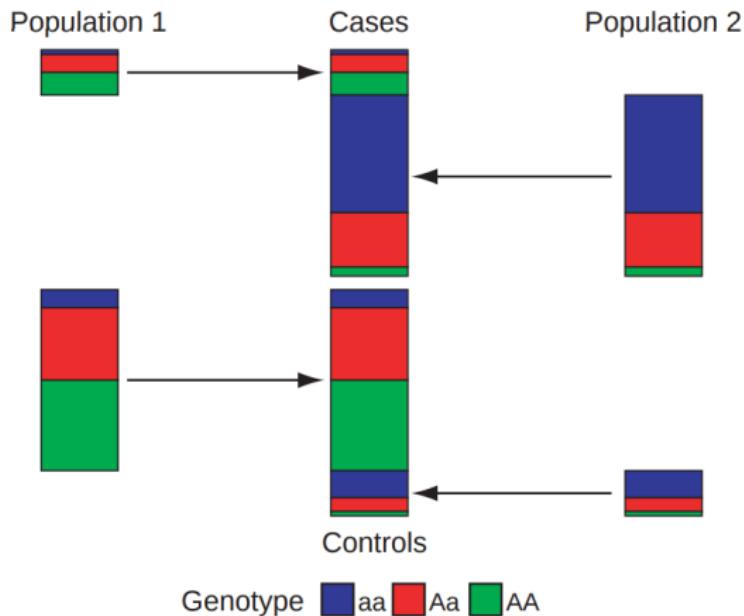


WHEN GROUPING INFORMATION IS UNKNOWN

- In our applications, the grouping information is unknown
- It must be estimated from the data



MOTIVATION



¹Marchini et al. Nature genetics (2004)