

# The Strong Screening Rule for SLOPE

## Mathematical Methods of Modern Statistics 2

Johan Larsson<sup>1</sup>   Małgorzata Bogdan<sup>1,2</sup>   Jonas Wallin<sup>1</sup>

<sup>1</sup>Department of Statistics, Lund University,

<sup>2</sup>Department of Mathematics, University of Wrocław

June 16, 2020



**LUND**  
UNIVERSITY

# Sorted L-One Penalized Estimation (SLOPE)

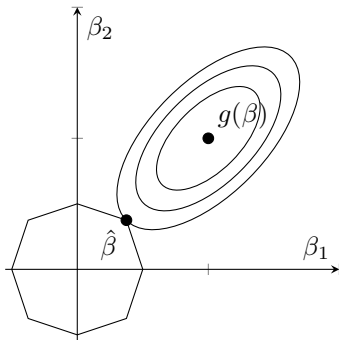
The SLOPE (bogdan2015) estimate is

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^p} \{g(\beta) + J(\beta; \lambda)\}$$

where  $J(\beta; \lambda) = \sum_{i=1}^p \lambda_i |\beta|_{(i)}$  is the **sorted**  $\ell_1$  **norm**, where

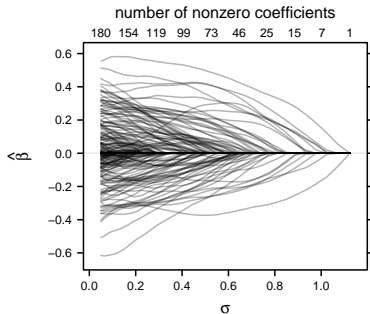
$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_p \geq 0, \quad |\beta|_{(1)} \geq |\beta|_{(2)} \geq \cdots \geq |\beta|_{(p)}.$$

Equivalent to an  
inequality-constrained convex  
optimization problem



# Motivation for screening rules

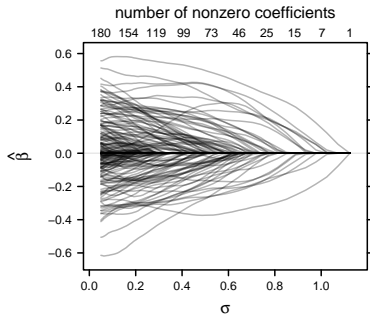
- we are interested in a **path** of penalties  $\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(m)}$ , many of which will lead to **sparse** solutions



**Figure 1:** SLOPE path with  $n = 200$ ,  $p = 20000$ .  $\sigma$  indicates strength of regularization.

# Motivation for screening rules

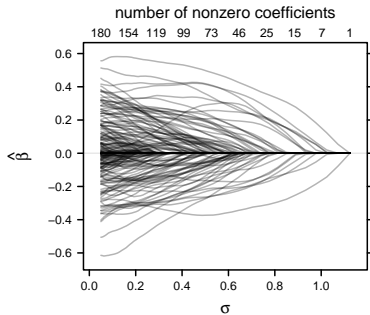
- we are interested in a **path** of penalties  $\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(m)}$ , many of which will lead to **sparse** solutions
- **basic idea**: what if we could, based on a relatively **cheap** test, determine which predictors will be inactive before fitting the model?



**Figure 1:** SLOPE path with  $n = 200$ ,  $p = 20000$ .  $\sigma$  indicates strength of regularization.

# Motivation for screening rules

- we are interested in a **path** of penalties  $\lambda^{(1)}, \lambda^{(2)}, \dots, \lambda^{(m)}$ , many of which will lead to **sparse** solutions
- **basic idea**: what if we could, based on a relatively **cheap** test, determine which predictors will be inactive before fitting the model?
- it turns out that we can, using screening rules!



**Figure 1:** SLOPE path with  $n = 200$ ,  $p = 20000$ .  $\sigma$  indicates strength of regularization.

## Strong screening rule for SLOPE

Assume that we have a solution for  $\lambda^{(k-1)}$  and want the solution at  $\lambda^{(k)}$ .

Using the optimality condition for the SLOPE problem,

$$\mathbf{0} \in \nabla g(\beta(\lambda^{(k)})) + \partial J(\beta(\lambda^{(k)}); \lambda^{(k)}),$$

where  $\partial J$  is the subgradient, we can determine which predictors will be active at  $\lambda^{(k)}$ .

## Strong screening rule for SLOPE

Assume that we have a solution for  $\lambda^{(k-1)}$  and want the solution at  $\lambda^{(k)}$ .

Using the optimality condition for the SLOPE problem,

$$\mathbf{0} \in \nabla g(\beta(\lambda^{(k)})) + \partial J(\beta(\lambda^{(k)}); \lambda^{(k)}),$$

where  $\partial J$  is the subgradient, we can determine which predictors will be active at  $\lambda^{(k)}$ .

### **our contribution**

approximate  $\nabla g(\beta(\lambda^{(k)}))$  and use optimality criterion (as if  $\nabla g(\beta(\lambda^{(k)}))$  was known) to discard predictors

## Strong screening rule for SLOPE

Assume that we have a solution for  $\lambda^{(k-1)}$  and want the solution at  $\lambda^{(k)}$ .

Using the optimality condition for the SLOPE problem,

$$\mathbf{0} \in \nabla g(\beta(\lambda^{(k)})) + \partial J(\beta(\lambda^{(k)}); \lambda^{(k)}),$$

where  $\partial J$  is the subgradient, we can determine which predictors will be active at  $\lambda^{(k)}$ .

### our contribution

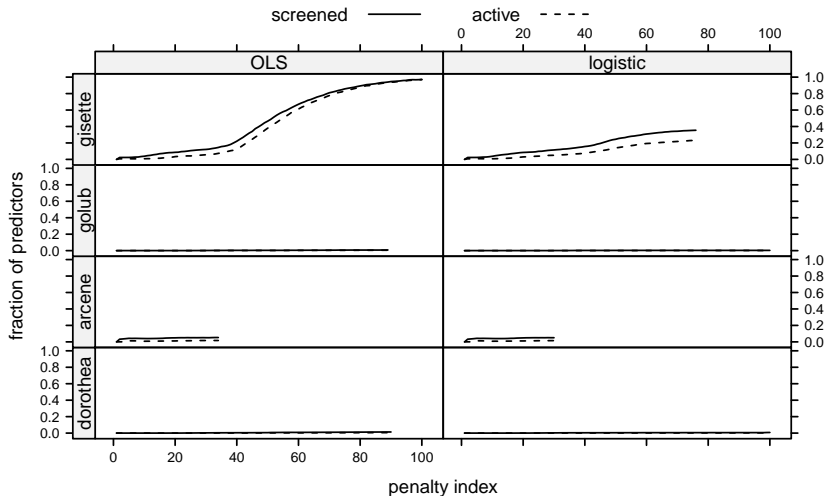
approximate  $\nabla g(\beta(\lambda^{(k)}))$  and use optimality criterion (as if  $\nabla g(\beta(\lambda^{(k)}))$  was known) to discard predictors

### violations

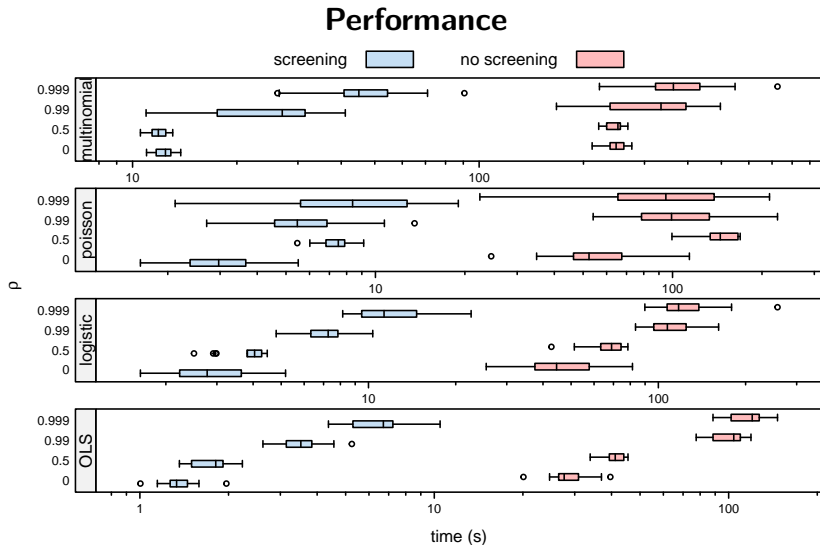
- violations—incorrectly discarding predictors—may occur
- can always be caught by checking optimality conditions (and refitting if present)
- are **so rare** that, in practice, the benefits from using the rule far outweigh the costs it incurs



## Efficiency for real data



**Figure 2:** Efficiency for real data sets. The dimensions of the predictor matrices are  $100 \times 9920$  (arcene),  $800 \times 88119$  (dorothea),  $6000 \times 4955$  (gisette), and  $38 \times 7129$  (golub).



**Figure 3:** Performance benchmarks for various generalized linear models with  $X \in \mathbb{R}^{200 \times 20000}$ . Predictors are autocorrelated through an AR(1) process with correlation  $\rho$ .

# References I