# Convexification

Lecture 3

# Recap from last lecture

- **Model:** $y_i = \langle \beta^*, x_i \rangle + \varepsilon_i$, $i = 1, \ldots, m$

  with $\varepsilon_i \overset{iid}{\sim} N(0, \sigma^2)$

- **Notation:**

$$Y := \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \; ; \quad f^* := \begin{bmatrix} f^*(x_1) \\ \vdots \\ f^*(x_m) \end{bmatrix} \; ; \quad \mathcal{E} = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{bmatrix}$$

and $X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} \in \mathbb{R}^{m \times p}$

$$\Rightarrow \quad Y = X\beta^* + \mathcal{E} = f^* + \mathcal{E}.$$

- **Hidden structure:** We assume that

$$|\beta^*|_0 := \text{card} \{ j : \beta_j^* \neq 0 \} \text{ is small}$$

$\underbrace{\qquad\qquad\qquad\qquad}_{\text{Coordinate sparse assumption.}}$

- For $S \subset \{1, \ldots, p\}$, we set

$$\overline{S} = \{ X\beta : \text{supp}(\beta) \subset S \}$$

and $\hat{f}^{(s)} := X\hat{\beta}^{(s)}$ with

$$\hat{\beta}^{(s)} \in \underset{\text{supp}(\beta) \subset S}{\arg\min} \| Y - X\beta \|^2$$

- **Structure learning:**

$$\hat{S} \in \underset{S \subset \{1, \ldots, p\}}{\arg\min} \{ \| Y - \hat{f}^{(s)} \|^2 + \text{pen}(s) \sigma^2 \} \quad \text{(ns)}$$

with $\text{pen}(s) = K \, |S| \, \log \dfrac{e \, p}{|S|}$

$\quad\quad\quad\quad\quad \underset{\text{constant} \, \gtrsim 2}{\uparrow}$

fulfills

$$R(\hat{f}^{(\hat{S})}) \leq \frac{\sigma^2}{m} |\beta^*|_0 \, \log \frac{e \, p}{|\beta^*|_0}$$

$\quad\quad\quad\quad \underbrace{\qquad\qquad\qquad}_{\text{minimax optimal}}$

- **Main issue:**

prohibitive computational complexity.

# Solution(s)?

→ convex proxy for the minimisation problem (ПS)

⤳ this lecture

→ greedy / iterative approximate minimisation

⤳ next lecture.

# Our goal today:

→ explain and discuss the convexification paradigm in the coordinate sparse setting

→ highlights the strengths and weaknesses of this approach.

---

To avoid normalizing issues, we assume in the following that the columns $X_{:j}$ of $X$ have been normalized $\|X_{:j}\| = 1$.

# Lasso estimator

Let us consider the approximate version of (ПS)

$$\hat{S} \in \underset{S \subset \{1, \dots, p\}}{\arg\min} \left\{ \|Y - \hat{f}^{(s)}\|^2 + \lambda |S| \right\}, \quad \text{with } \lambda = K \sigma^2 \log p$$

$\uparrow$ constant.

Since $\hat{f}^{(s)} = X \hat{\beta}^{(s)}$, with $\hat{\beta}^{(s)} \in \underset{\beta : \, \text{Supp}(\beta) \subset S}{\arg\min} \|Y - X\beta\|^2$, we have

$$\hat{S} \in \underset{S \subset \{1 \dots p\}}{\arg\min} \quad \underset{\beta : \, \text{Supp}(\beta) = S}{\min} \left\{ \|Y - X\beta\|^2 + \lambda |\beta|_0 \right\}$$

and

$$\hat{\beta}^{(\hat{S})} \in \underset{\beta \in \mathbb{R}^p}{\arg\min} \left\{ \underbrace{\|Y - X\beta\|^2}_{\text{nicely convex} \; \smile} + \underbrace{\lambda |\beta|_0}_{\text{highly non-convex} \; \frown} \right\}$$

# Recipe:

→ **constrained version**

$$\min \ \|Y - X\beta\|^2$$
$$|\beta|_0 \leq D$$

→ **convexification**

$$|\beta|_0 \leq D \rightsquigarrow |\beta|_1 \leq R$$

- $p = 2$
- $|S| = 1$



$\{|\beta|_0 \leq 1\} \cap \{\|\beta\| \leq 1\}$

$\{|\beta|_1 \leq 1\}$

# Lasso estimator:

- $\hat{\beta}^{(\lambda)} \in \underset{\beta \in \mathbb{R}^p}{\arg\min} \ \left\{ \|Y - X\beta\|^2 + \lambda |\beta|_1 \right\}$

- $\hat{f}^{(\lambda)} := X \hat{\beta}^{(\lambda)}$

## Geometric interpretation:

constrained version

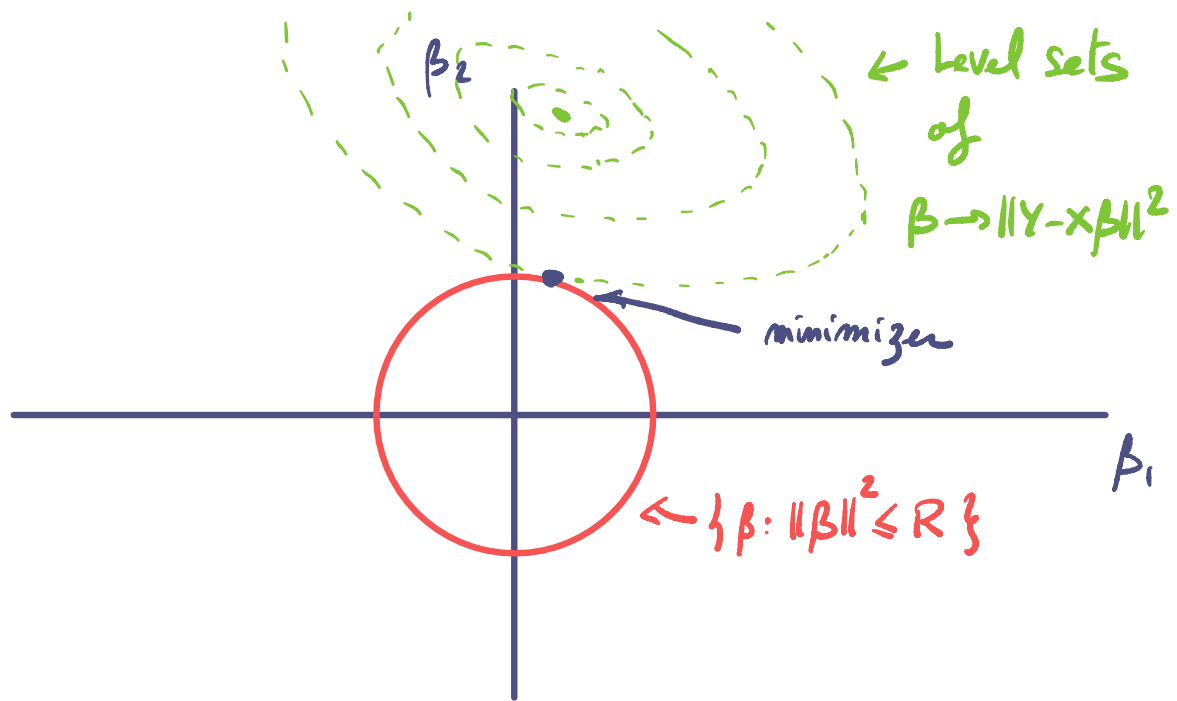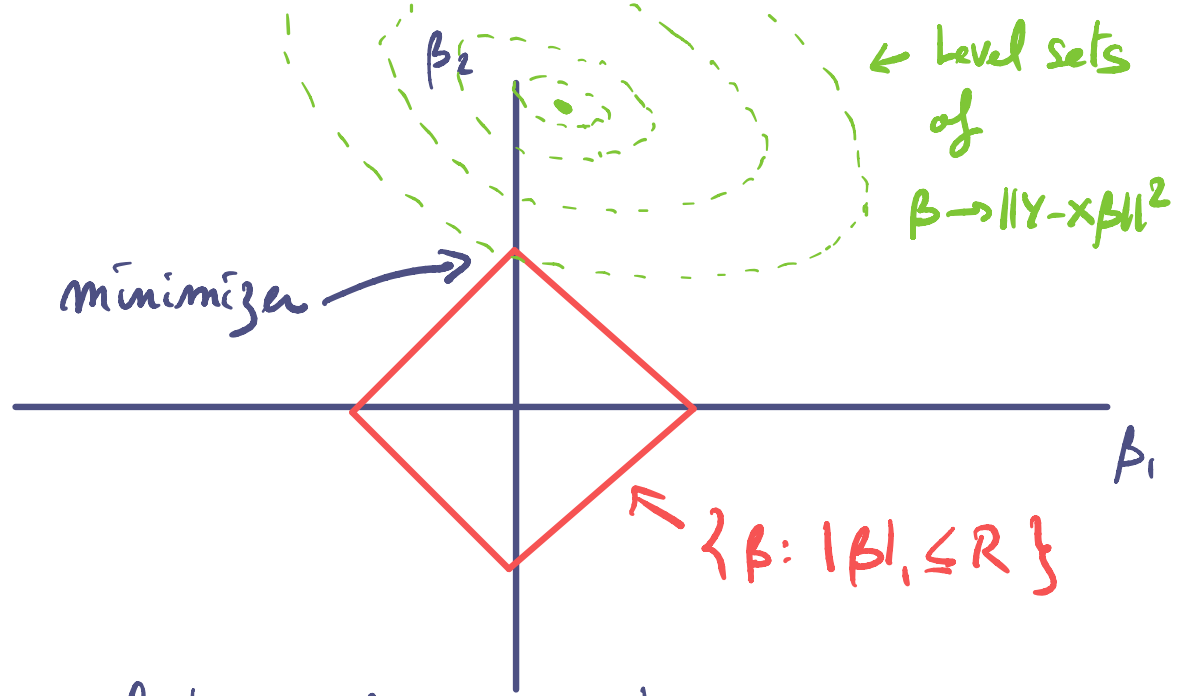$$\min_{|\beta|_1 \leq R} \quad \|Y - X\beta\|^2$$

Level sets of $\beta \to \|Y - X\beta\|^2$

minimizer →

$\{\beta : |\beta|_1 \leq R\}$

$\beta_2$

$\beta_1$

## Remark:

singularities of $\{|\beta|_1 \leq R\}$ $\longleftrightarrow$ selection of coordinates

## Ridge: $\ell^1$- ball $\leadsto$ $\ell^2$- ball

$$\min_{\|\beta\|^2 \leq R} \quad \|Y - X\beta\|^2$$

$\leadsto$ no selection occurs.

Level sets of $\beta \to \|Y - X\beta\|^2$

minimizer

$\{\beta : \|\beta\|^2 \leq R\}$

$\beta_2$

$\beta_1$

# Analytic analysis

The objective function
$$\mathcal{L}_\lambda(\beta) = \|Y - X\beta\|^2 + \lambda |\beta|_1$$
is convex but not differentiable
where $\beta_j = 0$ for some $j \in \{1, \dots, p\}$

## Subdifferential:

$$\partial |\beta|_1 = \left\{ 3 \in \mathbb{R}^p : \begin{array}{l} \cdot \, 3_j = \text{sign}(\beta_j) \text{ if } \beta_j \neq 0 \\ \cdot \, 3_j \in [-1, 1] \text{ if } \beta_j = 0 \end{array} \right\}$$

So
$$\partial \mathcal{L}_\lambda(\beta) = \left\{ -2X^T(Y - X\beta) + \lambda 3 : 3 \in \partial |\beta|_1 \right\}$$

Since $0 \in \partial \mathcal{L}_\lambda(\hat{\beta}^{(\lambda)})$, $\exists \, \hat{3} \in \partial |\hat{\beta}^{(\lambda)}|_1$
such that
$$\underbrace{X^T X \hat{\beta}^{(\lambda)} = X^T Y}_{\hookleftarrow \text{ least square}} - \underbrace{\frac{\lambda}{2} \hat{3}}_{\text{selection}}$$

Set $X_{\hat{S}} := X[\cdot, \hat{S}]$, where $\hat{S} = \text{supp}(\hat{\beta}^{(\lambda)})$

Then, since $\hat{3}_{\hat{S}} = \text{sign}(\hat{\beta}_{\hat{S}}^{(\lambda)})$
we have
$$X_{\hat{S}}^T X_{\hat{S}} \hat{\beta}_{\hat{S}}^{(\lambda)} = X_{\hat{S}}^T Y - \frac{\lambda}{2} \text{sign}(\hat{\beta}_{\hat{S}}^{(\lambda)})$$

so
$$\hat{\beta}_{\hat{S}}^{(\lambda)} = \underbrace{(X_{\hat{S}}^T X_{\hat{S}})^{-1} X_{\hat{S}}^T Y}_{= \hat{\beta}^{(\hat{S})}} - \underbrace{\frac{\lambda}{2}(X_{\hat{S}}^T X_{\hat{S}})^{-1} \text{sign}(\hat{\beta}_{\hat{S}}^{(\lambda)})}_{\text{bias term}}$$

(least square estimator on $\hat{S}$ )

bias term induced by the $\ell^1$ constraint

Remark: We cannot get an explicit expression for $\hat{\beta}^{(\lambda)}$, but in the case where $X$ has orthogonal columns:
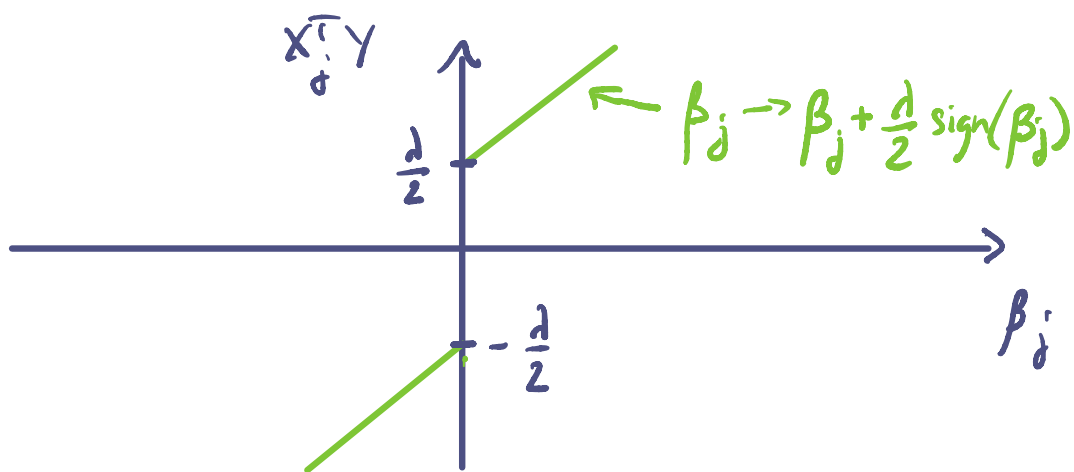$$X^T X = I_p.$$

Case $X^T X = I_p$ :

then $\hat{\beta} = X^T Y - \frac{\lambda}{2} \hat{g}$ with $\hat{g} \in \partial |\hat{\beta}|_1$

$\rightarrow$ if $\hat{\beta}_j \neq 0$: then

$$\hat{\beta}_j = X_j^T Y - \frac{\lambda}{2} \text{sign}(\hat{\beta}_j) \quad \text{ie}$$

$$X_j^T Y = \hat{\beta}_j + \frac{\lambda}{2} \text{sign}(\hat{\beta}_j)$$

$\leadsto$ only possible if $|X_j^T Y| > \lambda/2$



$\rightarrow$ if $\hat{\beta}_j = 0$ : then

$$\hat{g}_j = \frac{2}{\lambda} X_j^T Y \in [-1, 1]$$

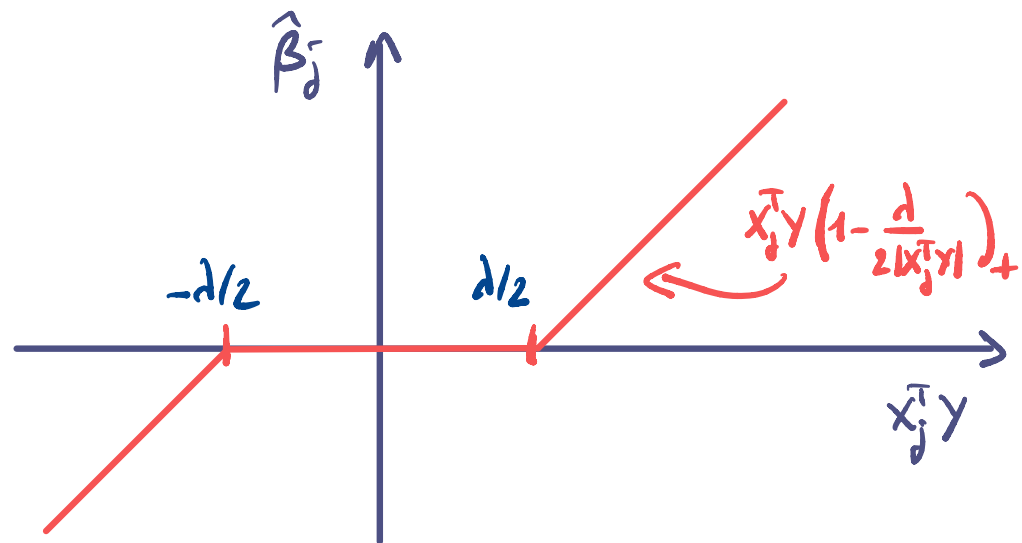$\leadsto$ only possible if $|X_j^T Y| \leq \lambda/2$

So

$$\hat{\beta}_j = \begin{cases} 0 & \text{if } |X_j^T Y| \leq \lambda/2 \\ X_j^T Y - \frac{\lambda}{2} \text{sign}(X_j^T Y) & \text{if } |X_j^T Y| > \lambda/2 \end{cases}$$

$$= \underbrace{X_j^T Y}_{\text{least square}} \underbrace{\left(1 - \frac{\lambda}{2|X_j^T Y|}\right)_+}_{\text{selects and shrinks}}$$

# Comparaison between Lasso and (ΠS): when $X^T X = I_p$

(ΠS): $\hat{\beta}^{(\Pi S)} \in \arg\min \{ \|Y - X\beta\|^2 + \tau^2 |\beta|_0 \}$

(Lasso): $\hat{\beta}^{(Lasso)} \in \arg\min \{ \|Y - X\beta\|^2 + 2\tau |\beta|_1 \}$



$\hat{\beta}^{(\Pi S)}_j = X_j^T Y \cdot \mathbb{1}_{|X_j^T Y| > \tau}$

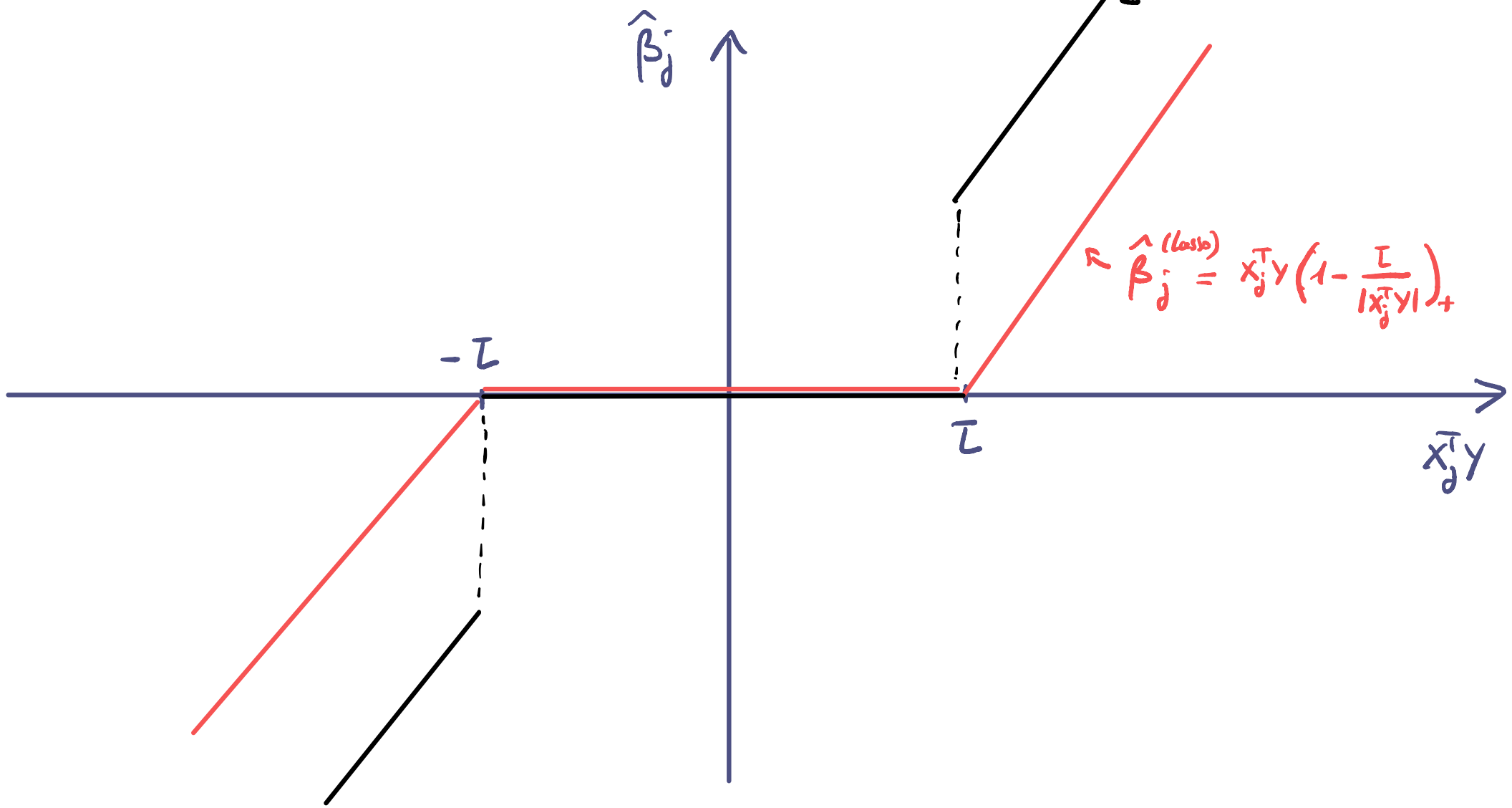$\hat{\beta}^{(Lasso)}_j = X_j^T Y \left(1 - \frac{\tau}{|X_j^T Y|}\right)_+$

# Theoretical garanties

Can we compare the performance of $\hat{\beta}^{(Lasso)}$ to $\hat{\beta}^{(ns)}$ ?

## Compatibility constant:

$$\mathcal{K}(\beta^*) = \min\left\{ \frac{\sqrt{|\beta^*|_0}\,\|X\sigma\|}{|\sigma_S|_1} : \sigma \in \mathcal{C}(\beta^*)\right\}$$

where
- $S = supp(\beta^*)$
- $\mathcal{C}(\beta^*) = \left\{\sigma \in \mathbb{R}^p : 5|\sigma_S|_1 > |\sigma_{S^c}|_1\right\}$

$\rightsquigarrow$ account for (local) orthogonality

__Fact__ : $\mathcal{K}(\beta) \geq d_{min}(X^T X)^{1/2}$

__Proof__ : for any $\sigma \in \mathbb{R}^p$:

$$\|X\sigma\|^2 \geq d_{min}\|\sigma\|^2 \geq d_{min}\|\sigma_S\|^2$$
$$\underset{c.s.}{\geq} d_{min}\frac{|\sigma_S|_1^2}{|S|^2}$$

$\square$

---

__Theorem:__ for $\lambda = 3\sigma\sqrt{2K\log p}$, $\underset{>1}{\underbrace{\phantom{2K}}}$

we have with probability $\geq 1 - \frac{1}{p^{K-1}}$

$$d_m\left(\hat{f}^{(\lambda)}, f^*\right)^2 \leq C_K\, \frac{\sigma^2|\beta^*|_0 \log p}{m\,\mathcal{K}(\beta^*)^2}$$

price to pay for computational tractability

$\left(\frac{1}{\mathcal{K}(\beta^*)^2}\right.$ can be huge$\left.\right)$

__Proof:__ see Theorem S.1.   $\square$

---

Bias of Lasso

estimators

# Example

- We have $m = 60$ noisy observations of $f^* : [0,1] \to \mathbb{R}$

$$y_i = f^*\left(\frac{i}{m}\right) + \varepsilon_i, \qquad i = 1, \dots, m$$

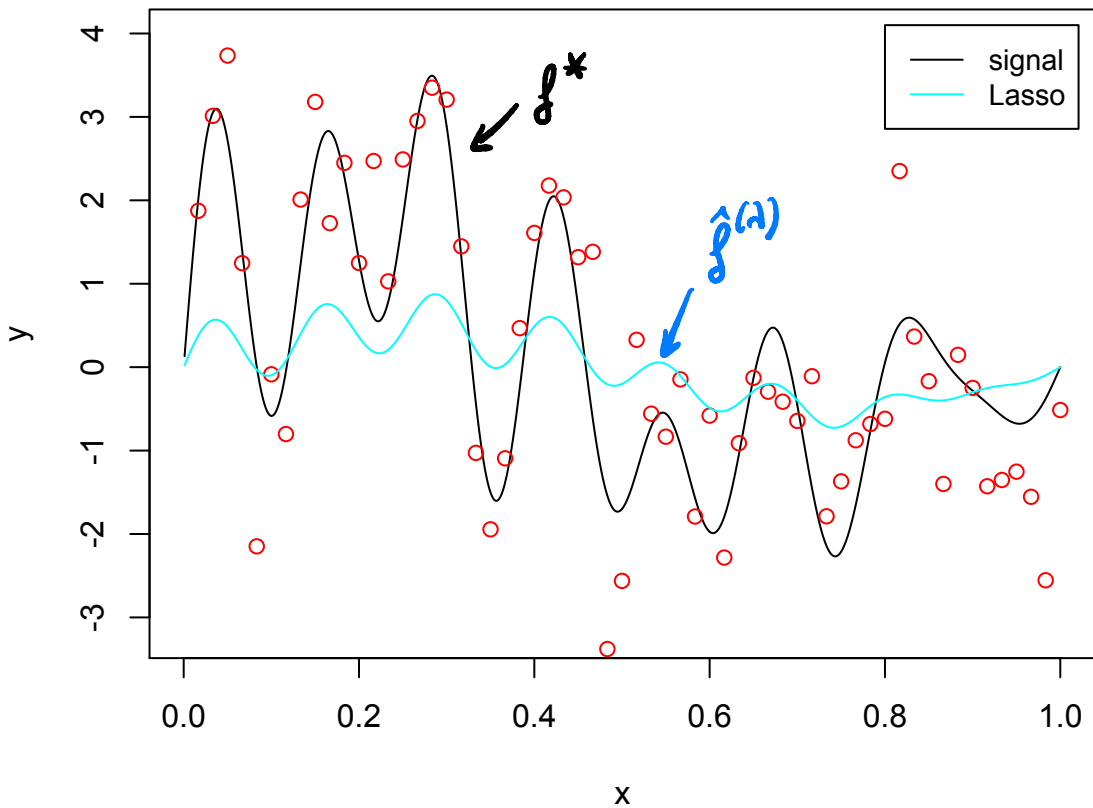- We expand $f^*$ on the Fourier basis $\{\varphi_j : j \geq 0\}$

$$f^*\left(\frac{i}{m}\right) = \sum_j \beta_j^* \underbrace{\varphi_j(i/m)}_{=: X_{ij}}$$

- We compute $\hat{\beta}^{Lasso}$ and plot the estimator

$$\hat{f}(x) = \sum_j \hat{\beta}_j^{Lasso} \varphi_j(x)$$

# Lasso

# Why?

We have:

$$\hat{\beta}^{(\lambda)} \in \underset{\beta \in \mathbb{R}^P}{\operatorname{argmin}} \quad \|Y - X\beta\|^2 + \lambda |\beta|_1$$

promotes small norm solutions

It can be seen in the formula

$$\hat{\beta}^{(\lambda)}_{\hat{S}} = \left(X_{\hat{S}}^T X_{\hat{S}}\right)^{-1} X_{\hat{S}}^T Y - \frac{\lambda}{2}\left(X_{\hat{S}}^T X_{\hat{S}}\right)^{-1} \operatorname{sign}\left(\hat{\beta}^{(\lambda)}_{\hat{S}}\right)$$

unbiased least square estimator on $\hat{S}$

bias induced by the $\ell^1$ penalty.

# Gauss-Lasso estimator

- Compute $\hat{\beta}^{(\lambda)}$ lasso estimator and set $\hat{S} = \text{supp}(\hat{\beta}^{(\lambda)})$

- Fit the least square estimator on $\hat{S}$ :
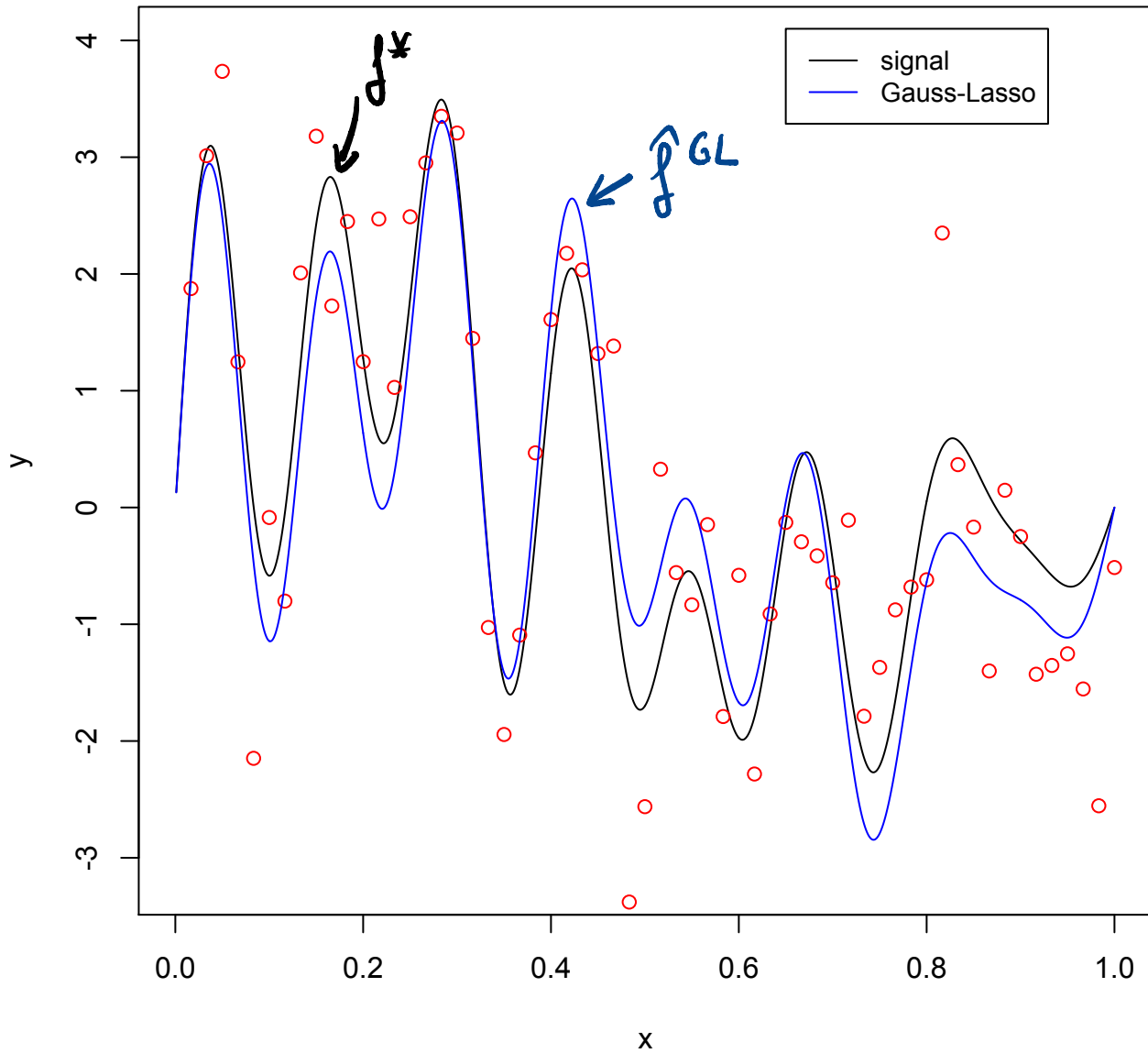
  - $\hat{\beta}^{GL}_{\hat{S}^c} \equiv 0$

  - $\hat{\beta}^{GL}_{\hat{S}} = (X^T_{\hat{S}} X_{\hat{S}})^{-1} X^T_{\hat{S}} Y$

$\rightsquigarrow$ it removes the shrinkage bias $\frac{\lambda}{2}(X^T_{\hat{S}} X_{\hat{S}})^{-1} \text{sign}(\hat{\beta}^{(\lambda)}_{\hat{S}})$

$\triangle$ For selecting $\lambda$, apply cross-validation to $\hat{\beta}^{GL}$, not $\hat{\beta}^{(\lambda)}$.

**Gauss-Lasso**

# The bias creates a deletere noise

Ref: W. Su, M. Bogdan, E. Candès
"False discoveries occur early on the Lasso path" (2016)

Setting:

- $x_{:j} \sim \mathcal{N}(0, \frac{1}{m} I_p)$

  (hence $\mathbb{E}[\|x_{:j}\|^2] = 1$)

- $\beta^*_j \overset{iid}{\sim} \alpha \delta_0 + (1-\alpha) \mathcal{V}$

  with $\mathcal{V}(0) = 0$ and $\int x^2 d\mathcal{V}(x) < +\infty$

  so $S^* := \text{supp}(\beta^*)$ fulfills

  $\mathbb{E}[|S^*|] = \alpha p$

- $m = \delta p$, with $\delta > \alpha$

- $\hat{\beta}^{(\lambda)} \in \underset{\beta}{\text{argmin}} \left\{ \|Y - X\beta\|^2 + \lambda|\beta|_1 \right\}$

  $\hat{S}^{(\lambda)} = \text{support}(\hat{\beta}^{(\lambda)})$

## Theorem (informal)

Even if $\sigma = 0$ (no noise),

$$\frac{|\hat{S}^{(\lambda)} \setminus S^*|}{|\hat{S}^{(\lambda)}|} \geq \text{something} > 0$$

with high probability

So, even when there is no noise (!), for any $\lambda > 0$, a positive fraction of the variables selected by the Lasso are **not** in $S^*$.

Why? The bias induces a pseudo-noise blurring the residuals

**Hand-waving proof:** $\sigma^2 = 0$ so that
$$Y = X\beta^* = X_{S^*}\beta_{S^*}^*$$

- We have
$$(X^T X)\hat{\beta} \stackrel{(*)}{=} X^T Y - \frac{\lambda}{2}\hat{z}$$

where $\hat{z}_j = \begin{cases} \text{sign}(\hat{\beta}_j) & \text{if } \hat{\beta}_j \neq 0 \\ \in [-1,1] & \text{otherwise} \end{cases}$

From $(*)$, we also have
$$\hat{z}_j = \frac{2}{\lambda} X_{\cdot j}^T (X\hat{\beta} - \overbrace{X_{S^*}\beta_{S^*}^*}^{Y})$$

- Let $\tilde{\beta}_{S^*} = \text{Lasso}(Y, X_{S^*})$

so that
$$X_{S^*}^T X_{S^*} \tilde{\beta}_{S^*} = X_{S^*}^T X_{S^*}\beta_{S^*}^* - \frac{\lambda}{2}\tilde{z}_{S^*}$$

i.e
$$\tilde{\beta}_{S^*} - \beta_{S^*}^* = -\frac{\lambda}{2}(X_{S^*}^T X_{S^*})^{-1}\tilde{z}_{S^*}$$

---

**hand-waving claim:** — — — — — — —

$$\text{if } \left| \frac{2}{\lambda} X_{\cdot j}^T \left( X_{S^*}\tilde{\beta}_{S^*} - X_{S^*}\beta_{S^*}^* \right) \right| > 1$$

then variable $j$ is selected in $\hat{S}^{(\lambda)}$.

- for $j \notin S^*$, we have conditionally on $X_{S^*}$
$$\frac{2}{\lambda} X_{\cdot j}^T X_{S^*} (\tilde{\beta}_{S^*} - \beta_{S^*}^*) \sim \mathcal{N}(0, v^2)$$

where $v^2 = \frac{4}{m\lambda^2} \| X_{S^*}(\tilde{\beta}_{S^*} - \beta_{S^*}^*) \|^2$

$|S^*| = \frac{k}{\delta}m \rightarrow \quad \simeq \frac{4}{m\lambda^2} \|\tilde{\beta}_{S^*} - \beta_{S^*}^*\|^2$

$\underbrace{}_{< 1}$

$$\simeq \frac{1}{m} \| (X_{S^*}^T X_{S^*})^{-1}\tilde{z}_{S^*} \|^2$$

$$\simeq \frac{1}{m} \|\tilde{z}_{S^*}\|^2$$

$$\geq \frac{1}{m} \times |S^* \cap \hat{S}^{(\lambda)}|$$

$$= \underbrace{\frac{|S^*|}{m}}_{= \alpha/\delta} \times \frac{|S^* \cap \hat{S}^{(\lambda)}|}{|S^*|}$$

Hence $v^2 \geq$ constant $> 0$, when $|\hat{S}^{(\lambda)} \cap S^*| \geq$ constant $|S^*|$.

So $\quad \mathbb{P}[\, j \text{ selected}\,] \geq$ constant $> 0$, when $|\hat{S}^{(\lambda)} \cap S^*| \geq$ constant $|S^*|$

and around constant $\times \underbrace{(p - |S^*|)}_{(1-\alpha)p}$ variables $\notin S^*$ are selected,

leading to a non-vanishing False Discovery Proportion.

This mis-selection is due to the non-vanishing bias

$$ X_{S^*}(\tilde{\beta}_{S^*} - \hat{\beta}_{S^*}) = -\frac{\lambda}{2} X_{S^*}(X_{S^*}^T X_{S^*})^{-1} \tilde{z}_{S^*} \quad \text{which is correlated} $$

with the $X_{:j}$, $j \notin S^*$.

Remark: when $\quad |S^*| \leq c \frac{n}{\log p}$, the bias is not strong enough

in order to create so many false positives.

# Adaptive-Lasso Estimator

• Take $\hat{\beta}^{init}$ be a first rough estimator of $\beta^*$, for example for example a Ridge estimator

for $\beta \approx \hat{\beta}^{init}$, we have $|\beta|_0 \approx \sum_j \frac{|\beta_j|}{|\hat{\beta}_j^{init}|}$

$\underbrace{\qquad\qquad}_{\text{Convex in } \beta}$

$$\hat{\beta}^{adaptive} \in \underset{\beta \in \mathbb{R}^p}{\arg\min} \left\{ \|Y - X\beta\|^2 + \lambda \sum_{j=1}^{p} \frac{|\beta_j|}{|\hat{\beta}_j^{init}|} \right\}$$

⟶ it is still convex and it reduces the bias problems

# Adaptive-Lasso