



LUND  
UNIVERSITY

# Optimization and Algorithms in Sparse Regression

Screening Rules, Coordinate Descent, and Normalization

---

Johan Larsson

Department of Statistics, Lund University

June 28, 2024

# Synopsis

## Overarching Theme

Optimization methods and numerical algorithms (screening rules) that help speed up the estimation of sparse regression models.

# Synopsis

## Overarching Theme

Optimization methods and numerical algorithms (screening rules) that help speed up the estimation of sparse regression models.

## Overview

**I–III** Screening rules for sparse regression

**IV** Benchmarking optimization methods

**V** Coordinate descent (an optimization method) for SLOPE

**VI** Normalization in penalized regression

## The Strong Screening Rule for SLOPE (Paper I)

---

# The Strong Screening Rule for SLOPE

Published and presented at NeurIPS 2020<sup>1</sup>

Co-written with supervisors Małgorzata Bogdan and Jonas Wallin.



---

<sup>1</sup>Johan Larsson, Małgorzata Bogdan, and Jonas Wallin (Dec. 6–12, 2020). “The Strong Screening Rule for SLOPE”. In: *Advances in Neural Information Processing Systems 33*. 34th Conference on Neural Information Processing Systems (NeurIPS 2020). Ed. by Hugo Larochelle et al. Vol. 33. Virtual: Curran Associates, Inc., pp. 14592–14603. ISBN: 978-1-71382-954-6

# Motivation

- Data sets are growing in size.
- Fitting models to sets of millions (maybe billions) of features is computationally costly.
- Hyperparameter tuning increases costs further.
- In sparse regression, we can solve a reduced problem if we somehow knew which features were active.
- Can we somehow figure this out **before** fitting the model?

# Screening Rules

## General Idea

1. If  $p \gg n$ , the solution will be sparse.<sup>2</sup>
2. Estimate feature importance before fitting (correlations, etc) and solve a reduced problem.
3. Maybe not so bad to miss some features—just add them back and update the fit.

---

<sup>2</sup>The lasso can for instance at most select  $\min(n, p)$  features.

# Screening Rules

## General Idea

1. If  $p \gg n$ , the solution will be sparse.<sup>2</sup>
2. Estimate feature importance before fitting (correlations, etc) and solve a reduced problem.
3. Maybe not so bad to miss some features—just add them back and update the fit.

Turn out to be a pretty good idea!

---

<sup>2</sup>The lasso can for instance at most select  $\min(n, p)$  features.

# Optimality Conditions and the Gradient

## Optimality Condition

$\beta$  is optimal if

$$\mathbf{0} \in \nabla g(\beta) + \partial h(\beta),$$

where  $\partial h(\beta)$  is the subdifferential of  
the penalty term  $h$ .

# Optimality Conditions and the Gradient

## Optimality Condition

$\beta$  is optimal if

$$\mathbf{0} \in \nabla g(\beta) + \partial h(\beta),$$

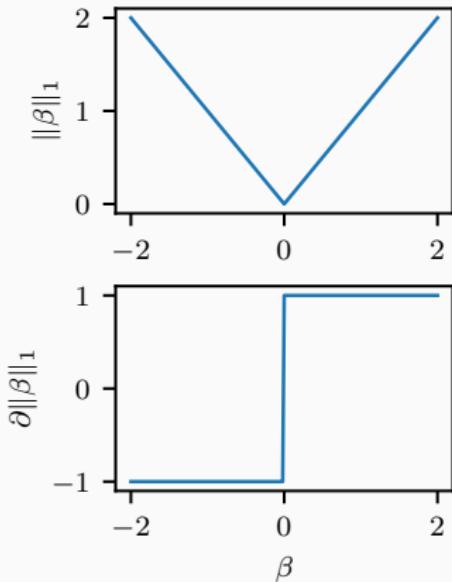
where  $\partial h(\beta)$  is the subdifferential of the penalty term  $h$ .

## The Lasso

Here we have

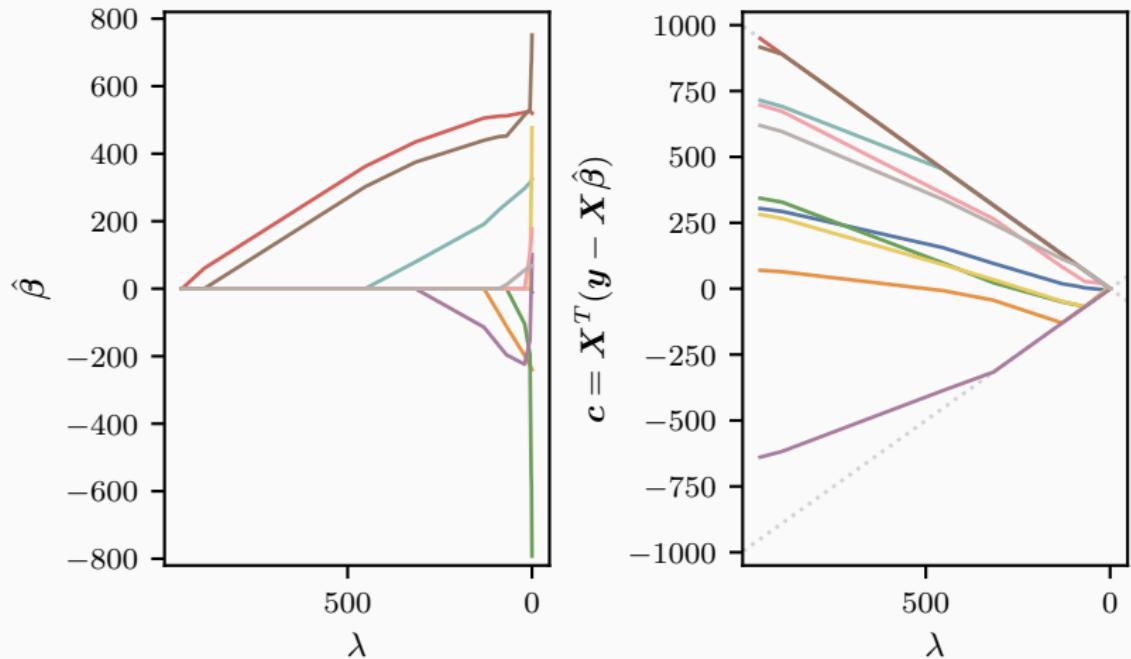
$$\partial h(\beta)_j = \begin{cases} \{\lambda \text{sign}(\beta_j)\} & \text{if } \beta_j \neq 0, \\ [-\lambda, \lambda] & \text{otherwise.} \end{cases}$$

Hence, if  $|\nabla g(\beta)_j| < \lambda$ , then  $\beta_j^* = 0$ .



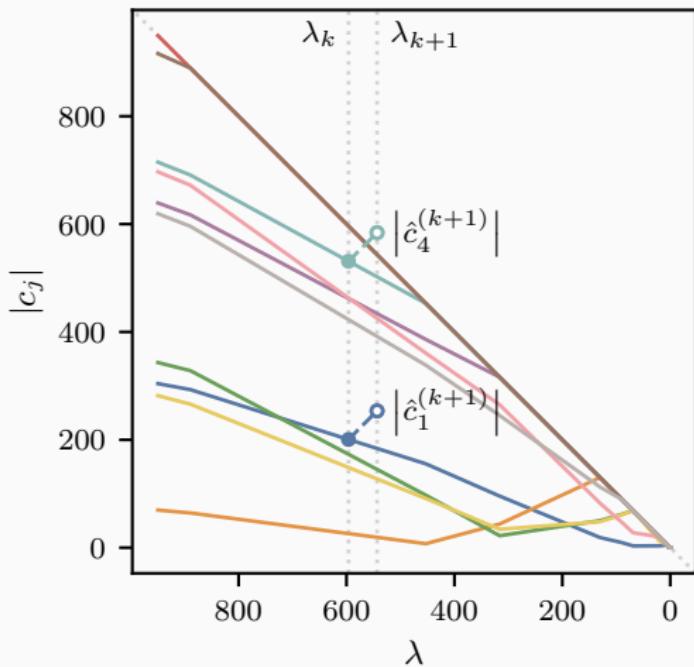
**Figure 1:** The lasso penalty ( $h(\beta) = \lambda\|\beta\|_1$ ) and its subdifferential

# Correlations and Coefficients



**Figure 2:** Coefficients and correlations along the lasso path

# Strong Rule for the Lasso



**Figure 3:** The strong rule for the lasso (Tibshirani et al. 2012)

# Strong Rule for SLOPE

Same idea as for lasso strong rule! Just need the subdifferential.

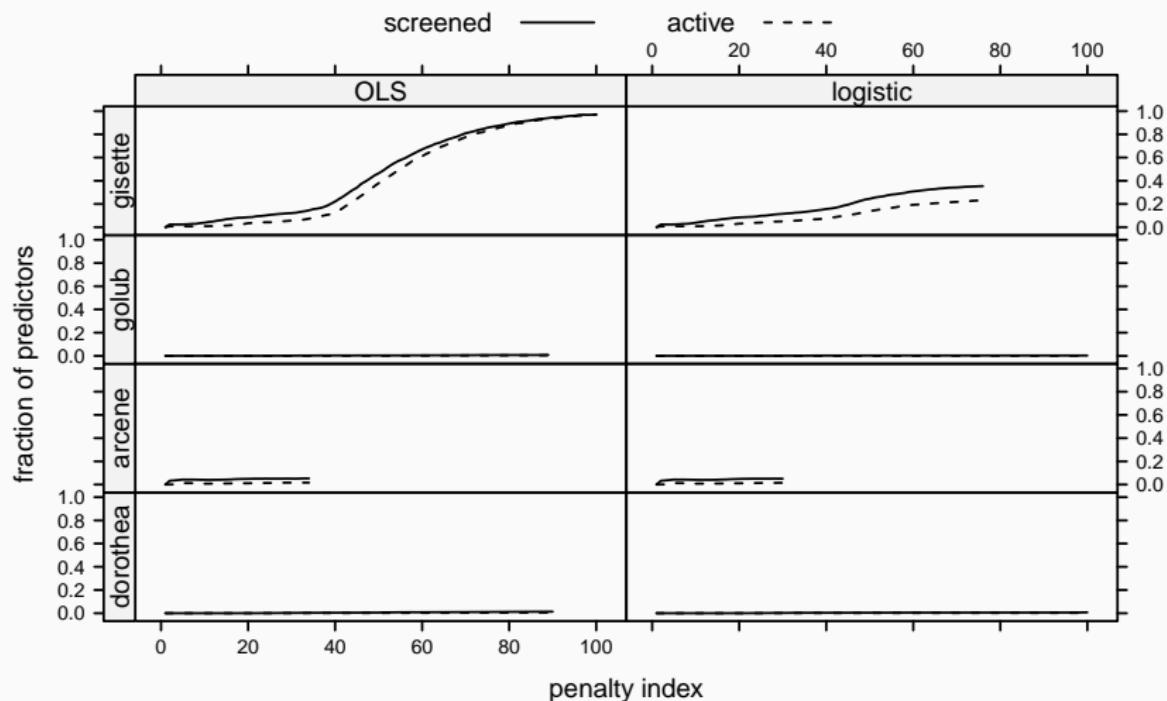
## SLOPE Subdifferential

The set of all  $g \in \mathbb{R}^p$  such that

$$g_{\mathcal{A}_i} = \left\{ s \in \mathbb{R}^{|\mathcal{A}_i|} \mid \begin{cases} \text{cumsum}(|s|_{\downarrow} - \lambda_{R(s)_{\mathcal{A}_i}}) \preceq \mathbf{0} & \text{if } \beta_{\mathcal{A}_i} = \mathbf{0}, \\ \text{cumsum}(|s|_{\downarrow} - \lambda_{R(s)_{\mathcal{A}_i}}) \preceq \mathbf{0} \\ \text{and } \text{sign}(\beta_{\mathcal{A}_i}) = \text{sign } s \\ \text{and } \sum_{j \in \mathcal{A}_i} (|s_j| - \lambda_{R(s)_j}) = 0 & \text{otherwise.} \end{cases} \right\}$$

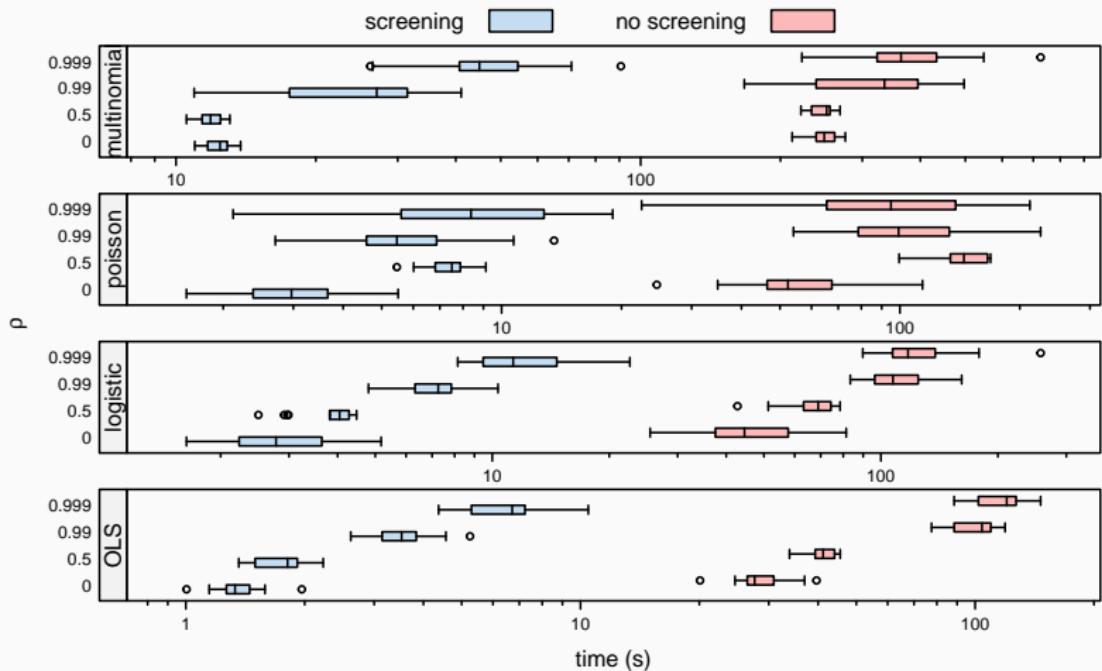
- $\mathcal{A}_i$  defines a **cluster** containing indices of coefficients equal in absolute value.
- $R(x)$  is an operator that returns the **ranks** of elements in  $|x|$ .
- $|x|_{\downarrow}$  returns the absolute values of  $x$  sorted in non-increasing order.

## Results: Effectiveness



**Figure 4:** Effectiveness for some real data sets

## Results: Speed



**Figure 5:** Time to fit a full SLOPE path with and without the strong rule

## Heuristic Screening and Violations

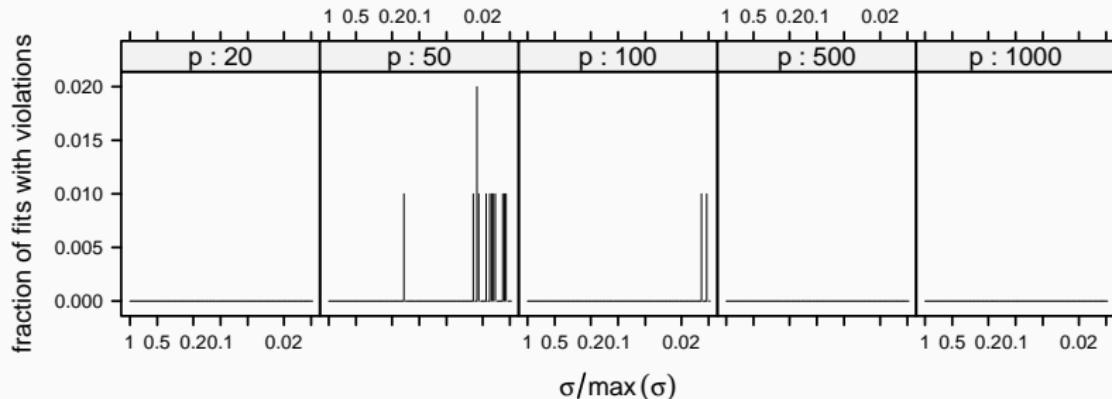
- The rule is **heuristic**<sup>3</sup>, which means that **violations** may occur.
- But they are so rare and checking for and refitting with them so cheap that it rarely matters.

---

<sup>3</sup>As opposed to *safe*.

# Heuristic Screening and Violations

- The rule is **heuristic**<sup>3</sup>, which means that **violations** may occur.
- But they are so rare and checking for and refitting with them so cheap that it rarely matters.



**Figure 6:** Violations of the strong rule for SLOPE

<sup>3</sup>As opposed to *safe*.

# Summary

## Contributions

- The first screening rule for SLOPE
- An efficient algorithm with which to apply it
- A new formulation of the subdifferential for SLOPE
- Implemented in the R package SLOPE<sup>4</sup>

---

<sup>4</sup><https://doi.org/10.32614/CRAN.package.SLOPE>

# Summary

## Contributions

- The first screening rule for SLOPE
- An efficient algorithm with which to apply it
- A new formulation of the subdifferential for SLOPE
- Implemented in the R package SLOPE<sup>4</sup>

## Limitations

- Somewhat conservative

---

<sup>4</sup><https://doi.org/10.32614/CRAN.package.SLOPE>

## Look-Ahead Screening Rules for the Lasso (Paper II)

---

# Look-Ahead Screening Rules for the Lasso

Published in EYSM 2021<sup>5</sup>



<sup>5</sup> Johan Larsson (Sept. 6, 2021). "Look-Ahead Screening Rules for the Lasso". In: *22nd European Young Statisticians Meeting - Proceedings*. 22nd European Young Statisticians Meeting. Ed. by Andreas Makridis et al. Athens, Greece: Panteion university of social and political sciences, pp. 61–65. ISBN: 978-960-7943-23-1

# Motivation

- Even if screening rules improve the speed remarkably, they still involve costs (e.g. gradient computations).

## Motivation

- Even if screening rules improve the speed remarkably, they still involve costs (e.g. gradient computations).
- Previous screening rules just look one-step ahead.

# Motivation

- Even if screening rules improve the speed remarkably, they still involve costs (e.g. gradient computations).
- Previous screening rules just look one-step ahead.
- But information from the current step is useful for more than the next one, so why not look further ahead?

## The Dual

Complementary to the primal problem. For the lasso, it is

$$\underset{\boldsymbol{\theta} \in \mathbb{R}^n}{\text{maximize}} \left( D(\boldsymbol{\theta}; \lambda) = \frac{1}{2} \mathbf{y}^\top \mathbf{y} - \frac{\lambda^2}{2} \left\| \boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda} \right\|_2^2 \right)$$

# The Dual

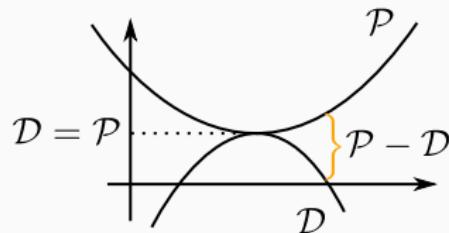
Complementary to the primal problem. For the lasso, it is

$$\underset{\boldsymbol{\theta} \in \mathbb{R}^n}{\text{maximize}} \left( D(\boldsymbol{\theta}; \lambda) = \frac{1}{2} \mathbf{y}^\top \mathbf{y} - \frac{\lambda^2}{2} \left\| \boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda} \right\|_2^2 \right)$$

## Duality Gap

Difference between the primal and dual objectives and is tight at the optimum, that is

$$P(\beta^*; \lambda) - D(\boldsymbol{\theta}^*; \lambda) = 0.$$



# The Dual

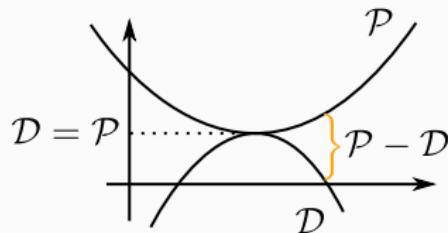
Complementary to the primal problem. For the lasso, it is

$$\underset{\boldsymbol{\theta} \in \mathbb{R}^n}{\text{maximize}} \left( D(\boldsymbol{\theta}; \lambda) = \frac{1}{2} \mathbf{y}^\top \mathbf{y} - \frac{\lambda^2}{2} \left\| \boldsymbol{\theta} - \frac{\mathbf{y}}{\lambda} \right\|_2^2 \right)$$

## Duality Gap

Difference between the primal and dual objectives and is tight at the optimum, that is

$$P(\beta^*; \lambda) - D(\boldsymbol{\theta}^*; \lambda) = 0.$$



## Dual–Primal Relationship

The two problems are related via

$$\mathbf{y} = \mathbf{X}\beta(\lambda) + \lambda\boldsymbol{\theta}(\lambda).$$

## Gap Safe Screening

Recall the optimality conditions for the lasso:

$$\text{if } |\nabla g(\boldsymbol{\beta}^*)_j| < \lambda \iff |\mathbf{x}_j^\top \boldsymbol{\theta}^*| < 1 \text{ then } \beta_j^* = 0.$$

# Gap Safe Screening

Recall the optimality conditions for the lasso:

$$\text{if } |\nabla g(\boldsymbol{\beta}^*)_j| < \lambda \iff |\mathbf{x}_j^\top \boldsymbol{\theta}^*| < 1 \text{ then } \beta_j^* = 0.$$

## Idea

Replace the inequality  $|\mathbf{x}_j^\top \boldsymbol{\theta}^*| < 1$  with an inequality that offers same conclusion, but doesn't involve  $\boldsymbol{\theta}^*$ .

# Gap Safe Screening

Recall the optimality conditions for the lasso:

$$\text{if } |\nabla g(\boldsymbol{\beta}^*)_j| < \lambda \iff |\mathbf{x}_j^\top \boldsymbol{\theta}^*| < 1 \text{ then } \beta_j^* = 0.$$

## Idea

Replace the inequality  $|\mathbf{x}_j^\top \boldsymbol{\theta}^*| < 1$  with an inequality that offers same conclusion, but doesn't involve  $\boldsymbol{\theta}^*$ .

Gap Safe screening (Fercoq, Gramfort, and Salmon 2015) uses the **duality gap** to define such a region, and discards the  $j$ th predictor if

$$|\mathbf{x}_j^\top \boldsymbol{\theta}_\lambda| + \|\mathbf{x}_j\|_2 \sqrt{\frac{1}{\lambda_+^2} (\mathcal{P}(\boldsymbol{\beta}_\lambda; \lambda_+) - \mathcal{D}(\boldsymbol{\theta}_\lambda; \lambda_+))} < 1$$

where  $\boldsymbol{\theta}_\lambda$  is the **scaled** dual,

$$\boldsymbol{\theta}_\lambda = \frac{\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_\lambda}{\max(\|\mathbf{X}^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_\lambda)\|_\infty, \lambda)}.$$

## Look-Ahead Screening

Inequality on last slide is **quadratic**, which means we can find the next critical point easily:

$$\lambda_+ = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

where

$$a = (1 - |\mathbf{x}_j^\top \boldsymbol{\theta}_\lambda|)^2 - \frac{1}{2} \boldsymbol{\theta}_\lambda^\top \boldsymbol{\theta}_\lambda \|\mathbf{x}_j\|_2^2,$$

$$b = (\boldsymbol{\theta}_\lambda^\top \mathbf{y} - \|\boldsymbol{\beta}_\lambda\|_1) \|\mathbf{x}_j\|_2^2,$$

$$c = -\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_\lambda\|_2^2 \|\mathbf{x}_j\|_2^2.$$

## Look-Ahead Screening

Inequality on last slide is **quadratic**, which means we can find the next critical point easily:

$$\lambda_+ = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

where

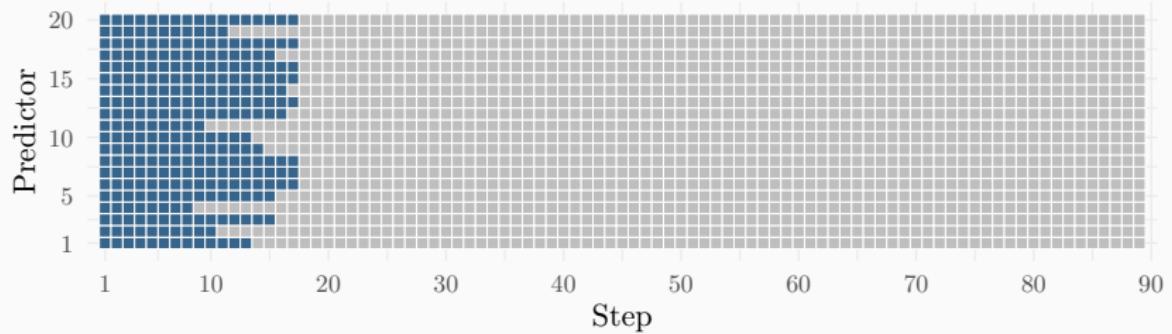
$$a = (1 - |\mathbf{x}_j^\top \boldsymbol{\theta}_\lambda|)^2 - \frac{1}{2} \boldsymbol{\theta}_\lambda^\top \boldsymbol{\theta}_\lambda \|\mathbf{x}_j\|_2^2,$$

$$b = (\boldsymbol{\theta}_\lambda^\top \mathbf{y} - \|\boldsymbol{\beta}_\lambda\|_1) \|\mathbf{x}_j\|_2^2,$$

$$c = -\frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_\lambda\|_2^2 \|\mathbf{x}_j\|_2^2.$$

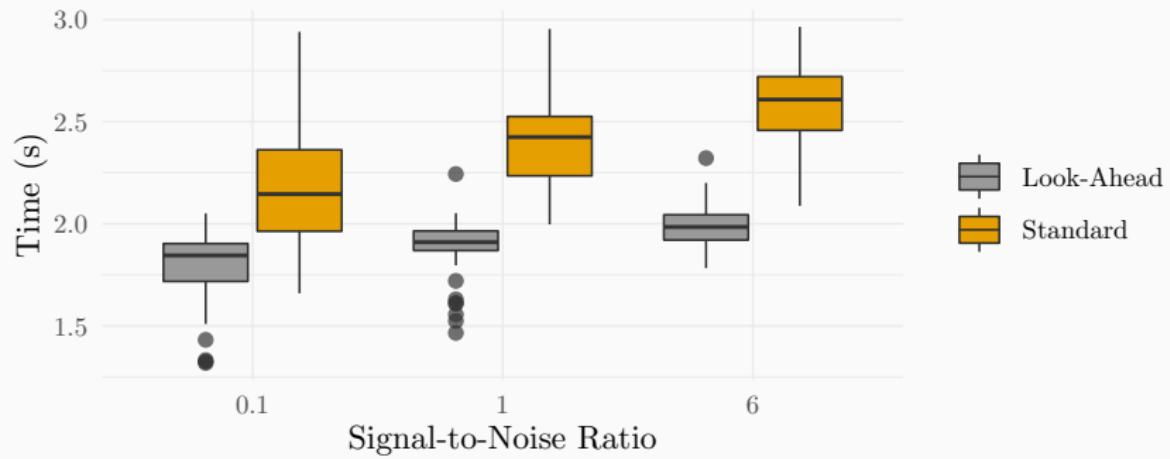
This allows us to screen predictors for **all** upcoming steps by simply checking the inequality for all  $\lambda_{k+1}, \lambda_{k+2}, \dots$

## Simple Example



**Figure 7:** The lasso path for a sample of 20 features from the leukemia data set. The squares show for which steps the feature can be discarded.

## Benchmarks



**Figure 8:** Benchmarks of time to fit the full lasso path with or without look-ahead screening.

# Summary

## Contributions

- Simple, safe, and essentially free screening rule for the lasso (and related problems)
- Moderate improvements in speed

# Summary

## Contributions

- Simple, safe, and essentially free screening rule for the lasso (and related problems)
- Moderate improvements in speed

## Limitations

- Needs a safe rule
- Most effective at start of path

## **The Hessian Screening Rule (Paper III)**

---

# The Hessian Screening Rule

Published and presented at  
NeurIPS 2022<sup>6</sup>



---

<sup>6</sup>Johan Larsson and Jonas Wallin (Nov. 28–Dec. 9, 2022). “The Hessian Screening Rule”. In: *Advances in Neural Information Processing Systems 35*. 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Ed. by Sanmi Koyejo et al. Vol. 35. New Orleans, USA: Curran Associates, Inc., pp. 15823–15835. ISBN: 978-1-71387-108-8

## Motivation

- Previous screening rules, like the strong rule, struggle with highly correlated features.
- Using second order information, we might both be able to improve screening, but also improve solving the problem.

# Screening Rules Seen As Gradient Estimates

Let  $c(\lambda) = -\nabla g(\beta(\lambda))$  be the correlation (negative gradient).

## Lasso Screening Rule Template

$0 \in \nabla g(\beta) + \lambda \partial$  suggests simple rule:

1. Replace  $c$  with an estimate  $\tilde{c}$ .
2. If  $|\tilde{c}_j| < \lambda$ , discard feature  $j$ .

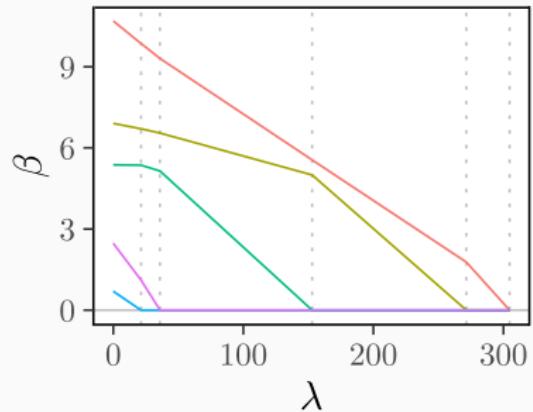
If  $\tilde{c}$  is accurate and not too conservative, we have a useful rule.

# The Ordinary Lasso

We now focus on the ordinary lasso,  $\ell_1$ -regularized least squares:

$$g(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2,$$

$$\nabla g(\beta) = \mathbf{X}^\top (\mathbf{X}\beta - \mathbf{y}).$$

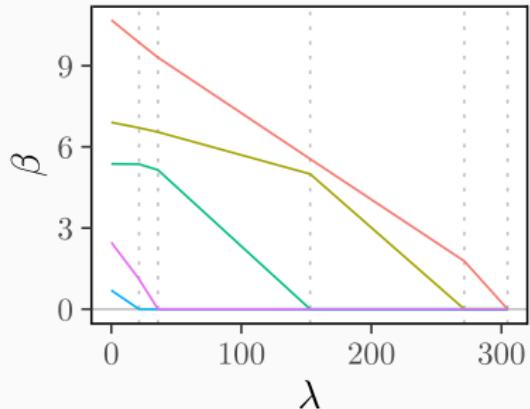


# The Ordinary Lasso

We now focus on the ordinary lasso,  $\ell_1$ -regularized least squares:

$$g(\beta) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\beta\|_2^2,$$

$$\nabla g(\beta) = \mathbf{X}^\top (\mathbf{X}\beta - \mathbf{y}).$$



It turns out that we can express the solution as a function of  $\lambda$ :

$$\hat{\beta}(\lambda) = (\mathbf{X}_{\mathcal{A}}^\top \mathbf{X}_{\mathcal{A}})^{-1} (\mathbf{X}_{\mathcal{A}}^\top \mathbf{y} - \lambda \text{sign}(\hat{\beta}_{\mathcal{A}})).$$

Holds for  $[\lambda_k, \lambda_{k+1}]$  if active set remains constant, and so

$$\hat{\beta}(\lambda_{k+1})_{\mathcal{A}} = \hat{\beta}(\lambda_k)_{\mathcal{A}} - (\lambda_k - \lambda_{k+1}) (\mathbf{X}_{\mathcal{A}}^\top \mathbf{X}_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}}).$$

# The Hessian Screening Rule

## Basic Form of the Rule

Take expression for  $\hat{\beta}$  on last slide and stick it into the gradient at step  $k + 1$ :

$$\begin{aligned}\tilde{\mathbf{c}}^H(\lambda_{k+1}) &= -\nabla g(\hat{\beta}(\lambda_{k+1})_{\mathcal{A}}) \\ &= \mathbf{c}(\lambda_k) + (\lambda_{k+1} - \lambda_k) \mathbf{X}^\top \mathbf{X}_{\mathcal{A}} (\mathbf{X}_{\mathcal{A}}^\top \mathbf{X}_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}}).\end{aligned}$$

# The Hessian Screening Rule

## Basic Form of the Rule

Take expression for  $\hat{\beta}$  on last slide and stick it into the gradient at step  $k + 1$ :

$$\begin{aligned}\tilde{\mathbf{c}}^H(\lambda_{k+1}) &= -\nabla g(\hat{\beta}(\lambda_{k+1})_{\mathcal{A}}) \\ &= \mathbf{c}(\lambda_k) + (\lambda_{k+1} - \lambda_k) \mathbf{X}^\top \mathbf{X}_{\mathcal{A}} (\mathbf{X}_{\mathcal{A}}^\top \mathbf{X}_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}}).\end{aligned}$$

- Exact expression for  $\mathbf{c}$  at  $k + 1$  if active set is constant.
- Heuristic rule, so needs checks for violations.

## Screening Rule Comparison

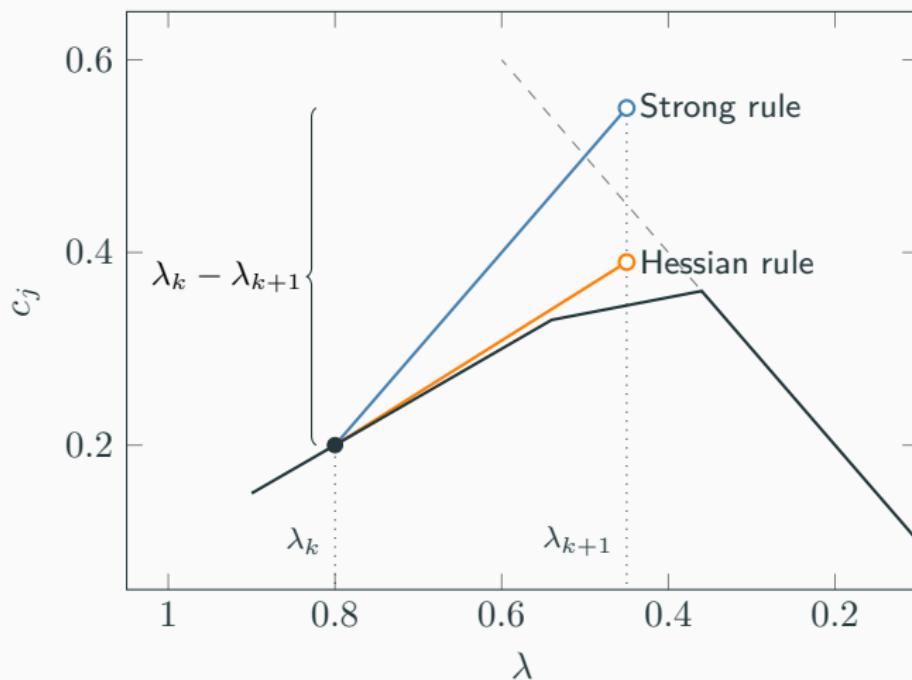
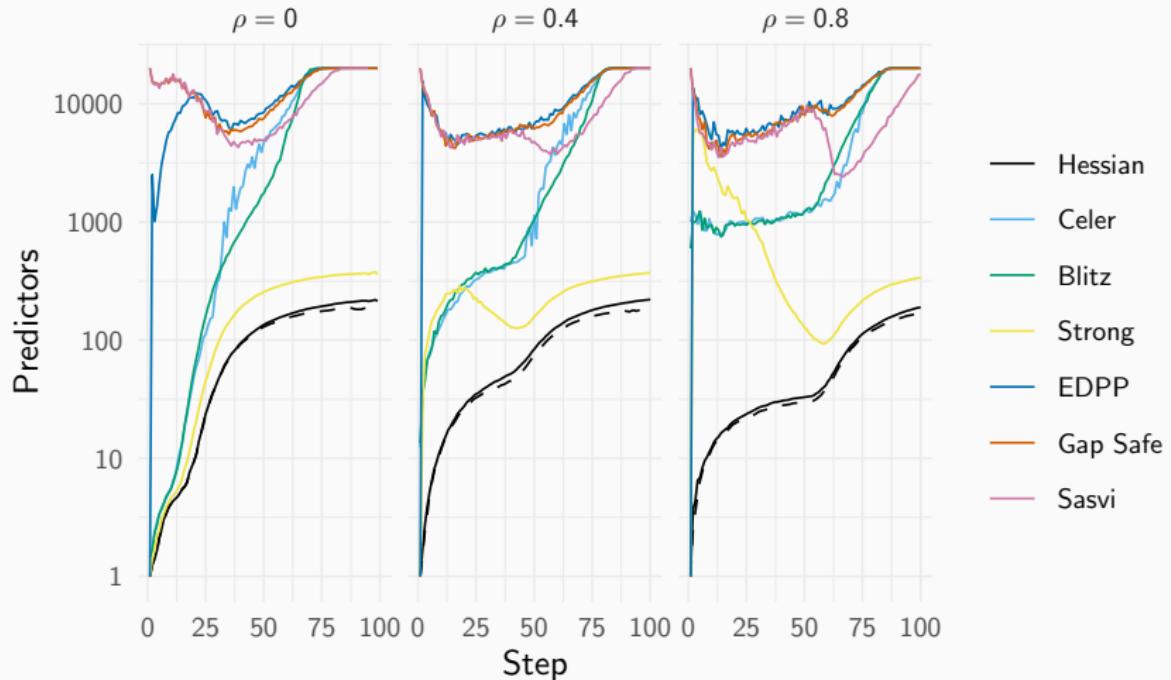


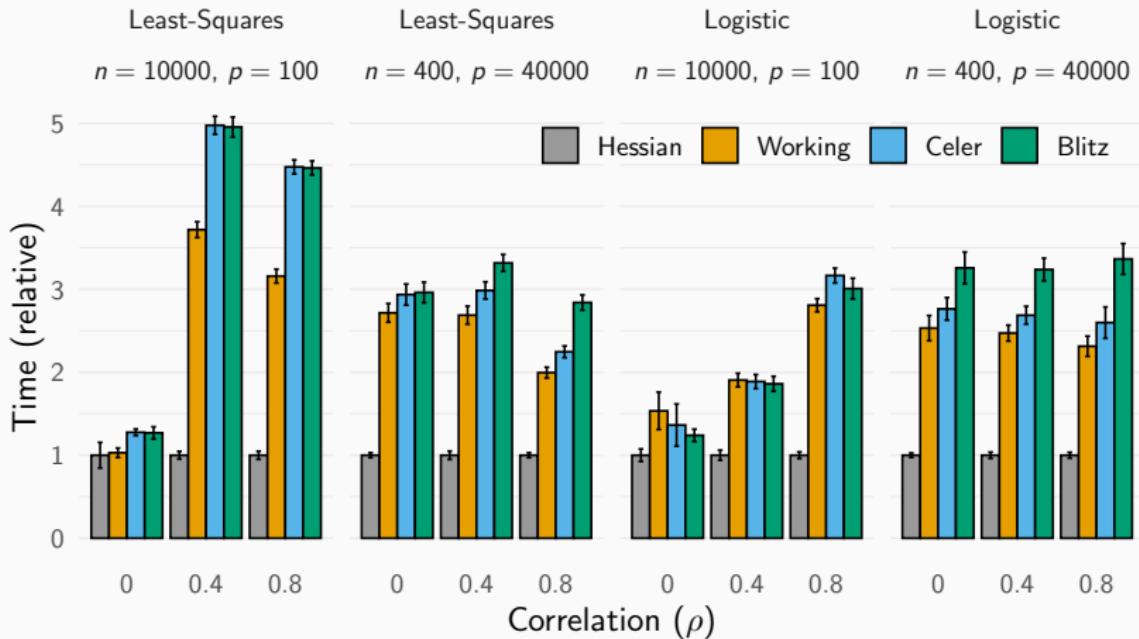
Figure 9: The strong and Hessian screening rules in action

## Results: Effectiveness



**Figure 10:** Features screened for a lasso path for  $\ell_1$ -regularized least-squares to a design with varying correlation ( $\rho$ ),  $n = 200$ , and  $p = 20000$ .

## Results: Simulated Data



**Figure 11:** Time to fit a full path of  $\ell_1$ -regularized least-squares and logistic regression to a design with  $n$  observations,  $p$  features, and pairwise correlation between features of  $\rho$ .

# Summary

## Contributions

- New screening rule that performs better in general and particularly for highly correlated features.
- Works well for both least-squares and logistic regression.
- Can be used in combination with approximate homotopy.

# Summary

## Contributions

- New screening rule that performs better in general and particularly for highly correlated features.
- Works well for both least-squares and logistic regression.
- Can be used in combination with approximate homotopy.

## Limitations

- Not well-tailored to problems with more complex Hessians.
- Implementation is a bit more involved.
- Although the rule works well with high correlation, speed-ups are not commensurate since our warm starts are worse.

# **Benchopt: Reproducible, Efficient and Collaborative Optimization Benchmarks (Paper IV)**

---

Published and presented at  
NeurIPS 2022<sup>7</sup>

Too many authors (20) to list  
here!



---

<sup>7</sup>Thomas Moreau et al. (Nov. 28–Dec. 9, 2022). “Benchopt: Reproducible, Efficient and Collaborative Optimization Benchmarks”. In: *Advances in Neural Information Processing Systems 35*. 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Ed. by S. Koyejo et al. Vol. 35. New Orleans, USA: Curran Associates, Inc., pp. 25404–25421. ISBN: 978-1-71387-108-8

## Motivation

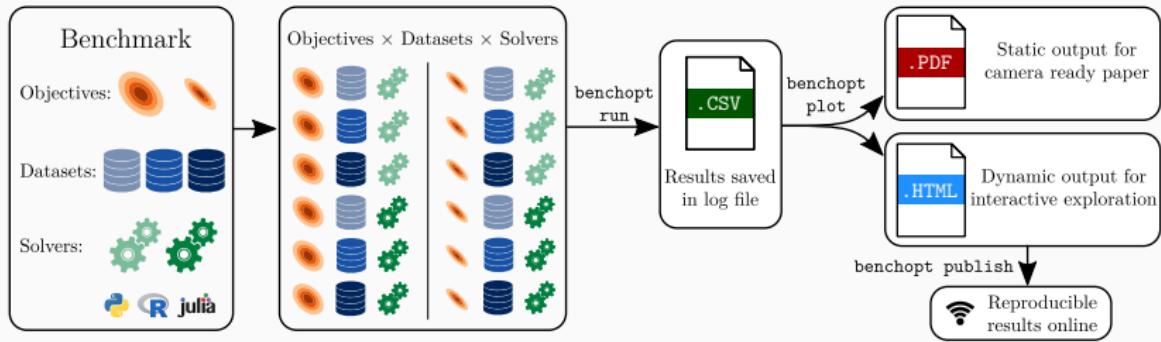
- Surging number of optimization methods
  - Best choice often depends on problem:
    - data set (dimensions, sparsity, conditioning, normalization)
    - hyper-parameters (regularization)
    - implementation (language)
  - Can we trust researchers to publish *fair* comparisons?

Benchpt attempts to help with all of this!

#### A: List of optimizers and schedules considered

Table 2: List of optimizers considered for our benchmark. This is only a subset of all existing methods for deep learning.

**Figure 12:** Methods benchmarked in Schmidt, Schneider, and Hennig (2021)



**Figure 13:** Schematic over how `benchopt` works

## Example: SLOPE

Install benchopt in an isolated environment.

```
1 conda create -n benchenv python  
2 conda activate benchenv  
3 pip install -U benchopt
```

## Example: SLOPE

Install `benchopt` in an isolated environment.

```
1 conda create -n benchenv python
2 conda activate benchenv
3 pip install -U benchopt
```

Install the benchmark (solvers and data sets).

```
1 benchopt install benchmark_slope -s rSLOPE -s PGD
```

## Example: SLOPE

Install `benchopt` in an isolated environment.

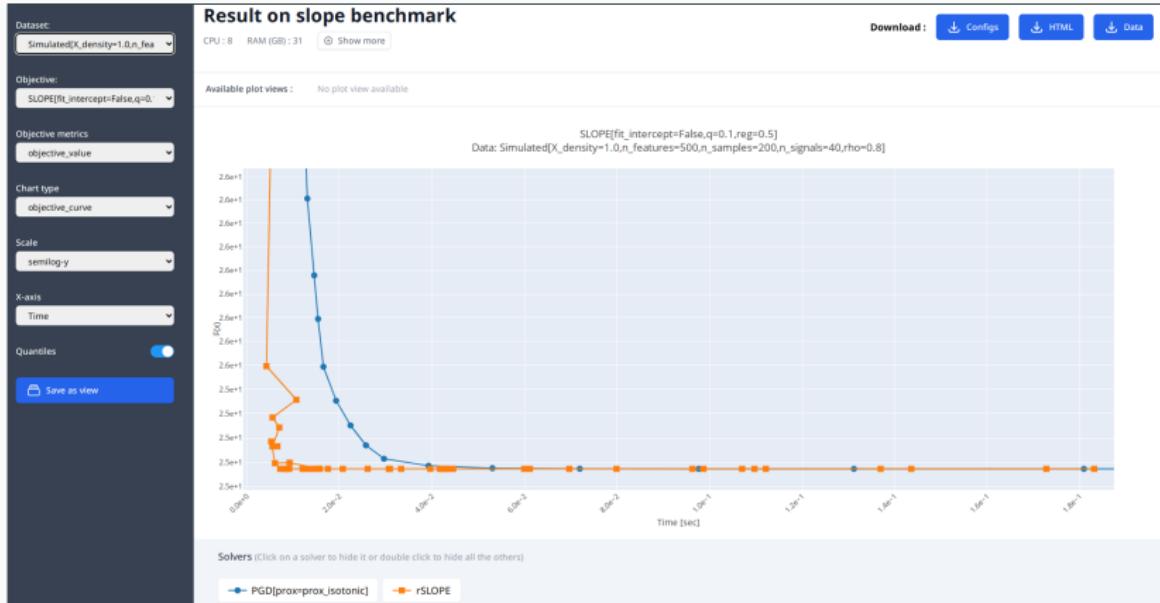
```
1 conda create -n benchenv python
2 conda activate benchenv
3 pip install -U benchopt
```

Install the benchmark (solvers and data sets).

```
1 benchopt install benchmark_slope -s rSLOPE -s PGD
```

And run it!

```
1 benchopt run benchmark_slope \
2   -o SLOPE[reg=0.5,fit_intercept=False,q=0.1] \
3   -s PGD[prox=prox_isotonic] -s rSLOPE \
4   -d Simulated[n_features=500,n_samples=200]
```



**Figure 14:** Benchopt benchmark results for SLOPE

# Benchmark Structure

A `benchopt` benchmark is a directory with

- an `objective.py` file with an `Objective`,
- a directory `solvers` with one file per solver, and
- a directory `datasets` with dataset generators/fetchers.

```
my_benchmark/
├── README.rst
└── datasets
    ├── simulated.py  # some dataset
    └── real.py  # some dataset
└── objective.py  # contains the definition of the objective
└── solvers
    ├── solver1.py  # some solver
    └── solver2.py  # some solver
```

# Summary

## Contributions

A framework that simplifies benchmarking of optimization methods:

- Automatic installation of dependencies
- Parametrized datasets, objectives and solvers
- Caching
- Interactive visualizations and easy publishing of results
- Mostly language agnostic

# Summary

## Contributions

A framework that simplifies benchmarking of optimization methods:

- Automatic installation of dependencies
- Parametrized datasets, objectives and solvers
- Caching
- Interactive visualizations and easy publishing of results
- Mostly language agnostic

## Limitations

- Doesn't directly solve the problem of varying implementations.
- Dealing with dependencies for 20+ solvers is still challenging.

## Coordinate Descent for SLOPE (Paper V)

---

# Coordinate Descent for SLOPE

Published and presented at AISTATS 2023<sup>8</sup>



Co-authored with Mathurin Massias, Quentin Bertrand, and Jonas Wallin.

---

<sup>8</sup>Johan Larsson, Quentin Klopfenstein, et al. (Apr. 25–27, 2023). “Coordinate Descent for SLOPE”. In: *Proceedings of the 26th International Conference on Artificial Intelligence and Statistics*. AISTATS 2023. Ed. by Francisco Ruiz, Jennifer Dy, and Jan-Willem van de Meent. Vol. 206. Proceedings of Machine Learning Research. Valencia, Spain: PMLR, pp. 4802–4821

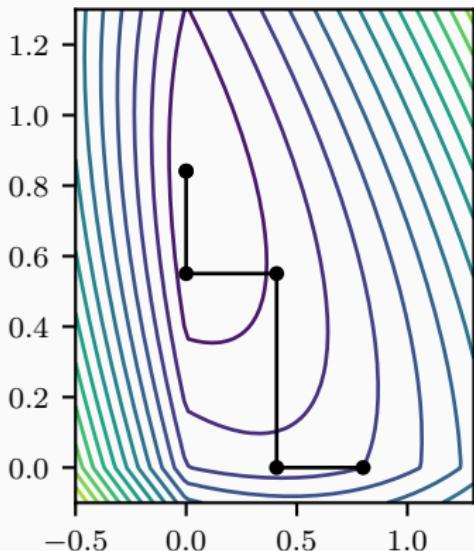
# Motivation

- SLOPE has many appealing properties but the best algorithms for solving SLOPE are slow (relative to the lasso)
- Lasso solvers are fast because they use coordinate descent
- Unfortunately, we cannot use coordinate descent directly for SLOPE
- But what if we could combine coordinate descent with another method?

# Coordinate Descent

## Simple Method!

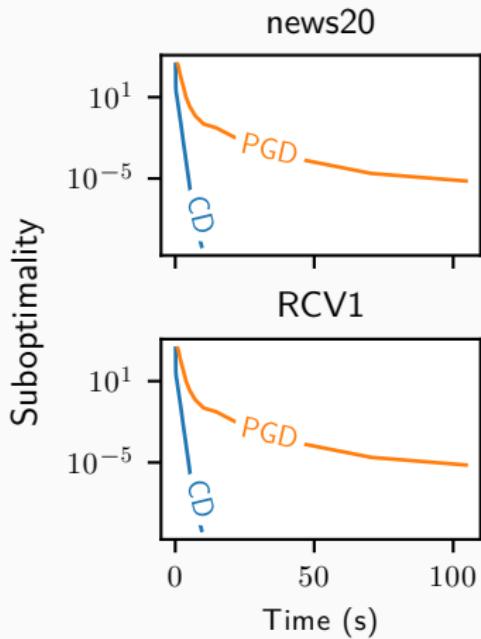
At each iteration, update a single coordinate (coefficient).



**Figure 15:** A basic example of coordinate descent in two dimensions

# Coordinate Descent Performance

Performs (surprisingly) well!

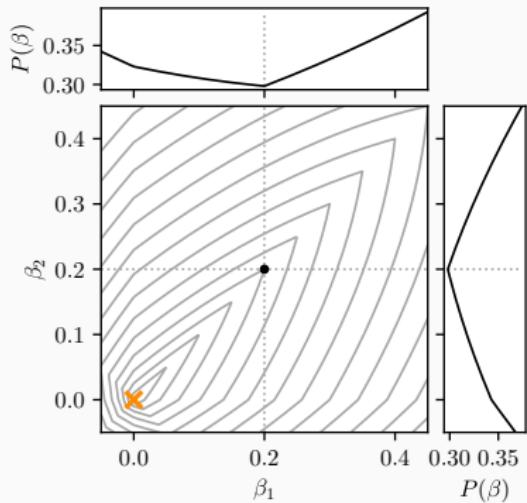


**Figure 16:** Coordinate descent versus proximal gradient descent for the lasso.

# Coordinate Descent and Separability

Cannot use basic coordinate descent for SLOPE since the sorted  $\ell_1$  norm is **inseparable**:

$$J(\beta) = \sum_{j=1}^p \lambda_j |\beta_{(j)}|.$$



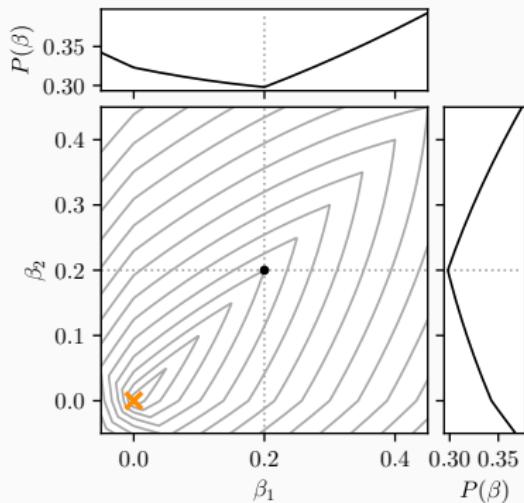
**Figure 17:** A *naive* coordinate descent algorithm cannot advance from the current iterate (●) to reach the optimum (✖).

# Coordinate Descent and Separability

Cannot use basic coordinate descent for SLOPE since the sorted  $\ell_1$  norm is **inseparable**:

$$J(\beta) = \sum_{j=1}^p \lambda_j |\beta_{(j)}|.$$

But if we could fix the clusters, we have separability!



**Figure 17:** A *naive* coordinate descent algorithm cannot advance from the current iterate (●) to reach the optimum (✖).

# Coordinate Descent and Separability

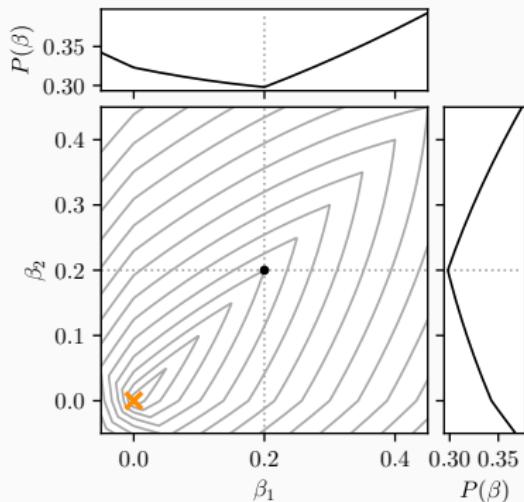
Cannot use basic coordinate descent for SLOPE since the sorted  $\ell_1$  norm is **inseparable**:

$$J(\beta) = \sum_{j=1}^p \lambda_j |\beta_{(j)}|.$$

But if we could fix the clusters, we have separability!

## Idea

Alternate between gradient descent and coordinate descent **on the clusters**.



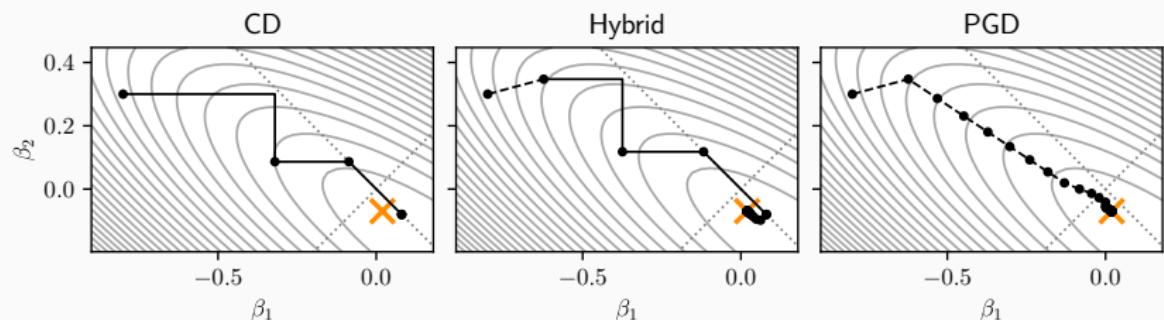
**Figure 17:** A *naive* coordinate descent algorithm cannot advance from the current iterate (●) to reach the optimum (✖).

## Hybrid Algorithm

- Every  $v$ th iteration, take a full proximal gradient step. This allows clusters to split (or merge).
- At all other iterations, take coordinate descent steps on the clusters.

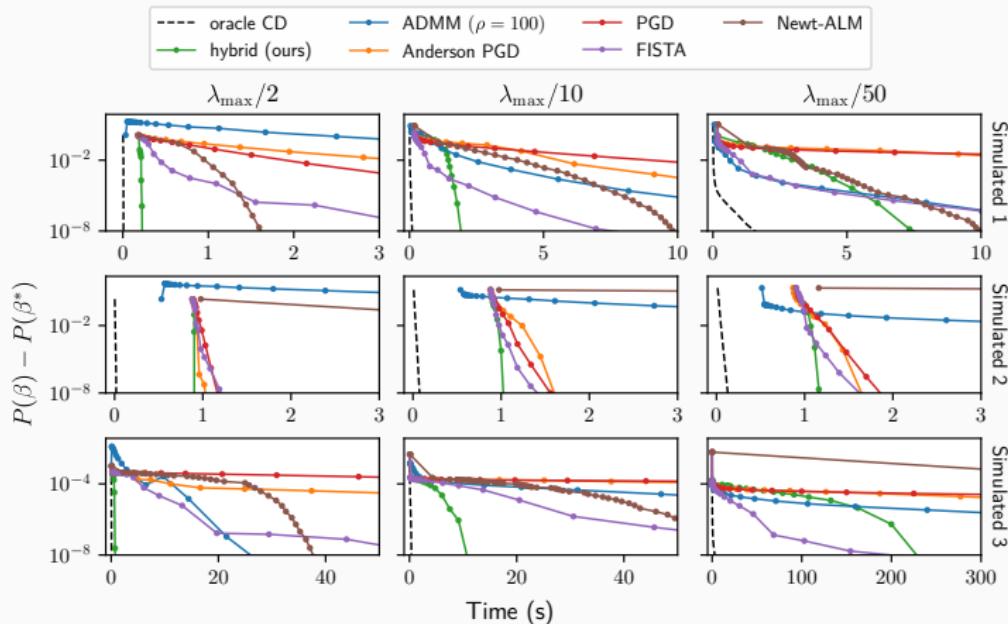
## Hybrid Algorithm

- Every  $v$ th iteration, take a full proximal gradient step. This allows clusters to split (or merge).
- At all other iterations, take coordinate descent steps on the clusters.



**Figure 18:** Our algorithm (hybrid) is a combination of CD and PGD.

# Experiments: Simulated Data



**Figure 19:** Benchmarks on simulated data. Scenario 1:  $n = 200$  and  $p = 20\,000$ ,  $X$ . Scenario 2:  $n = 20\,000$  and  $p = 200$ . Scenario 3:  $n = 200$ ,  $p = 200\,000$ , and sparse  $X$ .

# Summary

## Contributions

- A hybrid algorithm for SLOPE that brings the speed of coordinate descent to SLOPE.
- Faster than existing methods
- Implemented in Python package sortedl1<sup>9</sup>

---

<sup>9</sup><https://pypi.org/project/sortedl1/>

# Summary

## Contributions

- A hybrid algorithm for SLOPE that brings the speed of coordinate descent to SLOPE.
- Faster than existing methods
- Implemented in Python package sortedl1<sup>9</sup>

## Limitations

- Implementation is a bit complex.
- Not quite as fast as coordinate descent solvers for the lasso.

---

<sup>9</sup><https://pypi.org/project/sortedl1/>

# **The Lasso and Ridge Regression Yield Biased Estimates of Imbalanced Binary Features (Paper VI)**

---

# The Lasso and Ridge Regression Yield Biased Estimates of Imbalanced Binary Features

Not yet published, but will be soon™

Co-authored with Jonas.

# Motivation

- Most regularized methods (e.g. SLOPE and lasso) are scale-sensitive

# Motivation

- Most regularized methods (e.g. SLOPE and lasso) are scale-sensitive
- Straightforward when everything is normal, but about features that have other distributions, such as binary features?

# Motivation

- Most regularized methods (e.g. SLOPE and lasso) are scale-sensitive
- Straightforward when everything is normal, but about features that have other distributions, such as binary features?
- No literature on the effects of different normalization strategies

# The Elastic Net

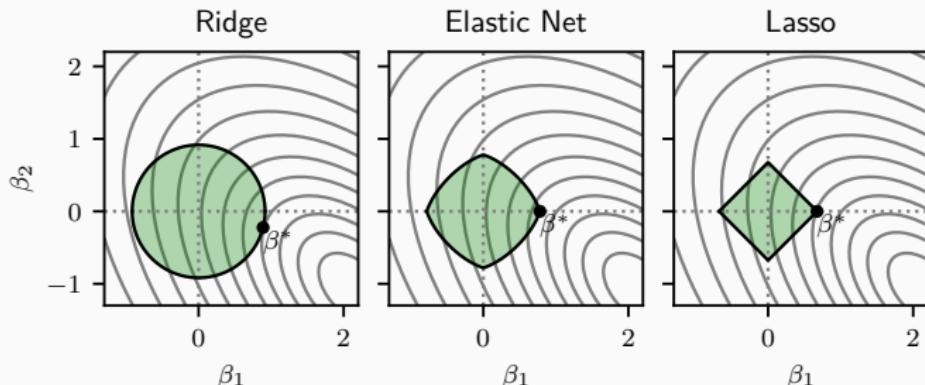
Linear regression plus a combination of the  $\ell_1$  and  $\ell_2$  penalties:

$$(\hat{\beta}_0, \hat{\boldsymbol{\beta}}) = \underset{\beta_0 \in \mathbb{R}, \boldsymbol{\beta} \in \mathbb{R}^p}{\arg \min} \left( \frac{1}{2} \|\mathbf{y} - \beta_0 - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda_1 \|\boldsymbol{\beta}\|_1 + \frac{\lambda_2}{2} \|\boldsymbol{\beta}\|_2^2 \right).$$

# The Elastic Net

Linear regression plus a combination of the  $\ell_1$  and  $\ell_2$  penalties:

$$(\hat{\beta}_0, \hat{\beta}) = \arg \min_{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^p} \left( \frac{1}{2} \|\mathbf{y} - \beta_0 - \mathbf{X}\beta\|_2^2 + \lambda_1 \|\beta\|_1 + \frac{\lambda_2}{2} \|\beta\|_2^2 \right).$$



**Figure 20:** The elastic net penalty is a combination of the lasso and ridge penalties. Here shown as a constrained problem.

## Sensitivity to Scale

Lasso and ridge penalties are **norms** of the coefficients—feature scales matter!

# Sensitivity to Scale

Lasso and ridge penalties are **norms** of the coefficients—feature scales matter!

## Example

Assume

$$\mathbf{X} \sim \text{Normal} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \right), \quad \beta^* = \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix}.$$

# Sensitivity to Scale

Lasso and ridge penalties are **norms** of the coefficients—feature scales matter!

## Example

Assume

$$\mathbf{X} \sim \text{Normal} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \right), \quad \beta^* = \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix}.$$

Model	$\hat{\beta}$	$\hat{\beta}_{\text{std}}$
OLS	$[0.50 \quad 1.00]^T$	$[1.00 \quad 1.00]^T$
Lasso	$[0.38 \quad 0.50]^T$	$[0.74 \quad 0.50]^T$
Ridge	$[0.37 \quad 0.41]^T$	$[0.74 \quad 0.41]^T$

# Sensitivity to Scale

Lasso and ridge penalties are **norms** of the coefficients—feature scales matter!

## Example

Assume

$$\mathbf{X} \sim \text{Normal} \left( \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \right), \quad \beta^* = \begin{bmatrix} \frac{1}{2} \\ 1 \end{bmatrix}.$$

Model	$\hat{\beta}$	$\hat{\beta}_{\text{std}}$
OLS	$[0.50 \quad 1.00]^T$	$[1.00 \quad 1.00]^T$
Lasso	$[0.38 \quad 0.50]^T$	$[0.74 \quad 0.50]^T$
Ridge	$[0.37 \quad 0.41]^T$	$[0.74 \quad 0.41]^T$

**Large** scale means **less** penalization because the size of  $\beta_j$  can be smaller for an equivalent effect (on  $y$ ).

# Normalization

- Scale sensitivity can be mitigated by normalizing the features.

## Normalization

- Scale sensitivity can be mitigated by normalizing the features.
- Let  $\tilde{\mathbf{X}}$  be the normalized feature matrix, with elements

$$\tilde{x}_{ij} = \frac{x_{ij} - c_j}{s_j},$$

where  $x_{ij}$  is an element of the (unnormalized) feature matrix and  $c_j$  and  $s_j$  are the *centering* and *scaling* factors respectively.

## Normalization

- Scale sensitivity can be mitigated by normalizing the features.
- Let  $\tilde{\mathbf{X}}$  be the normalized feature matrix, with elements

$$\tilde{x}_{ij} = \frac{x_{ij} - c_j}{s_j},$$

where  $x_{ij}$  is an element of the (unnormalized) feature matrix and  $c_j$  and  $s_j$  are the *centering* and *scaling* factors respectively.

- Usage of key terms ambiguous in literature.

## Normalization

- Scale sensitivity can be mitigated by normalizing the features.
- Let  $\tilde{\mathbf{X}}$  be the normalized feature matrix, with elements

$$\tilde{x}_{ij} = \frac{x_{ij} - c_j}{s_j},$$

where  $x_{ij}$  is an element of the (unnormalized) feature matrix and  $c_j$  and  $s_j$  are the *centering* and *scaling* factors respectively.

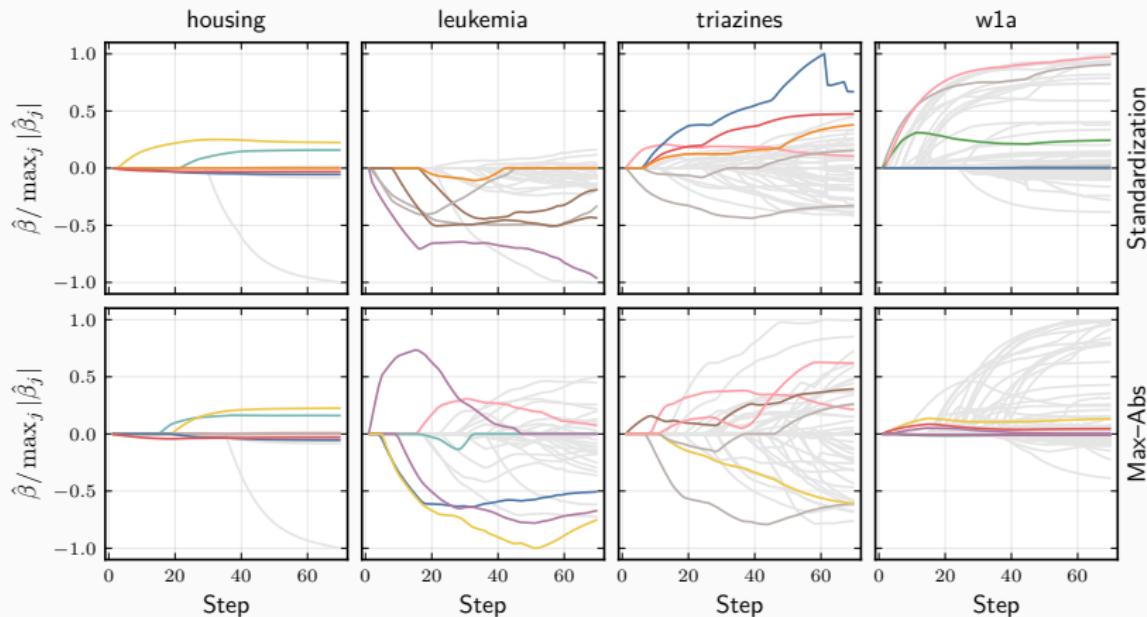
- Usage of key terms ambiguous in literature.
- After fitting, we transform the coefficients back to their original scale via

$$\hat{\beta}_j = \frac{\hat{\beta}_j^{(n)}}{s_j} \quad \text{for } j = 1, 2, \dots, p.$$

**Table 1:** Common ways to normalize  $\mathbf{X}$

Normalization	Centering ( $c_j$ )	Scaling ( $s_j$ )
Standardization	$\frac{1}{n} \sum_{i=1}^n x_{ij}$	$\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}$
Min–Max	$\min_i(x_{ij})$	$\max_i(x_{ij}) - \min_i(x_{ij})$
Unit Vector (L2)	0	$\sqrt{\sum_{i=1}^n x_{ij}^2}$
Max–Abs	0	$\max_i( x_{ij} )$
Adaptive Lasso	0	$\beta_j^{\text{OLS}}$

# The Type of Normalization Matters



**Figure 21:** Lasso paths under two different types of normalization (standardization and max-abs normalization). The union of the first ten features selected in any of the schemes are colored.

## Binary Features

For binary features (values 0 and 1 only), we have for the noiseless case

$$\hat{\beta}_j = \frac{S_{\lambda_1}(\tilde{\mathbf{x}}^\top \mathbf{y})}{s_j (\tilde{\mathbf{x}}_j^\top \tilde{\mathbf{x}}_j + \lambda_2)} = \frac{S_{\lambda_1} \left( \frac{\beta_j^* n(q-q^2)}{s_j} \right)}{s_j \left( \frac{n(q-q^2)}{s_j^2} + \lambda_2 \right)}.$$

## Binary Features

For binary features (values 0 and 1 only), we have for the noiseless case

$$\hat{\beta}_j = \frac{S_{\lambda_1}(\tilde{\mathbf{x}}^\top \mathbf{y})}{s_j (\tilde{\mathbf{x}}_j^\top \tilde{\mathbf{x}}_j + \lambda_2)} = \frac{S_{\lambda_1} \left( \frac{\beta_j^* n(q-q^2)}{s_j} \right)}{s_j \left( \frac{n(q-q^2)}{s_j^2} + \lambda_2 \right)}.$$

- Means that the elastic net estimator depends on class balance ( $q$ ).

## Binary Features

For binary features (values 0 and 1 only), we have for the noiseless case

$$\hat{\beta}_j = \frac{S_{\lambda_1}(\tilde{\mathbf{x}}^\top \mathbf{y})}{s_j (\tilde{\mathbf{x}}_j^\top \tilde{\mathbf{x}}_j + \lambda_2)} = \frac{S_{\lambda_1} \left( \frac{\beta_j^* n(q-q^2)}{s_j} \right)}{s_j \left( \frac{n(q-q^2)}{s_j^2} + \lambda_2 \right)}.$$

- Means that the elastic net estimator depends on class balance ( $q$ ).
- To remove the effect of  $q$ , an intuitive setting would be  $s_j = q - q^2$  (the variance of  $\mathbf{x}_j$ ) in the case of the lasso and  $s_j = \sqrt{q - q^2}$  (the standard deviation of  $\mathbf{x}_j$ ) in the case of the ridge.

## Binary Features

For binary features (values 0 and 1 only), we have for the noiseless case

$$\hat{\beta}_j = \frac{S_{\lambda_1}(\tilde{\mathbf{x}}^\top \mathbf{y})}{s_j (\tilde{\mathbf{x}}_j^\top \tilde{\mathbf{x}}_j + \lambda_2)} = \frac{S_{\lambda_1} \left( \frac{\beta_j^* n(q-q^2)}{s_j} \right)}{s_j \left( \frac{n(q-q^2)}{s_j^2} + \lambda_2 \right)}.$$

- Means that the elastic net estimator depends on class balance ( $q$ ).
- To remove the effect of  $q$ , an intuitive setting would be  $s_j = q - q^2$  (the variance of  $\mathbf{x}_j$ ) in the case of the lasso and  $s_j = \sqrt{q - q^2}$  (the standard deviation of  $\mathbf{x}_j$ ) in the case of the ridge.
- Suggests the parametrization

$$s_j = (q - q^2)^\delta, \quad \delta \geq 0.$$

## Binary Features

For binary features (values 0 and 1 only), we have for the noiseless case

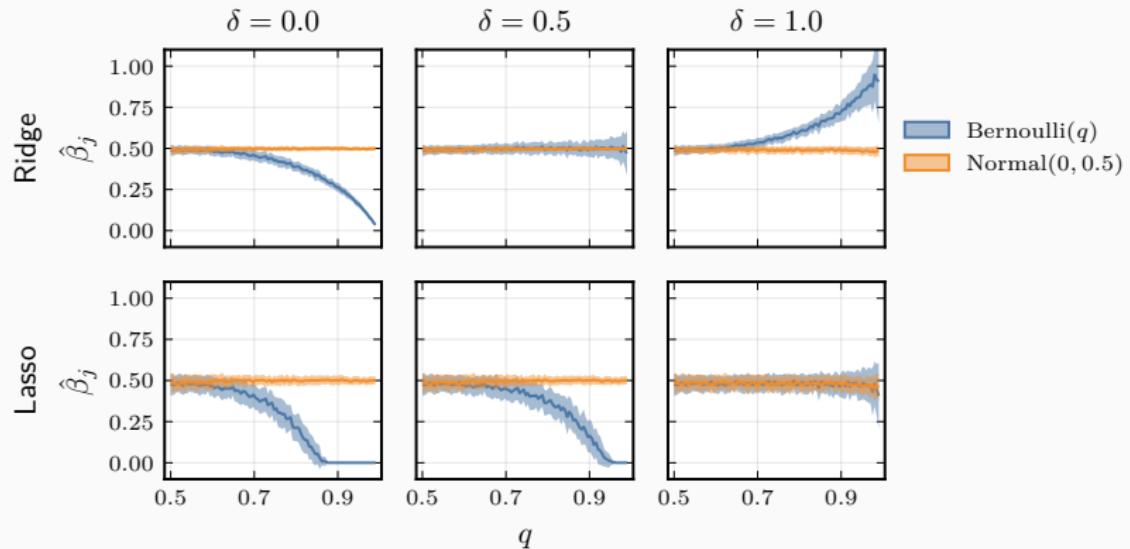
$$\hat{\beta}_j = \frac{S_{\lambda_1}(\tilde{x}^\top \mathbf{y})}{s_j (\tilde{x}_j^\top \tilde{x}_j + \lambda_2)} = \frac{S_{\lambda_1} \left( \frac{\beta_j^* n(q-q^2)}{s_j} \right)}{s_j \left( \frac{n(q-q^2)}{s_j^2} + \lambda_2 \right)}.$$

- Means that the elastic net estimator depends on class balance ( $q$ ).
- To remove the effect of  $q$ , an intuitive setting would be  $s_j = q - q^2$  (the variance of  $x_j$ ) in the case of the lasso and  $s_j = \sqrt{q - q^2}$  (the standard deviation of  $x_j$ ) in the case of the ridge.
- Suggests the parametrization

$$s_j = (q - q^2)^\delta, \quad \delta \geq 0.$$

- Indicates there might be no (simple)  $s_j$  that works for the elastic net.  
(Need to use weighted elastic net.)

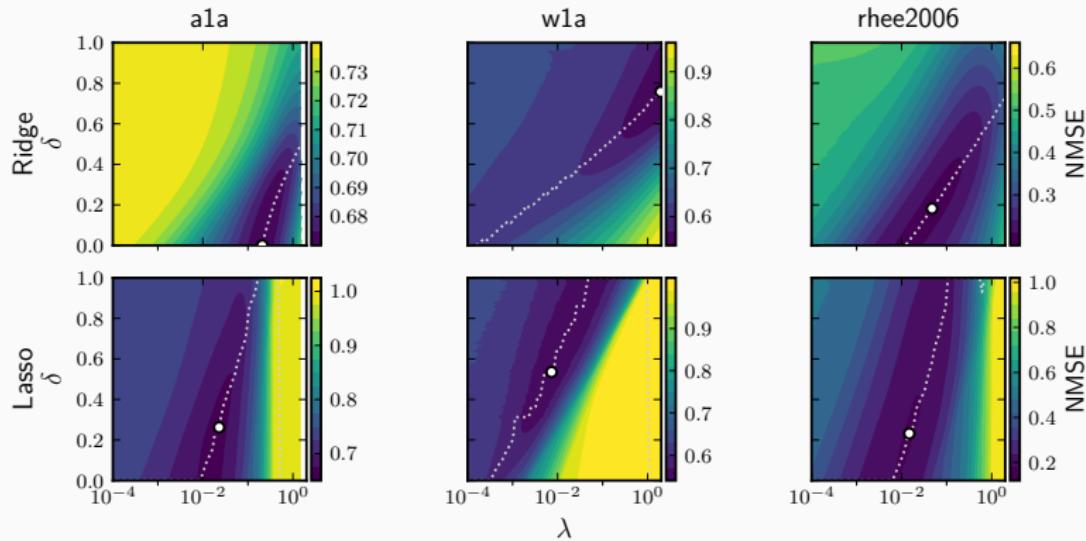
# Mixed Data



**Figure 22:** Comparison between lasso and ridge estimators for a data set with one binary and one quasi-normal feature.

# Hyperparameter Optimization

**Idea:** The choice of  $\delta$  affects the model, so let's optimize over it.



**Figure 23:** Contour plots of hold-out (validation set) error across a grid of  $\delta$  and  $\lambda$  values for the lasso and ridge.

# Summary

## Contributions

- As far as we know the first paper to investigate the interplay between normalization and regularization
- New scaling approach to deal with class-imbalanced binary features
- Discussion and suggestions for dealing with mixed data

# Summary

## Contributions

- As far as we know the first paper to investigate the interplay between normalization and regularization
- New scaling approach to deal with class-imbalanced binary features
- Discussion and suggestions for dealing with mixed data

## Limitations

- So far only theoretical results for limited cases:
  - Fixed data ( $X$ ), normal noise
  - Orthogonal features
  - Normal and binary features
- Only experimental results for interactions

# Recap

**Thank you!**