

Model Selection and Lasso



Yuxin Chen

Princeton University, Fall 2020

Outline

- Model selection
- Lasso estimator
- Risk inflation
- Minimax risk for sparse vectors

Asymptotic notation

- $f(n) \lesssim g(n)$ or $f(n) = O(g(n))$ means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \leq \text{const}$$

- $f(n) \gtrsim g(n)$ or $f(n) = \Omega(g(n))$ means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \geq \text{const}$$

- $f(n) \asymp g(n)$ or $f(n) = \Theta(g(n))$ means

$$\text{const}_1 \leq \lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} \leq \text{const}_2$$

- $f(n) = o(g(n))$ means

$$\lim_{n \rightarrow \infty} \frac{|f(n)|}{|g(n)|} = 0$$

Model selection

All models are wrong but some are useful.

— George Box

Basic linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta},$$

- $\mathbf{y} = [y_1, \dots, y_n]^\top \in \mathbb{R}^n$: observed data / response variables
- $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top$: design matrix / feature matrix (known)
 - assumed to be full rank
- $\boldsymbol{\beta} = [\beta_1, \dots, \beta_p]^\top \in \mathbb{R}^p$: **unknown signal** / regression coefficients
- $\boldsymbol{\eta} = [\eta_1, \dots, \eta_n]^\top \in \mathbb{R}^n$: noise

Throughout this lecture, we assume Gaussian noise

$$\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n)$$

Model selection / feature selection

Regression:

- find relationship between response y_i and explanatory variables $\mathbf{x}_{i,1}, \dots, \mathbf{x}_{i,p}$
- use the fitted model to **make prediction**

Question: which (sub)-set of variables / features should we include?
model

- Myth: nothing is lost by including every feature / variable available
- Paradoxically, we can often achieve better predictions by discarding a fraction of variables

Tradeoff

- Model too small \implies large bias (underfitting)
- Model too large \implies large variance and poor prediction (overfitting)

How to achieve a desired **tradeoff** between predictive accuracy and parsimony (model complexity)?

Underfitting

Recall that the least squares (LS) estimate is $\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

- Divide the design matrix into 2 parts: $\mathbf{X} = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}]$

- $\tilde{\mathbf{x}} = \begin{bmatrix} \tilde{\mathbf{x}}^{(1)} \\ \tilde{\mathbf{x}}^{(2)} \end{bmatrix}$: new data

- LS estimate based only on $\mathbf{X}^{(1)}$:

$$\hat{\beta}^{(1)} := (\mathbf{X}^{(1)\top} \mathbf{X}^{(1)})^{-1} \mathbf{X}^{(1)\top} \mathbf{y}$$

with prediction at $\tilde{\mathbf{x}}$ given by

$$\hat{y}_{\text{underfit}} = \tilde{\mathbf{x}}^{(1)\top} \hat{\beta}^{(1)}$$

- LS estimate based on true model

$$\hat{\beta} := (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

with prediction at $\tilde{\mathbf{x}}$ given by

$$\hat{y}_{\text{true}} = [\tilde{\mathbf{x}}^{(1)\top}, \tilde{\mathbf{x}}^{(2)\top}] \hat{\beta}$$

Bias due to underfitting

Suppose the ground truth is $\beta = \begin{bmatrix} \beta^{(1)} \\ \beta^{(2)} \end{bmatrix}$, then

$$\begin{aligned}\mathbb{E} \left[\hat{\beta}^{(1)} \right] &= \left(\mathbf{X}^{(1)\top} \mathbf{X}^{(1)} \right)^{-1} \mathbf{X}^{(1)\top} \left(\mathbf{X}^{(1)} \beta^{(1)} + \mathbf{X}^{(2)} \beta^{(2)} \right) \\ &= \beta^{(1)} + \underbrace{\left(\mathbf{X}^{(1)\top} \mathbf{X}^{(1)} \right)^{-1} \mathbf{X}^{(1)\top} \mathbf{X}^{(2)} \beta^{(2)}}_{\text{bias}}\end{aligned}$$

$\implies \hat{\beta}^{(1)}$ is a biased estimate of $\beta^{(1)}$

Prediction variance due to underfitting

Fact 8.1

$$\mathbf{Var} [\hat{y}_{\text{true}}] \geq \mathbf{Var} [\hat{y}_{\text{underfit}}]$$

- **Implications:** the “apparent” prediction variance tends to decrease when we adopt small models
- (Exercise): compute the prediction variance under overfitting

Proof of Fact 8.1

Observe that

$$\begin{aligned}\text{Cov}[\hat{\beta}] &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \text{Cov}[\mathbf{y}] \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \\ &= \sigma^2 \begin{bmatrix} (\mathbf{X}^{(1)\top} \mathbf{X}^{(1)})^{-1} + \mathbf{L} \mathbf{M} \mathbf{L}^\top & -\mathbf{L} \mathbf{M} \\ -\mathbf{M} \mathbf{L}^\top & \mathbf{M} \end{bmatrix} \quad (\text{matrix inversion identity})\end{aligned}$$

where $\mathbf{L} = (\mathbf{X}^{(1)\top} \mathbf{X}^{(1)})^{-1} \mathbf{X}^{(1)\top} \mathbf{X}^{(2)}$ and

$$\mathbf{M} = \left\{ \mathbf{X}^{(2)\top} \left(\mathbf{I} - \mathbf{X}^{(1)} (\mathbf{X}^{(1)\top} \mathbf{X}^{(1)})^{-1} \mathbf{X}^{(1)\top} \right) \mathbf{X}^{(2)} \right\}^{-1} \succeq \mathbf{0}.$$

This gives

$$\begin{aligned}\text{Var}[\hat{y}_{\text{true}}] &= \begin{bmatrix} \tilde{\mathbf{x}}^{(1)\top} & \tilde{\mathbf{x}}^{(2)\top} \end{bmatrix} \text{Cov}[\hat{\beta}] \begin{bmatrix} \tilde{\mathbf{x}}^{(1)} \\ \tilde{\mathbf{x}}^{(2)} \end{bmatrix} \\ &= \sigma^2 \tilde{\mathbf{x}}^{(1)\top} \left(\mathbf{X}^{(1)\top} \mathbf{X}^{(1)} \right)^{-1} \tilde{\mathbf{x}}^{(1)} + \sigma^2 \left(\mathbf{L}^\top \tilde{\mathbf{x}}^{(1)} - \tilde{\mathbf{x}}^{(2)} \right)^\top \mathbf{M} \left(\mathbf{L}^\top \tilde{\mathbf{x}}^{(1)} - \tilde{\mathbf{x}}^{(2)} \right) \\ &\geq \sigma^2 \tilde{\mathbf{x}}^{(1)\top} \left(\mathbf{X}^{(1)\top} \mathbf{X}^{(1)} \right)^{-1} \tilde{\mathbf{x}}^{(1)} = \text{Var}[\hat{y}_{\text{underfit}}]\end{aligned}$$

Model selection criteria

Choosing a subset of explanatory variables might improve prediction

Question: which subset shall we select?

One strategy

- (1) pick a criterion that measures how well a model performs
- (2) evaluate the criterion for each subset and pick the best

One popular choice: choose a model that predicts well

Prediction error and model error

- training set: \mathbf{y}, \mathbf{X}
- $\hat{\beta}$: an estimate based on training set
- **new** data: $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}\beta + \tilde{\boldsymbol{\eta}} \in \mathbb{R}^m$, where $\tilde{\boldsymbol{\eta}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)$
- Goal: use $\hat{\beta}$ to predict $\tilde{\mathbf{y}}$

One may assess the quality of the estimate based on its *prediction error* on $\tilde{\mathbf{y}}$, i.e.

$$\begin{aligned} \text{PE} &:= \mathbb{E} \left[\|\tilde{\mathbf{X}}\hat{\beta} - \tilde{\mathbf{y}}\|^2 \right] \\ &= \mathbb{E} \left[\|\tilde{\mathbf{X}}(\hat{\beta} - \beta)\|^2 \right] + 2\mathbb{E} \left[(\tilde{\mathbf{X}}(\hat{\beta} - \beta))^\top (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta) \right] + \mathbb{E} \left[\|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta\|^2 \right] \\ &= \underbrace{\mathbb{E} \left[\|\tilde{\mathbf{X}}(\hat{\beta} - \beta)\|^2 \right]}_{:=\text{ME (model error)}} + \underbrace{m\sigma^2}_{\text{variability of data}} \end{aligned}$$

Residual sum of squares (RSS)

We shall set $\tilde{\mathbf{X}} = \mathbf{X}$ (and hence $m = n$) out of simplicity

- the case where the structures of new and old data are the same

Unfortunately, we do not have access to PE (as we don't know β)

\implies need an operational criterion for estimating PE

- One candidate: estimate PE via residual sum of squares

$$\text{RSS} := \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2$$

\implies training error

Training error underestimates prediction error

Suppose $\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{\Pi}\mathbf{y}$ for some given $\mathbf{\Pi}$ with $\text{Tr}(\mathbf{\Pi}) > 0$ (e.g. LS), then

$$\text{PE} = \mathbb{E}[\text{RSS}] + 2\sigma^2\text{Tr}(\mathbf{\Pi}) > \mathbb{E}[\text{RSS}] \quad (8.1)$$

Proof:

$$\begin{aligned} \text{PE} - \mathbb{E}[\text{RSS}] &= \mathbb{E} \left[\|\tilde{\mathbf{y}} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 \right] - \mathbb{E} \left[\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2 \right] \\ &= \mathbb{E} \left[\|\tilde{\mathbf{y}}\|^2 + \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2 - 2\langle \tilde{\mathbf{y}}, \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] \\ &\quad - \mathbb{E} \left[\|\mathbf{y}\|^2 + \|\mathbf{X}\hat{\boldsymbol{\beta}}\|^2 - 2\langle \mathbf{y}, \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] \\ &= 2\mathbb{E} \left[\langle \mathbf{y} - \tilde{\mathbf{y}}, \mathbf{X}\hat{\boldsymbol{\beta}} \rangle \right] = 2\mathbb{E} \left[\langle \boldsymbol{\eta} - \tilde{\boldsymbol{\eta}}, \mathbf{\Pi}\mathbf{y} \rangle \right] \\ &= 2\mathbb{E} \left[\langle \boldsymbol{\eta}, \mathbf{\Pi}\boldsymbol{\eta} \rangle \right] \stackrel{(a)}{=} 2\text{Tr} \left(\mathbf{\Pi}\mathbb{E} \left[\boldsymbol{\eta}\boldsymbol{\eta}^\top \right] \right) \\ &= 2\sigma^2\text{Tr}(\mathbf{\Pi}), \end{aligned}$$

where (a) follows from the identity $\text{Tr}(\mathbf{A}^\top \mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{A}^\top)$.

Example: least squares (LS) estimator

The least squares solution is

$$\hat{\boldsymbol{\beta}}^{\text{ls}} := \arg \min_{\hat{\boldsymbol{\beta}}} \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2 = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

The fitted values $\hat{\mathbf{y}}^{\text{ls}}$ is given by

$$\hat{\mathbf{y}}^{\text{ls}} = \mathbf{\Pi}^{\text{ls}} \mathbf{y} := \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

$$\implies \text{PE} = \mathbb{E}[\text{RSS}] + 2\sigma^2 \text{Tr}(\mathbf{\Pi}^{\text{ls}}) = \mathbb{E}[\text{RSS}] + 2\sigma^2 p$$

LS estimator for a given model

Suppose the model (i.e. support of β) is $S \subseteq \{1, \dots, p\}$. The least squares solution given S is

$$\hat{\beta}_S := \arg \min_{\hat{\beta}} \{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 : \hat{\beta}_i = 0 \text{ for all } i \notin S\}$$

The fitted values $\hat{\mathbf{y}}$ is then given by

$$\hat{\mathbf{y}} = \mathbf{\Pi}_S \mathbf{y} := \mathbf{X}_S (\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \mathbf{y},$$

where \mathbf{X}_S is formed by the columns of \mathbf{X} at indices in S .

Mallows' C_p statistic

In view of (8.1),

$$\text{PE}(\hat{\beta}_S) = \mathbb{E}[\text{RSS}(\hat{\beta}_S)] + 2\sigma^2 \text{Tr}(\mathbf{\Pi}_S) = \mathbb{E}[\text{RSS}(\hat{\beta}_S)] + 2|S|\sigma^2,$$

since

$$\text{Tr}(\mathbf{\Pi}_S) = \text{Tr}(\mathbf{X}_S(\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top) = \text{Tr}((\mathbf{X}_S^\top \mathbf{X}_S)^{-1} \mathbf{X}_S^\top \mathbf{X}_S) = |S|.$$

Definition 8.2 (C_p statistic, Mallows '73)

$$C_p(S) = \underbrace{\text{RSS}(\hat{\beta}_S)}_{\text{training error}} + \underbrace{2\sigma^2|S|}_{\text{model complexity}}$$

C_p is an unbiased estimate of prediction error

Model selection based on C_p statistic

1. Compute $C_p(S)$ for each model S
2. Choose $S^* = \arg \min_S C_p(S)$

This is essentially an ℓ_0 -regularized least-squares problem

$$\text{minimize}_{\hat{\beta}} \quad \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 + \underbrace{2\sigma^2\|\hat{\beta}\|_0}_{\text{penalized by model complexity}} \quad (8.2)$$

Example: orthogonal design

Suppose $\mathbf{X} = \mathbf{I}$, then (8.8) reduces to

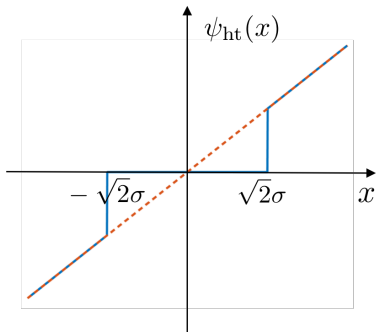
$$\text{minimize}_{\hat{\beta}} \quad \frac{1}{2} \sum_{i=1}^n (y_i - \hat{\beta}_i)^2 + 2\sigma^2 \mathbf{1}\{\hat{\beta}_i \neq 0\}$$

Solving this problem gives

$$\hat{\beta}_i = \begin{cases} 0, & |y_i| \leq \sqrt{2}\sigma \\ y_i, & |y_i| > \sqrt{2}\sigma \end{cases} \quad \text{hard thresholding}$$

- Keep large coefficients; discard small coefficients

Example: orthogonal design



$$\hat{\beta}_i = \psi_{\text{ht}}(y_i; \sqrt{2}\sigma) := \begin{cases} 0, & |y_i| \leq \sqrt{2}\sigma \\ y_i, & |y_i| > \sqrt{2}\sigma \end{cases} \quad \text{hard thresholding}$$

Hard thresholding preserves data outside threshold zone

Lasso estimator

Convex relaxation: Lasso (Tibshirani '96)

Lasso (Least absolute shrinkage and selection operator)

$$\text{minimize}_{\hat{\beta}} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 + \lambda \|\hat{\beta}\|_1 \quad (8.3)$$

for some regularization parameter $\lambda > 0$

- It is equivalent to

$$\begin{aligned} \text{minimize}_{\hat{\beta}} \quad & \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 \\ \text{s.t.} \quad & \|\hat{\beta}\| \leq t \end{aligned}$$

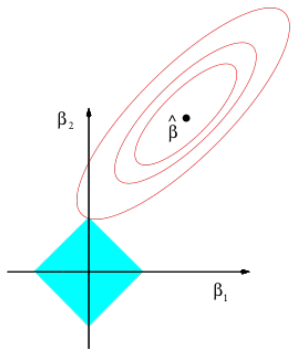
for some t that depends on λ

- a quadratic program (QP) with convex constraints

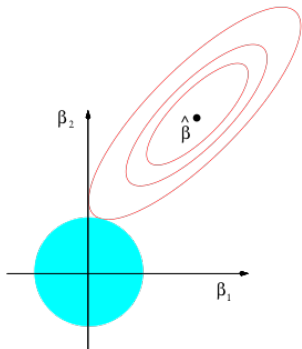
- λ controls model complexity: larger λ restricts the parameters more; smaller λ frees up more parameters

Lasso vs. MMSE (or ridge regression)

$\hat{\beta}$: least squares solution



$$\begin{aligned} &\text{minimize}_{\beta} && \|y - X\beta\|_2 \\ &\text{s.t.} && \|\beta\|_1 \leq t \end{aligned}$$

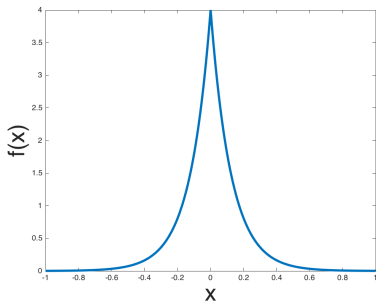


$$\begin{aligned} &\text{minimize}_{\beta} && \|y - X\beta\|_2 \\ &\text{s.t.} && \|\beta\|_2 \leq t \end{aligned}$$

Fig. credit: Hastie, Tibshirani, & Wainwright

A Bayesian interpretation

Orthogonal design: $\mathbf{y} = \boldsymbol{\beta} + \boldsymbol{\eta}$ with $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$.



Impose an i.i.d. prior on β_i to encourage sparsity (*Gaussian is not a good choice*):

(Laplacian prior) $\mathbb{P}(\beta_i = z) = \frac{\lambda}{2} e^{-\lambda|z|}$

A Bayesian interpretation of Lasso

Posterior of β :

$$\begin{aligned}\mathbb{P}(\beta | \mathbf{y}) &\propto \mathbb{P}(\mathbf{y}|\beta)\mathbb{P}(\beta) \propto \prod_{i=1}^n e^{-\frac{(y_i - \beta_i)^2}{2\sigma^2}} \frac{\lambda}{2} e^{-\lambda|\beta_i|} \\ &\propto \prod_{i=1}^n \exp\left\{-\frac{(y_i - \beta_i)^2}{2\sigma^2} - \lambda|\beta_i|\right\}\end{aligned}$$

\implies maximum *a posteriori* (MAP) estimator:

$$\arg \min_{\beta} \sum_{i=1}^n \left\{ \frac{(y_i - \beta_i)^2}{2\sigma^2} + \lambda|\beta_i| \right\} \quad (\text{Lasso})$$

Implication: Lasso is MAP estimator under Laplacian prior

Example: orthogonal design

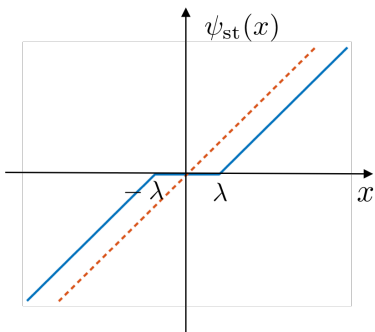
Suppose $\mathbf{X} = \mathbf{I}$, then Lasso reduces to

$$\text{minimize}_{\hat{\beta}} \quad \frac{1}{2} \sum_{i=1}^n (y_i - \hat{\beta}_i)^2 + \lambda |\hat{\beta}_i|$$

The Lasso estimate $\hat{\beta}$ is then given by

$$\hat{\beta}_i = \begin{cases} y_i - \lambda, & y_i \geq \lambda \\ y_i + \lambda, & y_i \leq -\lambda \\ 0, & \text{else} \end{cases} \quad \text{soft thresholding}$$

Example: orthogonal design

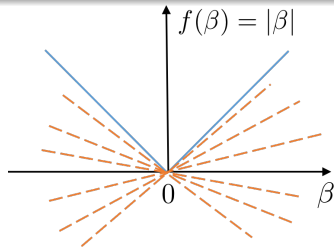


$$\hat{\beta}_i = \psi_{\text{st}}(y_i; \lambda) = \begin{cases} y_i - \lambda, & y_i \geq \lambda \\ y_i + \lambda, & y_i \leq -\lambda \\ 0, & \text{else} \end{cases} \quad \text{soft thresholding}$$

Soft thresholding shrinks data towards 0 outside threshold zone

Optimality condition for convex functions

For any convex function $f(\beta)$, β^* is an optimal solution iff $0 \in \partial f(\beta^*)$, where $\partial f(\beta)$ is the set of all subgradients at β



- s is a subgradient of $f(\beta) = |\beta|$ if

$$\begin{cases} s = \text{sign}(\beta), & \text{if } \beta \neq 0 \\ s \in [-1, 1], & \text{if } \beta = 0 \end{cases} \quad (8.4)$$

Optimality condition for convex functions

For any convex function $f(\beta)$, β^* is an optimal solution iff $0 \in \partial f(\beta^*)$, where $\partial f(\beta)$ is the set of all subgradients at β

- The subgradient of $f(\beta) = \frac{1}{2}(y - \beta)^2 + \lambda|\beta|$ can be written as

$$g = \beta - y + \lambda s \quad \text{with } s \text{ defined in (8.4)}$$

- We see that $\hat{\beta} = \psi_{\text{st}}(y; \lambda)$ by checking optimality conditions for two cases:
 - If $|y| \leq \lambda$, taking $\beta = 0$ and $s = y/\lambda$ gives $g = 0$
 - If $|y| > \lambda$, taking $\beta = y - \text{sign}(y)\lambda$ gives $g = 0$

Single-parameter setup

Consider the case where there is only a single parameter $\hat{\beta} \in \mathbb{R}$:

$$\text{minimize}_{\hat{\beta} \in \mathbb{R}} \quad \frac{1}{2} \|\mathbf{y} - \hat{\beta} \mathbf{z}\|_2^2 + \lambda |\hat{\beta}|.$$

Then one can verify that (homework)

$$\hat{\beta} = \psi_{\text{st}} \left(\frac{\mathbf{z}^\top \mathbf{y}}{\|\mathbf{z}\|_2^2}; \frac{\lambda}{\|\mathbf{z}\|_2^2} \right) = \begin{cases} \frac{\mathbf{z}^\top \mathbf{y}}{\|\mathbf{z}\|_2^2} - \frac{\lambda}{\|\mathbf{z}\|_2^2}, & \text{if } \mathbf{z}^\top \mathbf{y} > \lambda \\ 0, & \text{if } |\mathbf{z}^\top \mathbf{y}| \leq \lambda \\ \frac{\mathbf{z}^\top \mathbf{y}}{\|\mathbf{z}\|_2^2} + \frac{\lambda}{\|\mathbf{z}\|_2^2}, & \text{else} \end{cases}$$

Algorithm: coordinate descent

Idea: repeatedly cycle through the variables and, in each step, optimize only a single variable

- When updating $\hat{\beta}_j$, we solve

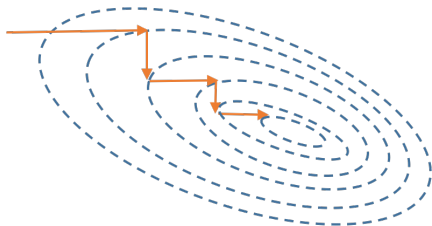
$$\text{minimize}_{\hat{\beta}_j \in \mathbb{R}} \frac{1}{2} \left\| \mathbf{y} - \sum_{i:i \neq j} \mathbf{X}_{:i} \hat{\beta}_i - \mathbf{X}_{:,j} \hat{\beta}_j \right\|_2^2 + \lambda |\hat{\beta}_j| + \lambda \sum_{i:i \neq j} |\hat{\beta}_i|$$

where $\mathbf{X}_{:,j}$ is j th column of \mathbf{X}

- This is exactly the single-parameter setting, and hence

$$\hat{\beta}_j \leftarrow \psi_{\text{st}} \left(\frac{\mathbf{X}_{:,j}^\top (\mathbf{y} - \sum_{i:i \neq j} \mathbf{X}_{:i} \hat{\beta}_i)}{\|\mathbf{X}_{:,j}\|^2}; \frac{\lambda}{\|\mathbf{X}_{:,j}\|^2} \right)$$

Algorithm: coordinate descent



Algorithm 8.1 Coordinate descent for Lasso

Repeat until convergence

for $j = 1, \dots, n$:

$$\hat{\beta}_j \leftarrow \psi_{\text{st}} \left(\frac{\mathbf{X}_{:,j}^\top (\mathbf{y} - \sum_{i:i \neq j} \mathbf{X}_{:,i} \hat{\beta}_i)}{\|\mathbf{X}_{:,j}\|_2^2}; \frac{\lambda}{\|\mathbf{X}_{:,j}\|_2^2} \right) \quad (8.5)$$

Risk inflation

Ideal risk: orthogonal design

$$y_i = \beta_i + \eta_i, \quad i = 1, \dots, n$$

Let's first select / fix a model and then estimate: for a fixed model $S \subseteq \{1, \dots, n\}$, the LS estimate $\hat{\beta}_S$ is

$$(\hat{\beta}_S)_i = \begin{cases} y_i, & \text{if } i \in S \\ 0, & \text{else} \end{cases}$$

- **Mean square estimation error for a fixed model S :**

$$\begin{aligned} \text{MSE}(\hat{\beta}_S, \beta) &:= \mathbb{E}[\|\hat{\beta}_S - \beta\|^2] = \sum_{i \in S} \mathbb{E}[(y_i - \beta_i)^2] + \sum_{i \notin S} \beta_i^2 \\ &= \underbrace{|S| \sigma^2}_{\text{variance due to noise}} + \underbrace{\sum_{i \notin S} \beta_i^2}_{\text{bias (since we don't estimate all coefficients)}} \end{aligned}$$

Ideal risk: orthogonal design

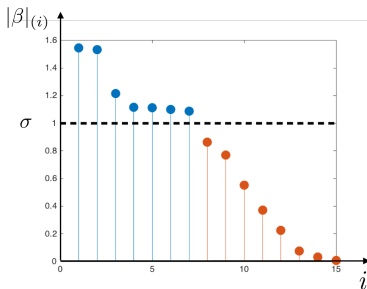
- **Smallest MSE for a fixed model size:** If we fix the model size $|S| = k$, then the best model achieves

$$\begin{aligned} \text{MSE}_k(\boldsymbol{\beta}) &:= \min_{S:|S|=k} \text{MSE}(\hat{\boldsymbol{\beta}}_S, \boldsymbol{\beta}) = k\sigma^2 + \min_{S:|S|=k} \sum_{i \notin S} \beta_i^2 \\ &= k\sigma^2 + \sum_{i=k+1}^n |\beta|_{(i)}^2 \end{aligned}$$

where $|\beta|_{(1)} \geq |\beta|_{(2)} \geq \dots \geq |\beta|_{(n)}$ are order statistics of $\{|\beta_i|\}$

Implication: good estimation is possible when $\boldsymbol{\beta}$ compresses well

Optimizing the risk over **all possible models**

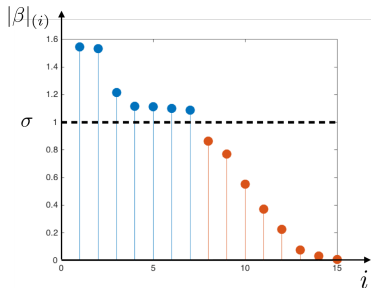


- **Ideal risk (smallest MSE over all models):** minimizing over all possible model size k gives

$$\begin{aligned} \text{MSE}^{\text{ideal}}(\beta) &:= \min_k \min_{S:|S|=k} \text{MSE}(\hat{\beta}_S, \beta) = \min_k \left\{ k\sigma^2 + \sum_{i=k+1}^n |\beta|_{(i)}^2 \right\} \\ &= \sum_{i=1}^n \min\{\sigma^2, \beta_i^2\} \end{aligned}$$

Oracle lower bound

$$\text{MSE}^{\text{ideal}}(\beta) = \sum_{i=1}^n \min\{\sigma^2, \beta_i^2\}$$



- β_i is worth estimating iff $|\beta_i| > \sigma^2$
- $\text{MSE}^{\text{ideal}}$ is the optimal risk *if an oracle reveals which variables are worth estimating and which can be safely ignored*
- With the oracle information, one can achieve $\text{MSE}^{\text{ideal}}$ via

$$\hat{\beta}_i^{\text{ideal}} = \begin{cases} y_i, & \text{if } |\beta_i| > \sigma, \\ 0 & \text{else} \end{cases} \quad (\text{eliminate irrelevant variables})$$

Risk inflation

- **Problem:** unfortunately, we do NOT know which model S is the best and hence cannot attain $\text{MSE}^{\text{ideal}}$...
- Instead, we shall treat it as an oracle lower bound, and consider the increase in estimation error **due to selecting** rather than knowing the correct model

Risk inflation

Definition 8.3 (risk inflation, Foster & George '94)

The risk inflation of an estimator $\hat{\beta}$ is

$$\text{RI}(\hat{\beta}) = \sup_{\beta} \frac{\text{MSE}(\hat{\beta}, \beta)}{\text{MSE}^{\text{ideal}}(\beta)},$$

where $\text{MSE}(\hat{\beta}, \beta) := \mathbb{E}[\|\beta - \hat{\beta}\|_2^2]$.

- Idea: calibrate the actual risk against the ideal risk for each β to better reflect the potential gains / loss
- Suggestion: *find a procedure that achieves low risk inflation!*

Risk inflation by soft / hard thresholding

Consider identity design $\mathbf{X} = \mathbf{I}$, and $\hat{\beta}_i = \psi_{\text{ht}}(y_i; \lambda)$ or $\hat{\beta}_i = \psi_{\text{st}}(y_i; \lambda)$ with threshold zone $[-\lambda, \lambda]$

- For the extreme case where $\beta = \mathbf{0}$,

$$\text{MSE}^{\text{ideal}}(\beta) = \sum_{i=1}^p \min\{\sigma^2, \beta_i^2\} = 0$$

- In order to control risk inflation, λ needs to be sufficiently large so as to ensure $\hat{\beta}_i \approx 0$ for all i . In particular,

$$\begin{aligned} \max_{1 \leq i \leq p} |y_i| &= \max_{1 \leq i \leq p} |\eta_i| \approx \sigma \sqrt{2 \log p} \quad (\text{exercise}) \\ \implies \lambda &\geq \sigma \sqrt{2 \log p} \end{aligned}$$

Risk inflation by soft / hard thresholding

Theorem 8.4 (Foster & George '94, Johnstone, Candes)

Let $\hat{\beta}$ be either a soft or hard thresholding procedure with threshold $\lambda = \sigma\sqrt{2\log p}$. Then

$$\text{MSE}(\hat{\beta}, \beta) \leq (2\log p + c) \left(\sigma^2 + \text{MSE}^{\text{ideal}}(\beta) \right)$$

where $c = 1$ for soft thresholding and $c = 1.2$ for hard thresholding.

For large p , one typically has $\text{MSE}^{\text{ideal}}(\beta) \gg \sigma^2$. Then Theorem 8.4 implies

$$\text{RI}(\hat{\beta}) \approx 2\log p$$

Proof of Theorem 8.4 for soft thresholding

WLOG, assume that $\sigma = 1$. The risk of soft thresholding for a single coordinate is

$$r_{\text{st}}(\lambda, \beta_i) := \mathbb{E}[(\psi_{\text{st}}(y_i; \lambda) - \beta_i)^2]$$

where $y_i \sim \mathcal{N}(\beta_i, 1)$.

1. There are 2 very special points that we shall single out: $\beta_i = 0$ and $\beta_i = \infty$. We start by connecting $r_{\text{st}}(\lambda, \beta_i)$ with $r_{\text{st}}(\lambda, 0)$ and $r_{\text{st}}(\lambda, \infty)$.

Lemma 8.5

$$r_{\text{st}}(\lambda, \beta) \leq r_{\text{st}}(\lambda, 0) + \beta^2 \quad (\text{quadratic upper bound})$$

$$r_{\text{st}}(\lambda, \beta) \leq r_{\text{st}}(\lambda, \infty) = 1 + \lambda^2$$

Proof of Theorem 8.4 for soft thresholding

2. The next step is to control $r_{\text{st}}(\lambda, 0)$

Lemma 8.6

$$r_{\text{st}}(\lambda, 0) \leq 2\phi(\lambda)/\lambda \stackrel{\lambda=\sqrt{2\log p}}{\ll} 1/p \quad (\text{very small})$$

where $\phi(z) := \frac{1}{\sqrt{2\pi}} \exp(-z^2/2)$.

3. With these lemmas in mind, we are ready to prove Theorem 8.4

$$\begin{aligned} \sum_{i=1}^p \mathbb{E}[(\beta_i - \hat{\beta}_i)^2] &\leq pr_{\text{st}}(\lambda, 0) + \sum_{i=1}^p \min\{\beta_i^2, \lambda^2 + 1\} \\ &< 1 + \sum_{i=1}^p \min\{\beta_i^2, 2\log p + 1\} \\ &\leq (2\log p + 1) \left[1 + \sum_{i=1}^p \min\{\beta_i^2, 1\}\right] \\ &= (2\log p + 1) \left(1 + \text{MSE}^{\text{ideal}}\right) \end{aligned}$$

Proof of Lemma 8.5

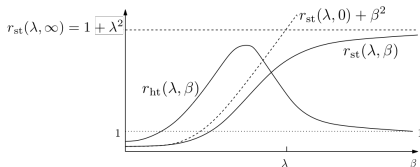


Figure adapted from
Johnstone '15

WLOG, assume $\beta \geq 0$.

(1) To prove $r_{\text{st}}(\lambda, \beta) \leq r_{\text{st}}(\lambda, 0) + \beta^2$, it suffices to show $\frac{\partial r_{\text{st}}}{\partial \beta} \leq 2\beta$, as

$$r_{\text{st}}(\lambda, \beta) - r_{\text{st}}(\lambda, 0) = \int_0^\beta \frac{\partial r_{\text{st}}(\lambda, \beta)}{\partial \beta} d\beta \leq \int_0^\beta 2\beta d\beta = \beta^2.$$

This follows since (exercise)

$$\frac{\partial r_{\text{st}}(\lambda, \beta)}{\partial \beta} = 2\beta \mathbb{P}(Z \in [-\lambda - \beta, \lambda - \beta]) \in [0, 2\beta], \quad (8.6)$$

with $Z \sim \mathcal{N}(0, 1)$

Proof of Lemma 8.5

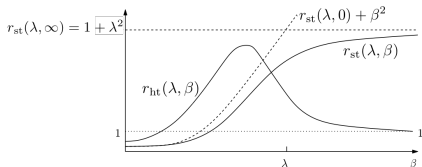


Figure adapted from
Johnstone '15

WLOG, assume $\beta \geq 0$.

(2) The identity (8.6) also shows r_{st} is increasing in $\beta > 0$, and hence

$$r_{st}(\lambda, \beta) \leq r_{st}(\lambda, \infty) = \mathbb{E}[\left((\beta + z - \lambda) - \beta\right)^2] = 1 + \lambda^2 \quad (8.7)$$

Proof of Lemma 8.6 (Candes)

$$\begin{aligned}\frac{r_{\text{st}}(\lambda, 0)}{2} &= \int_{\lambda}^{\infty} (y - \lambda)^2 \phi(y) dy \\ &= \int_{\lambda}^{\infty} (y - 2\lambda) y \phi(y) dy + \lambda^2 \int_{\lambda}^{\infty} \phi(y) dy \\ \stackrel{\text{(a)}}{=} & - \int_{\lambda}^{\infty} (y - 2\lambda) \phi'(y) dy + \lambda^2 \mathbb{P}\{Z > \lambda\} \\ \stackrel{\text{(b)}}{=} & - (y - 2\lambda) \phi(y) \Big|_{\lambda}^{\infty} + \int_{\lambda}^{\infty} \phi(y) dy + \lambda^2 \mathbb{P}\{Z > \lambda\} \\ &= -\lambda \phi(\lambda) + (1 + \lambda^2) \mathbb{P}\{Z > \lambda\} \\ \stackrel{\text{(c)}}{\leq} & -\lambda \phi(\lambda) + (1 + \lambda^2) \frac{\phi(\lambda)}{\lambda} = \frac{\phi(\lambda)}{\lambda},\end{aligned}$$

where (a) follows since $\phi'(y) = -y\phi(y)$, (b) follows from integration by parts, and (c) holds since $\mathbb{P}\{Z > \lambda\} \leq \frac{\phi(\lambda)}{\lambda}$.

Optimality

Theorem 8.7 (Foster & George '94, Jonestone)

$$\inf_{\hat{\beta}} \sup_{\beta} \frac{\text{MSE}(\hat{\beta}, \beta)}{\sigma^2 + \text{MSE}^{\text{ideal}}(\beta)} \geq (1 + o(1))2 \log p$$

- Soft and hard thresholding rules—depending only on available data without access to an oracle—can achieve the ideal risk up to the multiplicative factor $(1 + o(1))2 \log p$
- This $2 \log p$ factor is asymptotically optimal for **unrestricted** β

Comparison with canonical selection procedure

1. Minimax-optimal procedure w.r.t. risk inflation

$$\hat{\beta}_i = \psi_{\text{ht}} \left(y_i; \sigma \sqrt{2 \log p} \right)$$

2. Canonical selection based on C_p statistics

$$\hat{\beta}_i = \psi_{\text{ht}} \left(y_i; \sqrt{2} \sigma \right)$$

- Optimal procedure employs a much larger threshold zone
- **Reason:** $\min_S C_p(S)$ underestimates $\min_S \mathbb{E}[\text{PE}(S)]$ since

$$\mathbb{E} \left[\min_S C_p(S) \right] \underbrace{\leq}_{\text{sometimes} \ll} \min_S \mathbb{E} [C_p(S)] = \min_S \mathbb{E} [\text{PE}(S)]$$

- e.g. when $\beta = \mathbf{0}$, $\|\psi_{\text{ht}}(\mathbf{y}; \sqrt{2}\sigma) - \beta\|^2 \asymp n \gg 0$ with high prob.

More general models

Let's turn to a general design matrix \mathbf{X} :

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\eta} \quad \text{where } \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

One can take the ideal risk to be

$$\text{MSE}^{\text{ideal}} := \min_S \text{PE}(S) = \min_S \left\{ \underbrace{\mathbb{E}[\|\mathbf{X}_S \hat{\boldsymbol{\beta}}_S - \mathbf{X}\boldsymbol{\beta}\|_2^2]}_{\text{model error}} + |S|\sigma^2 \right\}$$

More general models

Consider the ℓ_0 -penalized selection procedure

$$\text{minimize}_{\hat{\beta}} \quad \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 + \lambda^2 \sigma^2 \|\hat{\beta}\|_0 \quad (8.8)$$

for some $\lambda \asymp \sqrt{\log p}$

Theorem 8.8 (Foster & George '94, Birge & Massart '01, Jonestone, Candes)

$$(\text{achievability}) \quad \text{MSE}(\hat{\beta}, \beta) \lesssim (\log p) \left\{ \sigma^2 + \text{MSE}^{\text{ideal}}(\beta) \right\}$$

$$(\text{minimax lower bound}) \quad \inf_{\hat{\beta}} \sup_{\beta} \frac{\text{MSE}(\hat{\beta}, \beta)}{\sigma^2 + \text{MSE}^{\text{ideal}}(\beta)} \gtrsim \log p$$

(8.8) is nearly minimax optimal for arbitrary designs!

Lasso is suboptimal for coherent design

$$\mathbf{X} = \begin{bmatrix} 1 & & & \epsilon \\ & 1 & & \vdots \\ & & \ddots & \epsilon \\ & & & 1 & 1 \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ -1/\epsilon \\ 1/\epsilon \end{bmatrix}, \quad \text{and} \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ 0 \end{bmatrix}$$

When $\epsilon \rightarrow 0$, solution to Lasso is

$$\hat{\boldsymbol{\beta}} = \begin{bmatrix} 1 \\ \vdots \\ 1 \\ 0 \end{bmatrix} \quad \text{far from the truth}$$

- Issue: the last 2 columns of \mathbf{X} are too similar / correlated

Minimax risk for sparse vectors

Asymptotic minimax risk for sparse vectors

So far we've considered risk without any restriction on β . Practically, prior knowledge (like sparsity of β) might be exploited to yield more accurate estimates.

Theorem 8.9

Suppose $X = I$. For any k -sparse β with $k \ll p$, the asymptotic minimax risk is

$$\inf_{\hat{\beta}} \sup_{\beta: \|\beta\|_0 \leq k} \text{MSE}(\hat{\beta}, \beta) = (1 + o(1))2\sigma^2 k \log(p/k)$$

Minimaxity of soft thresholding estimator

Consider $\hat{\beta}_i = \psi_{\text{st}}(y_i; \lambda)$ with $\lambda = \sigma\sqrt{2\log p}$ as before

- If $\beta = \mathbf{0}$, one has $\hat{\beta} \approx \mathbf{0}$ as discussed before
- If $\beta_1 \gg \beta_2 \gg \dots \gg \beta_k \gg \sigma$ and $\beta_{k+1} = \dots = \beta_p = 0$, then

$$\hat{\beta}_i \approx \begin{cases} y_i - \lambda, & \text{if } i \leq k \\ 0, & \text{else} \end{cases}$$

$$\begin{aligned} \implies \text{MSE}(\hat{\beta}, \beta) &\approx \sum_{i=1}^k \mathbb{E} \left[(y_i - \beta_i - \lambda)^2 \right] = k(\sigma^2 + \lambda^2) \\ &= k\sigma^2(2\log p + 1) > \underbrace{2k\sigma^2 \log(p/k)}_{\text{minimax risk}} \end{aligned}$$

- Need to pick a smaller threshold λ

Minimaxity of soft thresholding estimator

Theorem 8.10

Suppose $X = I$. For any k -sparse β with $k \ll p$, the soft thresholding estimator $\hat{\beta}_i = \psi_{\text{st}}(y_i; \lambda)$ with $\lambda = \sigma \sqrt{2 \log(p/k)}$ obeys

$$\text{MSE}(\hat{\beta}, \beta) \leq (1 + o(1))2\sigma^2 k \log(p/k)$$

- Threshold λ determined by sparsity level

Sanity check

If $\beta_1 \gg \cdots \gg \beta_k \gg \sigma$ and $\beta_{k+1} = \cdots = \beta_p = 0$, then

$$\hat{\beta}_i \approx \begin{cases} y_i - \lambda, & \text{if } i \leq k \\ 0, & \text{else} \end{cases}$$

$$\begin{aligned} \implies \text{MSE}(\hat{\beta}, \beta) &\approx k(\sigma^2 + \lambda^2) \quad (\text{as shown before}) \\ &= k\sigma^2(2 \log(p/k) + 1) \\ &\approx \underbrace{2k\sigma^2 \log(p/k)}_{\text{minimax risk}} \end{aligned}$$

Proof of Theorem 8.11

WLOG, suppose $\sigma = 1$. Under the sparsity constraint,

$$\begin{aligned} \text{MSE}(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}) &= \sum_{i=1}^p r_{\text{st}}(\lambda, \beta_i) = \sum_{i:\beta_i \neq 0} r_{\text{st}}(\lambda, \beta_i) + (p-k) r_{\text{st}}(\lambda, 0) \\ &\leq k r_{\text{st}}(\lambda, \infty) + (p-k) r_{\text{st}}(\lambda, 0) \end{aligned} \quad (8.9)$$

$$\leq k(1 + \lambda^2) + 2p \frac{\phi(\lambda)}{\lambda} \quad (8.10)$$

$$= (1 + o(1))2k \log(p/k) + \frac{k}{\sqrt{\pi \log(p/k)}}$$

$$= (1 + o(1))2k \log(p/k),$$

where (8.9) follows since $r_{\text{st}}(\lambda, \beta)$ is increasing in β , and (8.10) comes from (8.7) and Lemma 8.6

Adaptivity to unknown sparsity

- **Problem of “optimal” soft thresholding:** knowing the sparsity level *a priori* is often unrealistic
- **Question:** can we develop an estimator that is adaptive to unknown sparsity?

Adaptivity cannot be achieved via soft thresholding with **fixed thresholds**, but what if we adopt **data-dependent thresholds**?

Data-dependent thresholds

Let $|y|_{(1)} \geq \dots \geq |y|_{(p)}$ be the order statistics of $|y_1|, \dots, |y_p|$

Key idea: use a different threshold for y_i based on its rank

$$\hat{\beta}_i = \psi_{\text{st}}(y_i; \lambda_j) \quad \text{if } |y_i| = |y|_{(j)} \quad (8.11)$$

- originally due to Benjamini & Hochberg '95 for controlling false discovery rate

How to set thresholds? (non-rigorous)

Consider k -sparse vectors, and WLOG suppose $\sigma = 1$. Recall that when we use soft thresholding with $\lambda = \sqrt{2 \log(p/k)}$, the least favorable signal is

$$\beta_1 \gg \cdots \gg \beta_k \gg \sigma \quad \text{and} \quad \beta_{k+1} = \cdots = \beta_p = 0 \quad (8.12)$$

If we use data-dependent thresholds and if $\{\lambda_i\}_{i>k}$ are sufficiently large, then

$$\hat{\beta}_i \approx \begin{cases} y_i - \lambda_i, & \text{if } i \leq k \\ 0, & \text{else} \end{cases}$$

$$\text{MSE}(\hat{\beta}, \beta) \approx \sum_{i=1}^k \mathbb{E}[(y_i - \lambda_i - \beta_i)^2] = \sum_{i=1}^k (1 + \lambda_i^2)$$

How to set thresholds? (non-rigorous)

If the estimator is minimax for each k and if the worst-case β for each k is given by (8.12), then

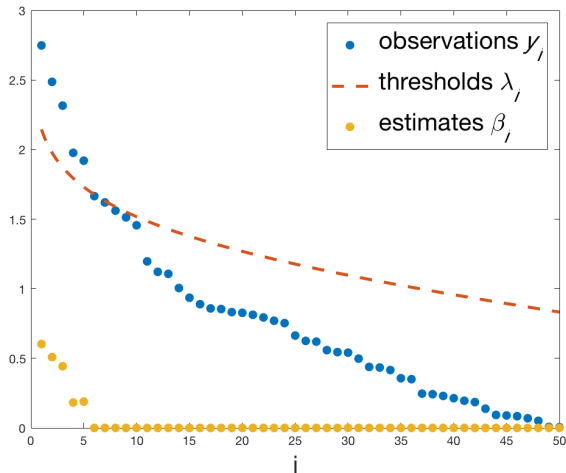
$$\text{MSE}(\hat{\beta}, \beta) \approx 2k \log(p/k), \quad k = 1, \dots, p$$

$$\implies \sum_{i=1}^k \lambda_i^2 \approx 2k \log(p/k), \quad k = 1, \dots, p$$

This suggests a choice (think of λ_i^2 as the derivative of $g(x) := 2x \log(p/x)$)

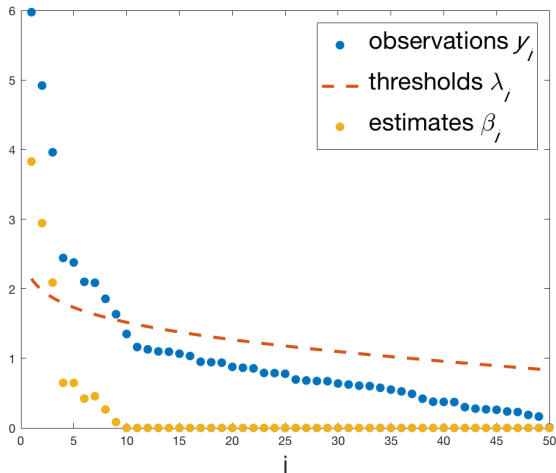
$$\lambda_i^2 \approx 2 \log(p/i) - 2 \approx 2 \log(p/i)$$

How to set thresholds? (non-rigorous)



When $\beta_1 = \dots = \beta_{50} = 0, \sigma = 1$

How to set thresholds? (non-rigorous)



When $\beta_1 = \dots = \beta_5 = 3, \beta_6 = \dots = \beta_{50} = 0, \sigma = 1$

Minimaxity

Theorem 8.11 (Abramovich '06, Su & Candes '16)

Suppose $\mathbf{X} = \mathbf{I}$, and $k \ll p$. The estimator (8.11) with $\lambda_i = \sigma \sqrt{2 \log(p/i)}$ is minimax, i.e.

$$\text{MSE}(\hat{\beta}, \beta) = (1 + o(1))2\sigma^2 k \log(p/k)$$

- Adaptive to unknown sparsity

Generalization to arbitrary design: SLOPE

SLOPE (Sorted L-One Penalized Estimation): a generalization of LASSO

$$\text{minimize}_{\hat{\beta} \in \mathbb{R}^p} \quad \frac{1}{2} \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2 + \lambda_1 |\hat{\beta}|_{(1)} + \lambda_2 |\hat{\beta}|_{(2)} + \cdots + \lambda_p |\hat{\beta}|_{(p)}$$

where $\lambda_i = \sigma \Phi^{-1}(1 - iq/(2p)) \approx \sigma \sqrt{2 \log(p/i)}$, $0 < q < 1$ is constant, and Φ is CDF of $\mathcal{N}(0, 1)$

- This is a convex program if $\lambda_1 \geq \cdots \geq \lambda_p \geq 0$ (homework)
- This can be computed efficiently via proximal methods
- SLOPE is minimax and adaptive to unknown sparsity under i.i.d. Gaussian design \mathbf{X}

Reference

- “*Lecture notes, Theory of Statistics (Stats 300C)*,” E. Candes.
- “*Statistical machine learning for high-dimensional data*,” J. Fan, R. Li, C. Zhang, H. Zou, 2018.
- “*Linear regression analysis*,” G. Seber and A. Lee, Wiley, 2003.
- “*Gaussian estimation: sequence and wavelet models*,” I. Johnstone, 2015.
- “*Some comments on C_p* ,” C. Mallows, *Technometrics*, 1973.
- “*The risk inflation criterion for multiple regression*,” D. Foster and E. George, *Annals of Statistics*, 1994.
- “*Regression Shrinkage and Selection via the lasso*,” R. Tibshirani, *Journal of the Royal Statistical Society*, 1996.
- “*Statistical learning with sparsity: the Lasso and generalizations*,” T. Hastie, R. Tibshirani, and M. Wainwright, 2015.

Reference

- “*Gaussian model selection*,” L. Birge and P. Massart, *Journal of the European Mathematical Society*, 2011.
- “*Controlling the false discovery rate: a practical and powerful approach to multiple testing*,” Y. Benjamini and Y. Hochberg, *Journal of the Royal Statistical Society*, 1995
- “*Adapting to unknown sparsity by controlling the false discovery rate*,” F. Abramovich, Y. Benjamini, D. Donoho and I. Johnstone, *Annals of Statistics*, 2006.
- “*SLOPE is adaptive to unknown sparsity and asymptotically minimax*,” W. Su and E. Candes, *Annals of Statistics*, 2016.