

# The Hessian Screening Rule and Adaptive Paths for the Lasso

Wroclaw Technical University Seminar

Johan Larsson

Department of Statistics, Lund University

November 16, 2021



# Overview

Preliminaries

The Hessian Screening Rule

Adaptive Lasso Path (Work in Progress)

# Preliminaries

# The Lasso

A type of penalized regression, represented by the following convex optimization problem:

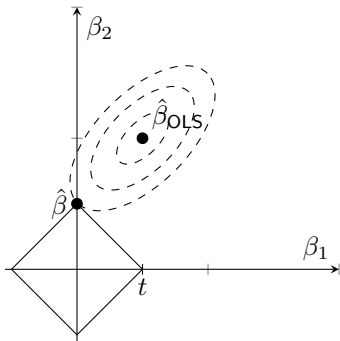
$$\underset{\beta \in \mathbb{R}^p}{\text{minimize}} \{f(\beta) + \lambda \|\beta\|_1\}$$

where  $f(\beta)$  is smooth and convex.

$f(\beta) = \frac{1}{2} \|y - X\beta\|_2^2$  leads to the ordinary lasso

$\lambda$  is a hyper-parameter that controls the level of **penalization**.

$\hat{\beta}(\lambda)$  is the solution to this problem for a given  $\lambda$ .



## The Lasso Path

Solving the lasso for  $\lambda \in [0, \lambda_{\max})$ , with

$$\lambda_{\max} := \max \{ \lambda \in \mathbb{R}^+ \mid \hat{\beta}(\lambda) = 0 \},$$

traces the set of all solutions for the lasso: the (exact) lasso path.

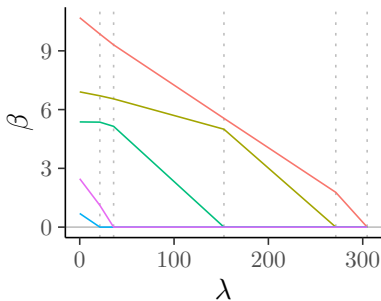
# The Lasso Path

Solving the lasso for  $\lambda \in [0, \lambda_{\max})$ , with

$$\lambda_{\max} := \max \{ \lambda \in \mathbb{R}^+ \mid \hat{\beta}(\lambda) = 0 \},$$

traces the set of all solutions for the lasso: the (exact) lasso path.

- Piece-wise linear with breaks wherever the active set changes
- The size of the active set cannot exceed  $n$ , the number of observations.



**Figure 1:** The lasso path for an example of the standard lasso

# Picking $\lambda$

## The Problem

Typically don't know the optimal value for  $\lambda$ . To tackle this, we use cross-validation to tune for  $\lambda$ .

# Picking $\lambda$

## The Problem

Typically don't know the optimal value for  $\lambda$ . To tackle this, we use cross-validation to tune for  $\lambda$ .

## Grid Search

For  $p \gg n$ , the standard procedure is to create a grid of  $\lambda$ s and solve the lasso numerically.



# Picking $\lambda$

## The Problem

Typically don't know the optimal value for  $\lambda$ . To tackle this, we use cross-validation to tune for  $\lambda$ .

## Grid Search

For  $p \gg n$ , the standard procedure is to create a grid of  $\lambda$ s and solve the lasso numerically.

But this is computationally demanding when  $p$  is large.

## Screening Rules

# Predictor Screening Rules

## Motivation

Many solutions along the regularization path are **sparse**, especially if  $p \gg n$ .

# Predictor Screening Rules

## Motivation

Many solutions along the regularization path are **sparse**, especially if  $p \gg n$ .

## Basic Idea

Say that we are at step  $k$  on the lasso path and are about to solve for step  $k + 1$ .

Intuitively, information at  $k$  should tell us something about which predictors are going to be non-zero at step  $k + 1$ .

Idea is to use this information to **discard** a subset of the predictors and fit the model to a smaller set of predictors—the screened set.

# Types of Screening Rules

## **Safe Rules**

Certifies that discarded predictors are inactive at the optimum.

# Types of Screening Rules

## Safe Rules

Certifies that discarded predictors are inactive at the optimum.

## Heuristic (Un-Safe) Rules

No guarantees! Result in **violations**: discarding predictors that actually will be active.

Need post-optimization checks of optimality conditions.

Checking the optimality conditions can be costly.

# Optimality Conditions

$\beta$  is a solution to the lasso problem if it satisfies the stationarity criterion

$$\mathbf{0} \in \nabla f(\beta) + \lambda \partial$$

where  $\partial$  is the subdifferential of  $\|\beta\|_1$ , defined as

$$\partial_j \in \begin{cases} \{\text{sign}(\beta_j)\} & \text{if } \beta_j \neq 0, \\ [-1, 1] & \text{otherwise.} \end{cases}$$

# Optimality Conditions

$\beta$  is a solution to the lasso problem if it satisfies the stationarity criterion

$$\mathbf{0} \in \nabla f(\beta) + \lambda \partial$$

where  $\partial$  is the subdifferential of  $\|\beta\|_1$ , defined as

$$\partial_j \in \begin{cases} \{\text{sign}(\beta_j)\} & \text{if } \beta_j \neq 0, \\ [-1, 1] & \text{otherwise.} \end{cases}$$

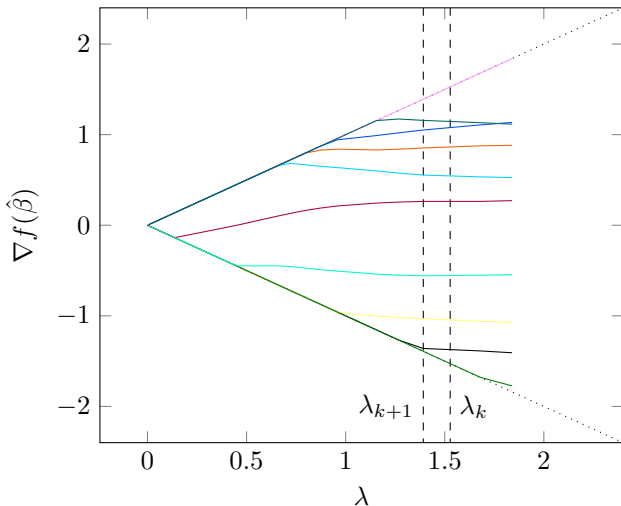
This means that

$$|\nabla f(\beta)_j| < \lambda \implies \hat{\beta}_j = 0.$$

Of course, we don't know  $\nabla f(\beta)$  prior to solving the problem.



## The Gradient Perspective of the Path



**Figure 2:** The gradient vector along the lasso path

## Screening Rules as a Gradient Estimate

Let  $c(\lambda) := -\nabla f(\beta(\lambda))$  be the so-called **correlation** vector.

$0 \in \nabla f(\beta) + \lambda \partial$  suggests a simple template for a screening rule:

1. replace  $-\nabla f(\beta)$  with an estimate  $\tilde{c}$
2. rule: if  $|\tilde{c}_j| < \lambda$ , discard predictor  $j$ .

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

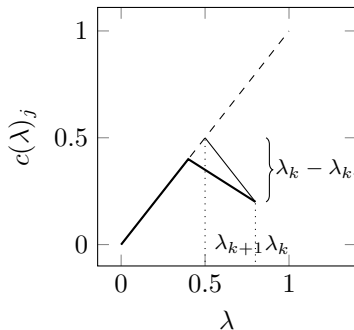
# The Strong Rule

The Strong Rule gradient estimate is

$$\tilde{c}^S(\lambda_{k+1}) = \underbrace{c(\lambda_k)}_{\text{previous gradient}} + \underbrace{(\lambda_k - \lambda_{k+1}) \text{sign}(c(\lambda_k))}_{\text{unit slope bound}}.$$

## simple idea

assume that the slope of the gradient is bounded by one (the unit slope bound) (Tibshirani et al. 2012)



# The Working Set Algorithm

The Strong Rule is conservative when predictors are correlated.

A better alternative is to use the predictors that have **ever been active** as a screened set and then

1. fit the lasso on the predictors in the screened set,
2. check the optimality conditions in the strong set, and then
3. check the optimality conditions for all predictors.

Whenever we encounter violations (in step 2 or 3), go back to step 1 and repeat.

# The Hessian Screening Rule

# The Ordinary Lasso

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

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

and

$$\nabla f(\beta) = X^T(X\beta - y).$$

# The Ordinary Lasso

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

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

and

$$\nabla f(\beta) = X^T(X\beta - y).$$

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

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

# The Ordinary Lasso

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

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

and

$$\nabla f(\beta) = X^T(X\beta - y).$$

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

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

This expression holds for an interval  $[\lambda_k, \lambda_{k+1}]$  in which no changes occur in the active set, which means we can retrieve any solution in this range via

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



## The Hessian Screening Rule

Take this expression and stick it into the gradient at step  $k + 1$ :

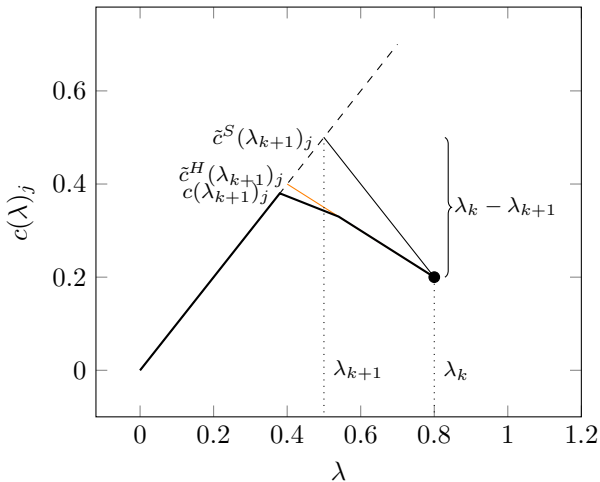
$$\begin{aligned}\tilde{c}^H(\lambda_{k+1}) &= -\nabla f(\hat{\beta}(\lambda_{k+1})_{\mathcal{A}}) \\ &= c(\lambda_k) + (\lambda_{k+1} - \lambda_k)X^T X_{\mathcal{A}}(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}}),\end{aligned}$$

which is the basic form of our screening rule: **The Hessian Screening Rule.**

Note that this is an exact expression for the correlation vector (negative gradient) at step  $k + 1$  if the activate set has remained unchanged.

The Hessian Screening Rule is a heuristic (un-safe) rule: it may result in violations.

## The Hessian and Strong Screening Rules



**Figure 3:** Conceptual comparison of screening rules

# Tweaks

## Avoiding Expensive Inner Products

The expression

$$\hat{c}^H(\lambda_{k+1}) = c(\lambda_k) + (\lambda_{k+1} - \lambda_k)X^T X_{\mathcal{A}}(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}}),$$

involves an expensive inner product with the full design matrix.

Instead, we replace  $X$  with the columns indexed by the strong rule.

---

<sup>1</sup>We've set  $\gamma = 0.01$  in our simulations, which has worked very well.

# Tweaks

## Avoiding Expensive Inner Products

The expression

$$\hat{c}^H(\lambda_{k+1}) = c(\lambda_k) + (\lambda_{k+1} - \lambda_k) X^T X_{\mathcal{A}} (X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}}),$$

involves an expensive inner product with the full design matrix.

Instead, we replace  $X$  with the columns indexed by the strong rule.

## Upwards Shift

We need some upwards bias on the estimate or else risk excessive numbers of violations. We add a fraction  $\gamma$  of the unit bound<sup>1</sup>.

---

<sup>1</sup>We've set  $\gamma = 0.01$  in our simulations, which has worked very well.

## Warm Starts

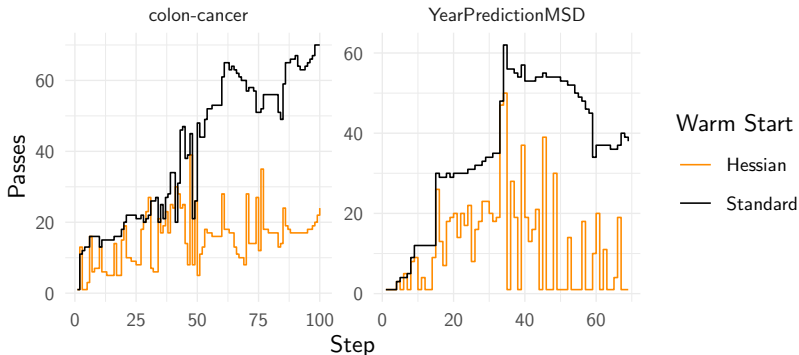
The availability of the Hessian inverse enables a better warm start:

$$\hat{\beta}(\lambda_{k+1})_{\mathcal{A}} := \hat{\beta}(\lambda_k)_{\mathcal{A}} + (\lambda_k - \lambda_{k+1})(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}})$$

## Warm Starts

The availability of the Hessian inverse enables a better warm start:

$$\hat{\beta}(\lambda_{k+1})_{\mathcal{A}} := \hat{\beta}(\lambda_k)_{\mathcal{A}} + (\lambda_k - \lambda_{k+1})(X_{\mathcal{A}}^T X_{\mathcal{A}})^{-1} \text{sign}(\hat{\beta}(\lambda_k)_{\mathcal{A}})$$



**Figure 4:** Number of passes of coordinate descent for two datasets using either Hessian warm starts or standard warm starts.

## Updating the Hessian

Computing the Hessian and its inverse naively is expensive:

$$\mathcal{O}(|\mathcal{A}|^3 + |\mathcal{A}|^2 n)$$

Fortunately, we can sweep columns of the Hessian and inverse in our out, yielding complexity

- $\mathcal{O}(|\mathcal{D}|^2 n + n|\mathcal{D}||\mathcal{E}| + |\mathcal{D}|^2|\mathcal{E}| + |\mathcal{D}|^3)$  when augmenting the Hessian and
- $\mathcal{O}(|\mathcal{C}|^3 + |\mathcal{C}|^2|\mathcal{E}| + |\mathcal{C}||\mathcal{E}|^2)$  when reducing it,

where

- $\mathcal{C} = \mathcal{A}_{k-1} \setminus \mathcal{A}_k$  (to-be deactivated)
- $\mathcal{D} = \mathcal{A}_k \setminus \mathcal{A}_{k-1}$  (to-be activated)
- $\mathcal{E} = \mathcal{A}_k \cap \mathcal{A}_{k-1}$  (still activate)

# General Loss Functions

The rule can be extended to many other loss functions in the family of generalized linear models.

Omit details here, but note that

- the gradient estimate now involves a matrix of weights—for logistic regression a diagonal matrix,
- the lasso path is no longer piece-wise linear, and
- updating the Hessian is no longer cheap.



## Results

# Setup

- Rows of the predictor matrix i.i.d. from  $\mathcal{N}(0, \Sigma)$
- Response generated from  $\mathcal{N}(X\beta, \sigma^2 I)$  with  $\sigma^2 = \beta^T \Sigma \beta / \text{SNR}$
- $s$  non-zero coefficients, equally spaced throughout the coefficient vector

## Scenario 1 (Low-Dimensional)

$n = 10\,000$ ,  $p = 100$ ,  $s = 5$ , and  $\text{SNR} = 1$

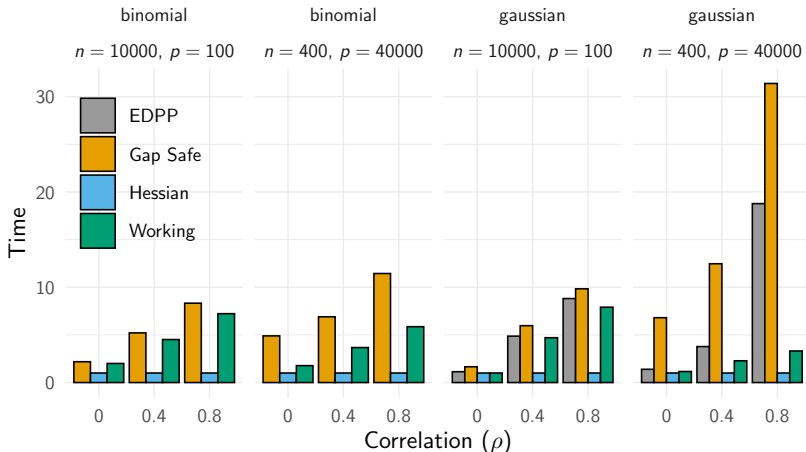
## Scenario 2 (High-Dimensional)

$n = 400$ ,  $p = 40\,000$ ,  $s = 20$ , and  $\text{SNR} = 2$

Code is located at

[github.com/jolars/HessianScreening](https://github.com/jolars/HessianScreening)

# Simulated Data



**Figure 5:** Time to fit a full regularization path for  $\ell_1$ -regularized least-squares and logistic regression to a design with  $n$  observations,  $p$  predictors, and pairwise correlation between predictors of  $\rho$ . Time is relative to the minimal value for each group.

# Real Data: Least-Squares Regression

**Table 1:** Time to fit a full regularization path of  $\ell_1$ -regularized least-squares regression to real data sets.

Dataset	$n$	$p$	Density	Time (s)			
				Gap Safe	Hessian	Working	EDPP
cadata	20 640	8	1	1.65	0.196	1.42	1.49
e2006-loglp	16 087	4 272 227	0.0014	13 200	194	204	942
e2006-tfidf	16 087	150 360	0.0083	565	29.3	66.8	337
YearPredMSD	463 715	90	1	196	46.9	159	142

## Real Data: Logistic Regression

**Table 2:** Time to fit a full regularization path of  $\ell_1$ -regularized logistic regression to real data sets.

Dataset	$n$	$p$	Density	Time (s)		
				Gap Safe	Hessian	Working
arcene	100	10 000	0.54	11	8.0	7.1
colon-cancer	62	2000	1	0.45	0.19	0.36
duke-breast-cancer	44	7129	1	0.49	0.20	0.38
ijcnn1	35 000	22	1	29	9.9	29
madelon	2000	500	1	620	94	580
news20	19 996	1 355 191	0.000 34	33 000	1700	2300
rcv1	20 242	47 236	0.0016	940	530	380

# Discussion

- simple, intuitive, idea
- performs well in our examples
- handles the highly-correlated case very well
- works for arbitrary loss functions that are twice differentiable
- works for other penalty functions too (SLOPE, lasso variations)
- a paper is submitted and under review, but a pre-print is available (Larsson and Wallin 2021).

## **Adaptive Lasso Path (Work in Progress)**

## Standard Grid Setup

The dominating choice for constructing the lasso path in the high-dimensional regime is to setup a log-spaced grid from  $\lambda_{\min}$  to  $\lambda_{\max}$  where

$$\lambda_{\min} = \lambda_{\max} \times \begin{cases} 10^{-2} & \text{if } p > n \\ 10^{-4} & \text{otherwise.} \end{cases}$$

Why this particular choice? We don't know.

### Why Not Use LARS?

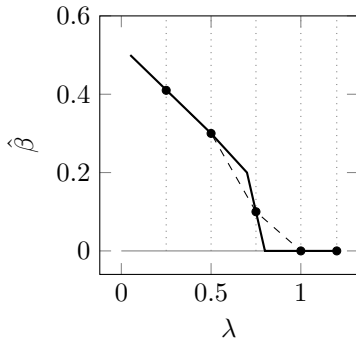
Homotopy methods scale poorly with  $p$ , in the worst case requiring  $(3^p + 1)/2$  steps



## The Problem with the Heuristic

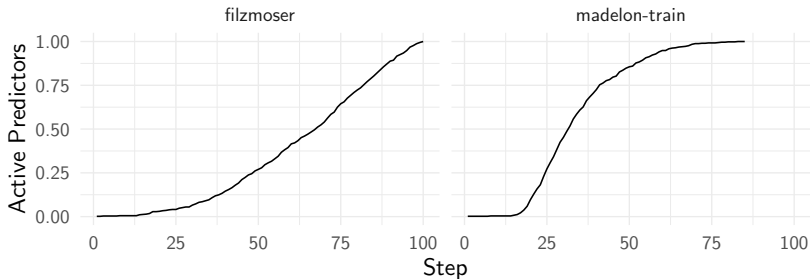
If  $\lambda$  values are spaced **coarsely**, interpolated values may fail to capture the path accurately.

On the other hand, if the values are spaced too densely where there are no critical points, then we are **wasting** resources.



## Lasso Path Profiles

When and how predictors are activated along the lasso path varies.



**Figure 6:** The number of included predictors at each step along the path for the standard grid path for the lasso. The number of included predictors varies considerably from data set to data set.

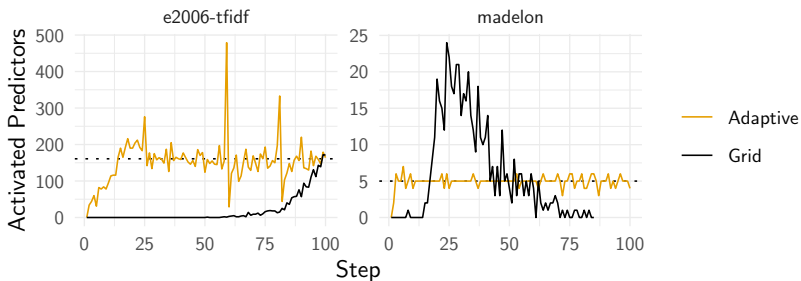
# The Adaptive Lasso Path

The gradient estimate from the Hessian screening rule can be used to predict at which  $\lambda$  values the predictors will enter the model.

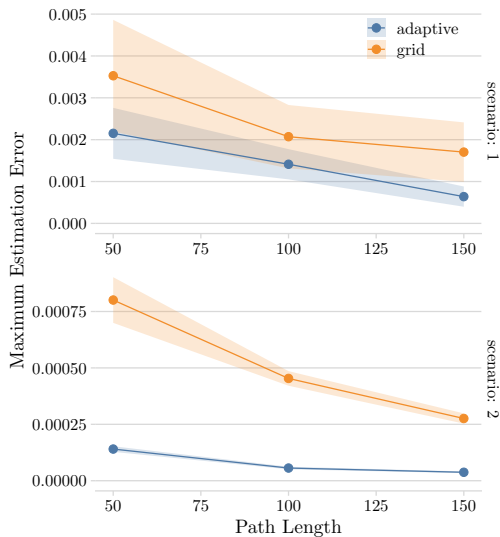
With the Adaptive Lasso Path, we use this information to find the next  $\lambda$  at which a desired number of predictors will enter the model

This lets us decide the resolution of the lasso path arbitrarily.

## An Example



**Figure 7:** The number of predictors entering the model at each step for either the adaptive path or the grid path.



**Figure 8:** Maximum error along the regularization path for the adaptive and grid methods. Scenario 1 is a high-dimensional setting and 2 a low-dimensional setup.

## Discussion

- The default grid heuristic is crude and sub-optimal in some cases.
- The adaptive lasso path adapts to the structure of the data and can be controlled by the user to tailor the resolution of the path.
- As  $n$  becomes smaller, the method converges towards a homotopy method.

# Thank You!

Thank you for listening! Questions? Thoughts?

# References I

- [1] Johan Larsson and Jonas Wallin. *The Hessian Screening Rule*. Apr. 27, 2021. arXiv: [2104.13026](https://arxiv.org/abs/2104.13026) [cs, stat]. URL: <http://arxiv.org/abs/2104.13026> (visited on 04/29/2021).
- [2] Robert Tibshirani et al. “Strong Rules for Discarding Predictors in Lasso-Type Problems”. In: *Journal of the Royal Statistical Society. Series B: Statistical Methodology* 74.2 (Mar. 2012), pp. 245–266. ISSN: 1369-7412. DOI: [10/c4bb85](https://doi.org/10.1111/j.1467-9868.2011.00855.x).