

# Convex relaxation

- Problem: solving the model selection minimization criterion

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \|Y - \hat{f}_m\|^2 + \operatorname{pen}(m) \sigma^2 \right\}$$

is impossible in practice when  $\mathcal{M}$  is very large.

Ex: coordinate sparse regression  $\mathcal{M} = \mathcal{P}(\{1, \dots, p\})$ , so we must  
 { evaluate  $|\mathcal{M}| = 2^p$  quantities  $\underbrace{\quad}_{!}$

- Today's recipe: modify the model selection criterion,  
 { in order to obtain a convex criterion, amenable to numerical  
 computations.

## ① Lasso estimator

- Sparse linear regression:

•  $Y = X\beta^* + \varepsilon$  with  $\|\beta^*\|_0$  small

• in all this lecture, we assume that the columns of  $X$  are normalized:  $\|X_j\| = 1$  for  $j=1, \dots, p$

- Model selection estimator:

•  $\mathcal{M} = \mathcal{P}(\{1, \dots, p\})$ ,  $S_m = \operatorname{Span}\{X_j : j \in m\}$ ,  $\hat{f}_m = \operatorname{Proj}_{S_m} Y$

•  $\pi_m = \exp(-|m| \log p)$

•  $\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \|Y - \hat{f}_m\|^2 + \lambda |m| \right\}$ , with  $\lambda = k(1 + \sqrt{2 \log p})^2 \sigma^2$

## Convexification:

$$\hat{m} \in \underset{m \in M}{\operatorname{argmin}} \{ \|Y - \hat{f}_m\|^2 + \lambda |m| \}$$

$$\text{we have } \hat{f}_m = X \hat{\beta}_m \text{ where } \hat{\beta}_m \in \underset{\beta: \operatorname{supp}(\beta) = m}{\operatorname{argmin}} \|Y - X\beta\|^2$$

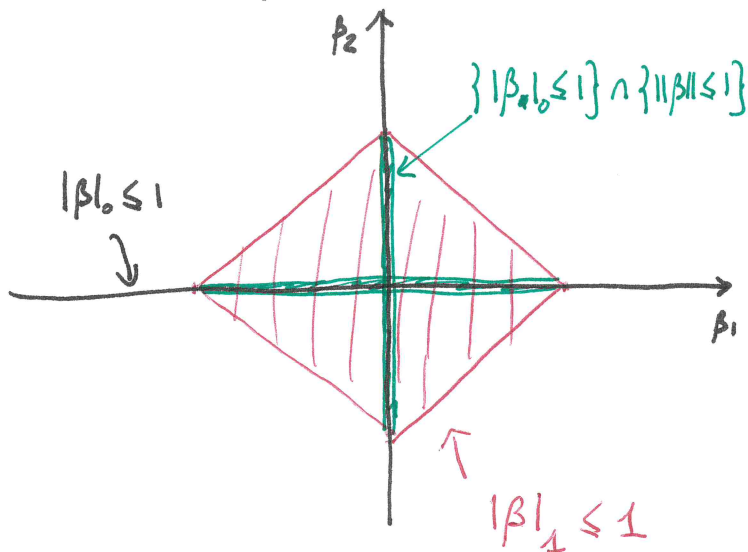
so

$$\hat{m} \in \underset{m \in M}{\operatorname{argmin}} \min_{\beta: \operatorname{supp}(\beta) = m} \{ \|Y - X\beta\|^2 + \lambda |\beta|_0 \}$$

and

$$\hat{\beta}_{\hat{m}} \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\{ \|Y - X\beta\|^2 \}}_{\text{convex } \checkmark} + \underbrace{\lambda |\beta|_0}_{\text{highly-non convex } \checkmark}$$

recipe



constrained version:

$$\min \|Y - X\beta\|^2$$

$$|\beta|_0 \leq D$$

convexification:

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

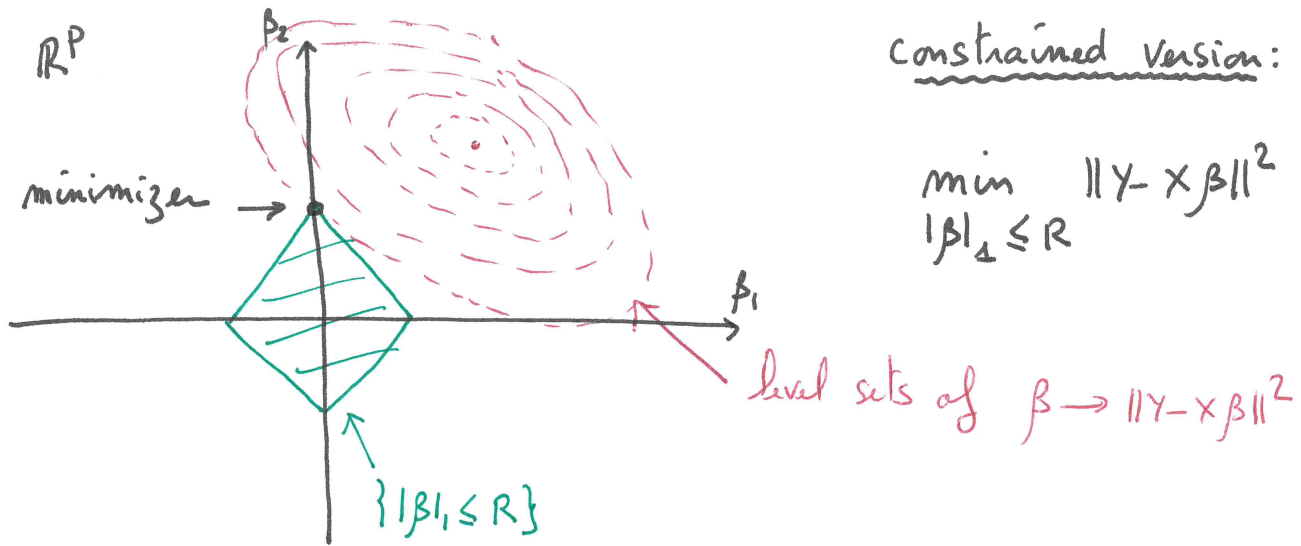
Lasso:

$$\hat{\beta}_\lambda \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \underbrace{\{ \|Y - X\beta\|^2 + \lambda |\beta|_1 \}}_{\text{convex}} \text{ for } \lambda > 0$$

$$=: \mathcal{L}_\lambda(\beta) \text{ convex}$$

$$\hat{f}_\lambda = X \hat{\beta}_\lambda$$

## Geometric interpretation



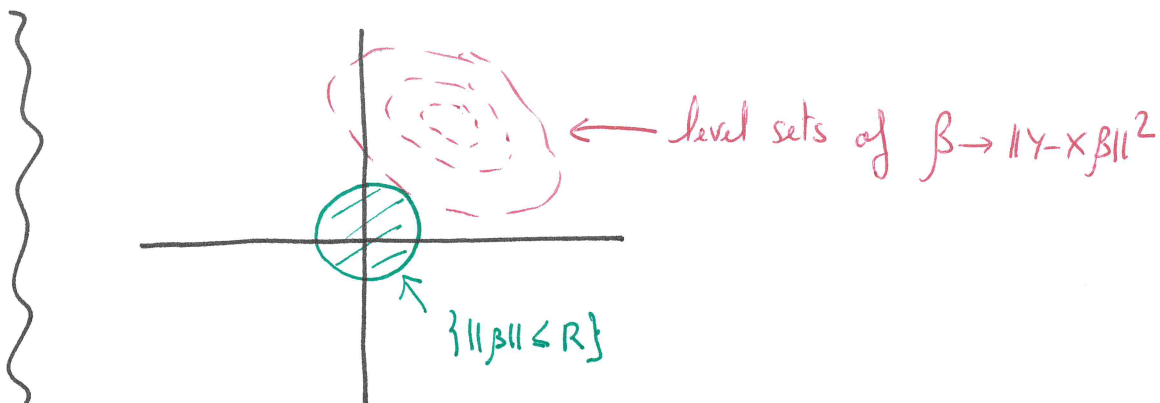
constrained version:

$$\min \|y - X\beta\|^2$$

$$\|\beta\|_1 \leq R$$

Singularities of  $\{\|\beta\|_1 \leq R\} \iff$  selection of variables

Remark: if we replace  $\|\beta\|_1$  by  $\|\beta\|^2$ , no selection occurs

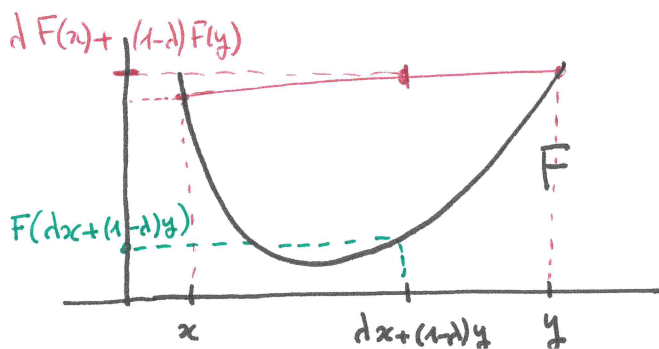


To do: exercise S.5.7 parts A) and B). (ridge & elastic net)

## Analytic interpretation

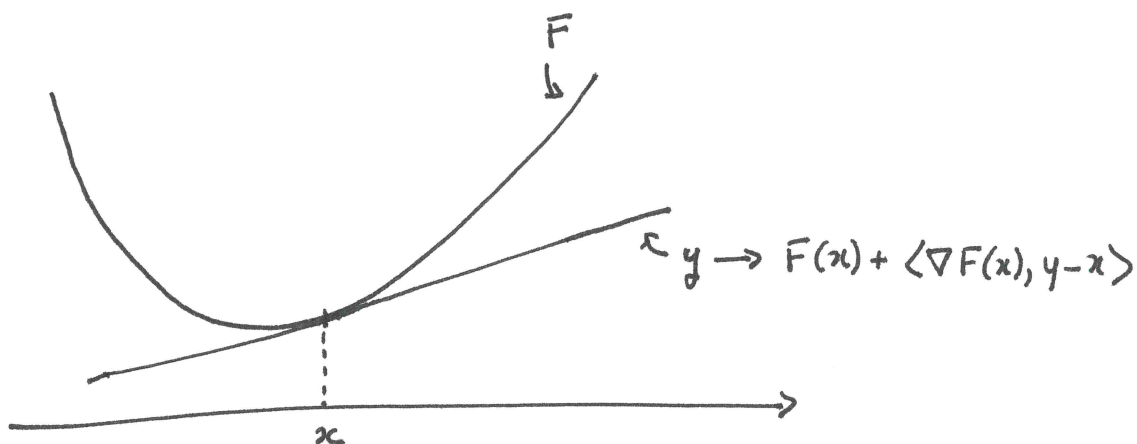
Reminder on convex functions and subdifferentials

$F: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex iff

$$\left\{ \begin{array}{l} F(\lambda x + (1-\lambda)y) \leq \lambda F(x) + (1-\lambda)F(y) \\ \forall \lambda \in [0, 1] \text{ and } \forall x, y \in \mathbb{R}^d \end{array} \right.$$


Lemma D1: if  $F$  is convex and differentiable in  $x$

$$F(y) \geq F(x) + \langle \nabla F(x), y-x \rangle, \quad \forall y \in \mathbb{R}^d$$



Subdifferential: if  $F$  is convex, we define the subdifferential

$$\partial F(x) = \left\{ \omega \in \mathbb{R}^d : F(y) \geq F(x) + \langle \omega, y-x \rangle \quad \forall y \in \mathbb{R}^d \right\}$$

↑  
subgradient

Lemma D2: Let  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  be a convex function

1.  $\partial F(x) \neq \emptyset$

2. if  $F$  is diff. in  $x$ :  $\partial F(x) = \{ \nabla F(x) \}$ .

Properties:  $F: \mathbb{R}^d \rightarrow \mathbb{R}$  convex

1. monotonicity:  $\forall \omega_x \in \partial F(x)$  and  $\forall \omega_y \in \partial F(y)$

$$\langle \omega_x - \omega_y, x - y \rangle \geq 0$$

2. minimum:

$$x^* \in \operatorname{argmin}_{x \in \mathbb{R}^d} F(x) \iff 0 \in \partial F(x^*)$$

Proof:

1. from the very definition:

$$F(y) \geq F(x) + \langle w_x, y-x \rangle$$

$$F(x) \geq F(y) + \langle w_y, x-y \rangle$$

(+)

---


$$0 \geq \langle w_x - w_y, y-x \rangle$$

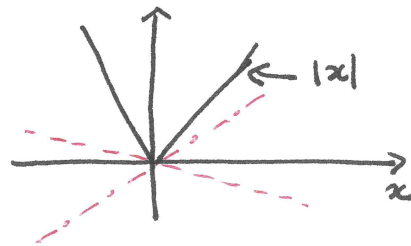
2- both statements are equivalent to

$$F(y) \geq F(x^*) + \langle 0, y-x^* \rangle \quad \forall y \in \mathbb{R}^d$$

□

Ex:  $\partial |x|_1$  ?

• dim = 1:  $|x|_1 = |x|$   
 so  $\partial |x| = \begin{cases} \text{sign}(x) & \text{if } x \neq 0 \\ [-1, 1] & \text{if } x = 0 \end{cases}$



• dim = d:  $|x|_1 = \sum_{j=1}^d |x_j|$

so  $\partial |x|_1 = \{z \in \mathbb{R}^d: z_j = \text{sign}(x_j) \text{ if } x_j \neq 0 \text{ and } z_j \in [-1, 1] \text{ if } x_j = 0\}$

• it is insightful to recover this result from a more principled way.

reminder:  $|x|_1 = \sup_{\|\phi\|_\infty \leq 1} \langle \phi, x \rangle$ .

we will prove  $\partial |x|_1 = \mathcal{D}_x$  where  $\mathcal{D}_x = \{\phi: \langle \phi, x \rangle = |x|_1, \|\phi\|_\infty \leq 1\}$

proof:

( $\supset$ ) for  $\phi \in \mathcal{D}_x$ :

$$|y|_1 \geq \underbrace{\langle \phi, y \rangle}_{\|\phi\|_\infty \leq 1} = |x|_1 + \underbrace{\langle \phi, y-x \rangle}_{\langle \phi, x \rangle = |x|_1}$$

Hence  $\phi \in \partial |x|_1$ .

(C) for  $\omega \in \partial |x|_1$ :

$$\left. \begin{array}{l} y = 2x: \quad 2|x|_1 \geq |x|_1 + \langle \omega, x \rangle \\ y = 0: \quad 0 \geq |x|_1 + \langle \omega, -x \rangle \end{array} \right\} \Rightarrow |x|_1 \leq \langle \omega, x \rangle \leq |x|_1$$

• we also have  $|\omega|_\infty = \langle \omega, z \rangle$  with  $|z|_1 = 1$

$$\text{so } |x|_1 + |z|_1 \stackrel{\Delta}{\geq} |x+z|_1 \stackrel{\omega \in \partial |x|_1}{\geq} |x|_1 + \underbrace{\langle \omega, z \rangle}_{= |\omega|_\infty} \Rightarrow |\omega|_\infty \leq 1$$

and hence  $\omega \in \mathcal{D}_x$ .

□

Lasso expression:

since  $\beta \rightarrow \mathcal{L}_\lambda(\beta) = \|Y - X\beta\|^2 + \lambda |\beta|_1$  is convex

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

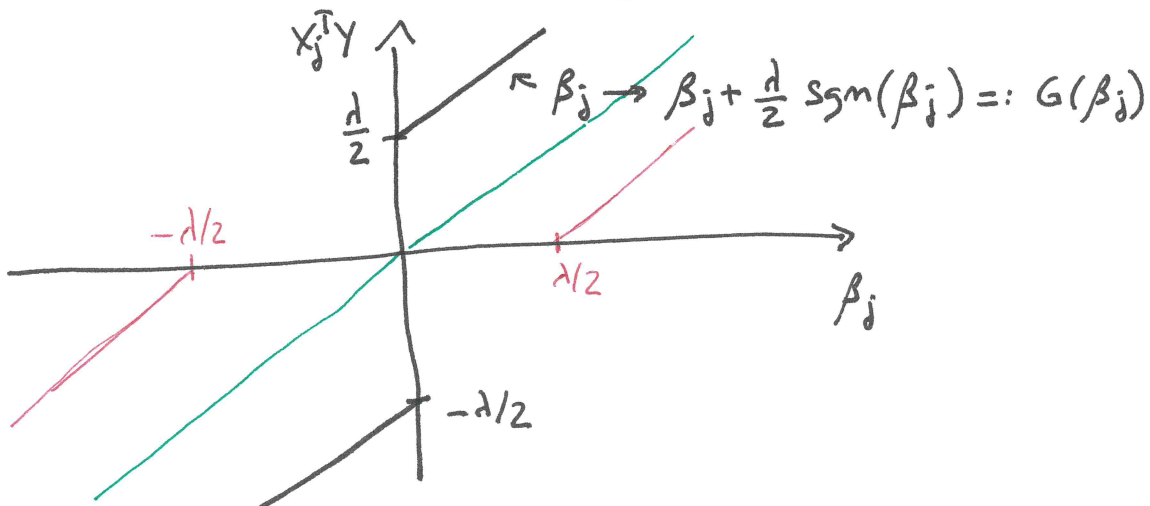
and  $\exists \hat{z} \in \partial |\hat{\beta}_\lambda|_1$  such that

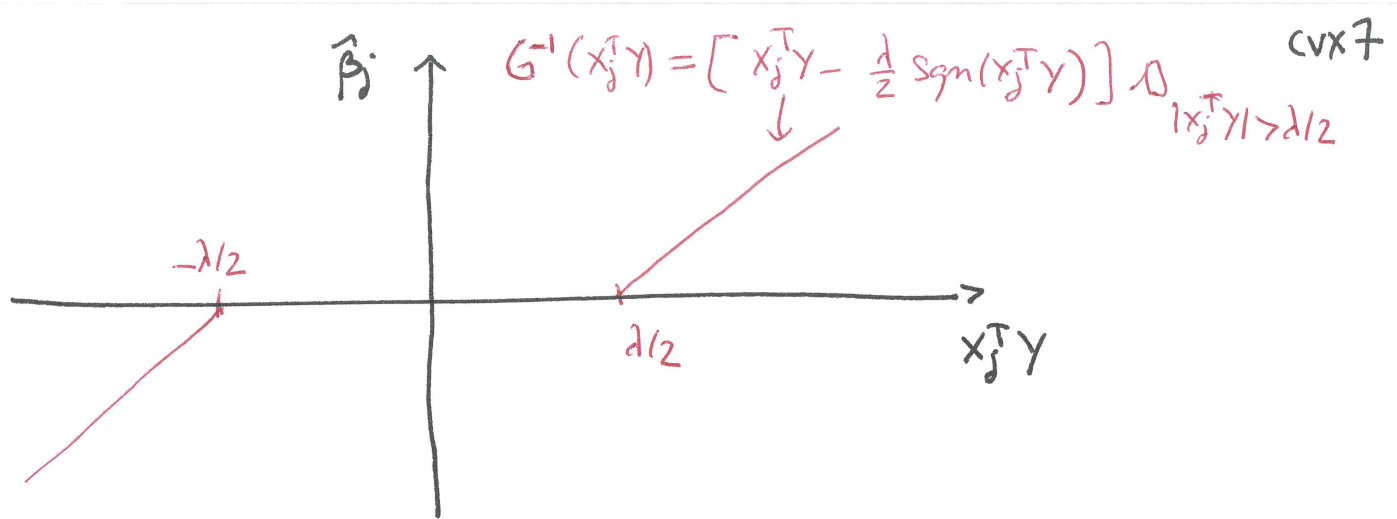
$$X^T X \hat{\beta}_\lambda = X^T Y - \frac{\lambda}{2} \hat{z}$$

no explicit expression, but in the orthogonal case.

• case  $X^T X = \text{Id}$ :  $\hat{\beta} = X^T Y - \frac{\lambda}{2} \hat{z}$ .

if  $\hat{\beta}_j \neq 0$ :  $\hat{\beta}_j = X_j^T Y - \frac{\lambda}{2} \text{sgn}(\hat{\beta}_j) \Rightarrow X_j^T Y = \hat{\beta}_j + \frac{\lambda}{2} \text{sgn}(\hat{\beta}_j)$





So: . if  $|x_j^T Y| > \lambda/2$ :  $\hat{\beta}_j = x_j^T Y - \frac{\lambda}{2} \text{sgn}(x_j^T Y)$   
 } . if  $|x_j^T Y| \leq \lambda/2$ :  $\hat{\beta}_j = 0$  and  $\hat{\gamma}_j = \frac{2}{\lambda} x_j^T Y$

compact formula: soft-thresholding

$$\hat{\beta}_j = x_j^T Y \left(1 - \frac{\lambda}{2|x_j^T Y|}\right)_+ \quad \text{for } j=1, \dots, p.$$

## ② Statistical analysis of Lasso estimator

Compatibility constant: account for (local) orthogonality

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

where .  $S = \text{supp}(\beta)$

.  $\mathcal{C}(\beta) = \{\sigma \in \mathbb{R}^p : \|\sigma_S\|_1 > \|\sigma_{S^c}\|_1\}$ .

Exercise: check that

1) if  $X^T X = I$  then  $\kappa(\beta) \geq 1$

2) we always have  $\kappa(\beta) \geq \lambda_{\min}(X^T X)^{1/2}$

### Theorem 5.1 Deterministic bound

For any  $\lambda > 3 \|X^T \varepsilon\|_\infty$ , we have

$$\|X \hat{\beta}_\lambda - X \beta^*\|^2 \leq \inf_{\beta \in \mathbb{R}^p} \left\{ \|X \beta - X \beta^*\|^2 + \frac{\lambda^2}{\kappa(\beta)^2} |\beta|_0 \right\}$$

### Corollary 5.3

Assume that  $\begin{cases} \cdot \|X_j\| = 1, \quad j=1, \dots, p \\ \cdot \varepsilon \sim \mathcal{N}(0, \sigma^2 I_m) \\ \cdot \lambda = 3\sigma \sqrt{2K \log p}, \quad \text{with } K > 1 \end{cases}$

then, with probability  $\geq 1 - \frac{1}{p^{K-1}}$

$$\begin{aligned} \|X \hat{\beta}_\lambda - X \beta^*\|^2 &\leq \inf_{\beta \in \mathbb{R}^p} \left\{ \|X \beta - X \beta^*\|^2 + \frac{18K \sigma^2 \log p}{\kappa(\beta)^2} |\beta|_0 \right\} \\ &\leq \inf_{m \in \{1, \dots, p\}} \left\{ \|X \hat{\beta}_m - X \beta^*\|^2 + \frac{18K \sigma^2 \log p}{\kappa(\hat{\beta}_m)^2} |m| \right\} \end{aligned}$$

↑  
price to pay for computational tractability

### Proof Cor. 5.3:

$$\cdot \|X^T \varepsilon\|_\infty = \max_{j=1, \dots, p} |X_j^T \varepsilon|, \quad \text{with } X_j^T \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

.. Hence

$$\mathbb{P} \left[ \|X^T \varepsilon\|_\infty > \sigma \sqrt{2K \log p} \right] \leq p \mathbb{P} \left[ |\mathcal{N}(0, 1)| > \sqrt{2K \log p} \right]$$

$$\begin{aligned} &\stackrel{\text{Lemma B.4.}}{\leq} p \exp\left(-\frac{1}{2}(2K \log p)\right) = \frac{1}{p^{K-1}} \end{aligned}$$

□



Proof of Theorem 5.1:

• optimality condition:  $0 \in \partial \mathcal{L}_\lambda(\hat{\beta}_\lambda)$

so  $\exists \hat{z} \in \partial |\beta|_1$  such that  $2X^T(X\hat{\beta}_\lambda - Y) + \lambda \hat{z} = 0$ .

•  $\forall \beta \in \mathbb{R}^p$ :

$$2 \langle X^T(X\hat{\beta}_\lambda - X\beta - \varepsilon), \hat{\beta}_\lambda - \beta \rangle + \lambda \langle \hat{z}, \hat{\beta}_\lambda - \beta \rangle = 0$$

i.e.  $2 \langle X\hat{\beta}_\lambda - X\beta^*, X\hat{\beta}_\lambda - X\beta \rangle - 2 \langle X^T\varepsilon, \hat{\beta}_\lambda - \beta \rangle + \lambda \langle \hat{z}, \hat{\beta}_\lambda - \beta \rangle = 0$

• Monotonicity:  $\forall z \in \partial |\beta|_1$ :  $\langle \hat{z}, \hat{\beta}_\lambda - \beta \rangle \geq \langle z, \hat{\beta}_\lambda - \beta \rangle$

so  $\forall \beta \in \mathbb{R}^p, \forall z \in \partial |\beta|_1$

$$\underbrace{2 \langle X(\hat{\beta}_\lambda - \beta^*), X(\hat{\beta}_\lambda - \beta) \rangle}_{=: \mathcal{A}(\hat{\beta}_\lambda)} \leq 2 \langle X^T\varepsilon, \hat{\beta}_\lambda - \beta \rangle - \lambda \langle z, \hat{\beta}_\lambda - \beta \rangle \quad (1)$$

• Lemma 5.2:  $S = \text{supp}(\beta)$

$$\mathcal{A}(\hat{\beta}_\lambda) \stackrel{(i)}{\leq} \frac{\lambda}{3} \left( 5 \underbrace{|\hat{\beta}_\lambda - \beta|_{S^c}}_{(ii)} - |\hat{\beta}_\lambda - \beta|_S \right) \leq 2\lambda |\hat{\beta}_\lambda - \beta|_S$$

Proof Lemma 5.2: • define  $z$  by: •  $z_j = \text{sign}(\beta_j)$  for  $j \in S$

•  $z_j = \text{sign}(\hat{\beta}_j - \beta_j)$  for  $j \in S^c$

• by construction  $z \in \partial |\beta|_1$ :

$$\mathcal{A}(\hat{\beta}_\lambda) \stackrel{(1)}{\leq} 2 \underbrace{\|X^T\varepsilon\|_\infty}_{\leq \lambda/3} \|\hat{\beta}_\lambda - \beta\|_1 - \lambda \langle z_S, (\hat{\beta}_\lambda - \beta)_S \rangle - \lambda \langle z_{S^c}, (\hat{\beta}_\lambda - \beta)_{S^c} \rangle$$

↑  
hypothesis of Theorem 5.1.

From the definition of  $\gamma$ :

CVX10

$$\begin{aligned} A(\hat{\beta}_\lambda) &\leq \frac{2\lambda}{3} \|\hat{\beta}_\lambda - \beta\|_1 + \lambda \|(\hat{\beta}_\lambda - \beta)_S\|_1 - \lambda \|(\hat{\beta}_\lambda - \beta)_{S^c}\|_1 \\ &= \frac{\lambda}{3} (5 \|(\hat{\beta}_\lambda - \beta)_S\|_1 - \|(\hat{\beta}_\lambda - \beta)_{S^c}\|_1) \quad \leftarrow \text{gives (i)} \\ &\leq 2\lambda \|(\hat{\beta}_\lambda - \beta)_S\|_1 \quad \leftarrow \text{gives (ii)}. \end{aligned}$$

□

We can conclude the proof of Theorem 5.1.

Al-Kashi

$$\rightarrow \text{if } A(\hat{\beta}_\lambda) \stackrel{\text{Al-Kashi}}{=} \|X(\hat{\beta}_\lambda - \beta^*)\|^2 + \|X(\hat{\beta}_\lambda - \beta)\|^2 - \|X(\beta - \beta^*)\|^2 \leq 0$$

$$\text{then } \|X(\hat{\beta}_\lambda - \beta^*)\|^2 \leq \|X(\beta - \beta^*)\|^2 \quad \therefore$$

$\rightarrow$  if  $A(\hat{\beta}_\lambda) > 0$ , then  $\hat{\beta}_\lambda - \beta \in \mathcal{C}(\beta)$  from (i), and from ~~the~~ (ii)

$$A(\hat{\beta}_\lambda) = \|X(\hat{\beta}_\lambda - \beta^*)\|^2 + \|X(\hat{\beta}_\lambda - \beta)\|^2 - \|X(\beta - \beta^*)\|^2$$

$$\stackrel{(ii)}{\leq} 2\lambda \|(\hat{\beta}_\lambda - \beta)_S\|_1$$

$$\begin{aligned} \text{definition of } \kappa(\beta) &\rightarrow \leq 2\lambda \frac{\sqrt{|\beta|_0}}{\kappa(\beta)} \|X(\hat{\beta}_\lambda - \beta)\| \leq \frac{\lambda^2 |\beta|_0}{\kappa(\beta)^2} + \|X(\hat{\beta}_\lambda - \beta)\|^2 \\ &\quad \uparrow \\ &\quad 2ab \leq a^2 + b^2 \end{aligned}$$

$$\Rightarrow \|X(\hat{\beta}_\lambda - \beta^*)\|^2 \leq \|X(\beta - \beta^*)\|^2 + \frac{\lambda^2 |\beta|_0}{\kappa(\beta)^2} \quad \therefore$$

□

### ③ Computing the Lasso

a/ Algebraic computation (LARS algorithm)

- $S_\lambda = \text{supp}(\hat{\beta}_\lambda)$

- optimality:  $\exists z_\lambda$  such that 
$$\begin{cases} [z_\lambda]_{S_\lambda} = \text{sign}([ \hat{\beta}_\lambda ]_{S_\lambda}) \\ \| [z_\lambda]_{S_\lambda^c} \|_\infty \leq 1 \end{cases}$$

and

$$X^T X \hat{\beta}_\lambda = X^T Y - \frac{\lambda}{2} z_\lambda$$

$$\Rightarrow \begin{cases} \text{on } S_\lambda: & X_{S_\lambda}^T X_{S_\lambda} [\hat{\beta}_\lambda]_{S_\lambda} = X_{S_\lambda}^T Y - \frac{\lambda}{2} \text{sign}([ \hat{\beta}_\lambda ]_{S_\lambda}) \\ \text{on } S_\lambda^c: & \| X_{S_\lambda^c}^T Y - X_{S_\lambda^c}^T X \hat{\beta}_\lambda \|_\infty \leq \frac{\lambda}{2} \end{cases}$$

- Since  $\lambda \rightarrow S_\lambda$  is piecewise constant,  $\lambda \rightarrow \hat{\beta}_\lambda$  is piecewise linear

$\leadsto$  see Figure 5.2

- LARS algorithm computes algebraically the break points  $\hat{\beta}_{\lambda_1}, \hat{\beta}_{\lambda_2}, \dots$

Then, for  $\lambda \in (\lambda_{k+1}, \lambda_k)$ ,  $\hat{\beta}_\lambda$  is computed by linear interpolation.



- while the computations are of algebraic nature, we do not have explicit formulas for  $\hat{\beta}_\lambda$



- matrix inversions are done up to some precision level

b/ Accelerated proximal method

Consider  $\min_{\beta} \{ F(\beta) + \lambda \|\beta\|_1 \}$  with  $F$  convex and smooth  
 (in our case  $F(\beta) = \|Y - X\beta\|^2$ )



$$F(\beta) = F(\beta^t) + \langle \nabla F(\beta^t), \beta - \beta^t \rangle + O(\|\beta - \beta^t\|^2)$$

iterate  $\Rightarrow \beta^{t+1} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ F(\beta^t) + \langle \nabla F(\beta^t), \beta - \beta^t \rangle + \frac{1}{2\eta} \|\beta - \beta^t\|^2 + \lambda \|\beta\|_1 \right\}$

Why easier?

$$1) \beta^{t+1} \in \operatorname{argmin}_{\beta} \left\{ \frac{1}{2\eta} \|\beta - \beta^t + 2 \nabla F(\beta^t)\|^2 + \underbrace{F(\beta^t) - \frac{2}{\eta} \|\nabla F(\beta^t)\|^2}_{\text{no } \beta \text{ here}} + \lambda \|\beta\|_1 \right\}$$

$$= \operatorname{argmin}_{\beta} \left\{ \frac{1}{2} \|\beta - (\beta^t - 2 \nabla F(\beta^t))\|^2 + \lambda \|\beta\|_1 \right\}$$

2) we have seen that  $S_{\lambda}(\alpha) := \begin{bmatrix} \alpha_1 (1 - \lambda/|\alpha_1|)_+ \\ \vdots \\ \alpha_p (1 - \lambda/|\alpha_p|)_+ \end{bmatrix}$  is solution

to  $\operatorname{argmin}_{\beta} \left\{ \frac{1}{2} \|\beta - \alpha\|^2 + \lambda \|\beta\|_1 \right\}$

Soft thresholding operator

$$\Rightarrow \boxed{\beta^{t+1} = S_{\lambda \eta}(\beta^t - 2 \nabla F(\beta^t))}$$

(here  $\nabla F(\beta^t) = 2X^T(X\beta^t - Y)$ )

Remarks:

- 1) What do you recognize for  $\lambda = 0$  ?
- 2) the scheme converges to  $\beta_{\lambda}$  for  $\eta$  small enough
- 3) acceleration: we can accelerate the convergence with Nesterov's acceleration scheme  $\rightarrow$  FISTA algorithm (see Section 5.2.4)

c/ Other algorithms

- There exists many other algorithms.
- A simple one: coordinate descent algorithm

→ To Do: exercise 5.5.7

④ Bias of the Lasso

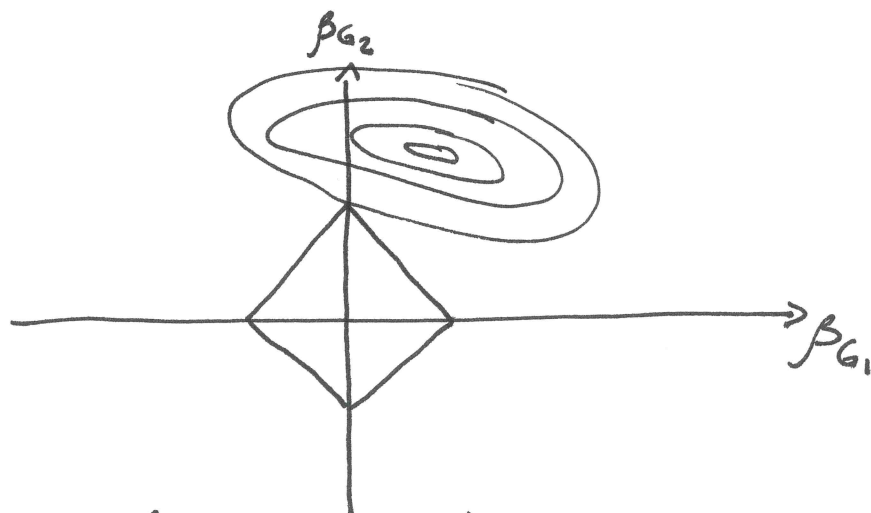
→ look at the slides

⑤ Other sparsity structures

Ex: group sparsity

$$\beta^* = \begin{bmatrix} \beta_{G_1}^* \\ \vdots \\ \beta_{G_n}^* \end{bmatrix} \left. \begin{array}{l} \} G_1 \\ \vdots \\ \} G_n \end{array} \right\} \text{with a few } \beta_{G_k}^* \neq 0$$

Coordinate sparse	group sparse
$\text{card}\{j: \beta_j \neq 0\}$ small	$\text{card}\{k: \beta_{G_k} \neq 0\}$ small
↓	↓
$\sum_{j=1}^p  \beta_j $ penalisation	$\sum_{k=1}^m \ \beta_{G_k}\ $



Group Lasso:  $\hat{\beta}_\lambda \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|^2 + \sum_k \lambda_k \|\beta_{G_k}\| \right\}$

→ again a convex criterion

→ usually  $d_k = \lambda \sqrt{|G_k|}$

• Exercise: following similar arguments as for  $l_2 l_1$  check that

$$\left. \partial \left( \sum_k d_k \|\beta_{G_k}\| \right) = \left\{ \beta \in \mathbb{R}^p : \begin{array}{l} \cdot \beta_{G_k} = d_k \frac{\beta_{G_k}}{\|\beta_{G_k}\|} \quad \text{if } \beta_{G_k} \neq 0 \\ \cdot \|\beta_{G_k}\| \leq d_k \quad \text{if } \beta_{G_k} = 0 \end{array} \right\} \right\}$$

## ⑥ Take home message

### Convex relaxation

→ a principled way to derive practical estimators from model selection

→ good theoretical guarantees

→ but suffers from some bias

---

Model selection is too hard?

Just relax!