# Lecture 7 bis : Methods for Regression
## Lasso estimator

K. Meziani

**Ðauphine** | PSL✻
UNIVERSITÉ PARIS

## Model

Consider the $n$-sample $(\{(x_i, y_i) \in \mathbb{R}^p \times \mathbb{R}\}_{i=1}^n$ such that

$$y_i = x_i^\top \beta^* + \epsilon_i \quad, i = 1, \ldots, n$$

where $\epsilon_i \overset{iid}{\sim} \mathcal{N}(0, \sigma^2)$, $x_i \in \mathbb{R}^p$ and $\beta^* \in \mathbb{R}^p$. Fix $p \geq 2$ and $n \geq 1$.

**Matrix form:**

$$Y = \sum_{j=1}^n \beta_j^* X_j + \epsilon = X\beta^* + \epsilon, \tag{1}$$

where $X_j \in \mathbb{R}^n$, $X = [X_1| \cdots |X_p]$ and $\epsilon \sim \mathcal{N}(\mathbf{0}_n, \sigma^2 \mathbb{I}_n)$.

# Remark

This lecture is based on the paper of Bickel, Ritov and Tsybakov,*Simultaneous analysis of Lasso and Dantzig selector* (2009, AOS).

**Sparsity set.** For all $\beta \in \mathbb{R}^p$, we denote by $J(\beta)$ the sparsity set, the subset of indices $\{1, \cdots, p\}$ where the vector $\beta$ has non-zero coordinates

$$J(\beta) = \{j : \beta_j \neq 0\}$$

**Sparsity of $\beta$.** The sparsity of the vector $\beta$ is characterized by the value $M(\beta)$, the cardinality of $J(\beta)$:

$$M(\beta) = \sum_{j=1}^{p} \mathbb{1}_{\beta_j \neq 0} = |J(\beta)|$$

## Notations

**Gram matrix:** $\Psi_n = \frac{X^T X}{n} = \frac{1}{n} \sum_{i=1}^{n} x_{i,j}^2$

**Some norms.** For all $a \in R^n$, $b \in \mathbb{R}^p$ and $J_0 \subseteq \{1, ..., p\}$, we denote

$$\|a\|_n^2 = \frac{1}{n} \sum_{i=1}^{n} a_i^2 \quad , \quad \|b\|_1 = \sum_{j=1}^{p} |b_j| \quad , \quad \|b\|_{2, J_0} = \sqrt{\sum_{j \in J_0} b_j^2} \quad \text{and}$$
$$\|b\|_{1, J_0} = \sum_{j \in J_0} |b_j|.$$

# Assumptions $\mathcal{H}$

- The Gram matrix $\Psi_n$ is such that its diagonal elements are equal to 1.

- The sparsity index is such that $M(\beta^*) \leq s$ for $1 \leq s \leq p$.

# Lasso

> **Definition** **LASSO**
>
> $$\widehat{\beta}^{\lambda,L} = \arg\min_{\beta \in \mathbb{R}^p} \left\{ \|Y - X\beta\|_n^2 + 2\lambda\|\beta\|_1 \right\}$$
>
> where the regularization parameter $\lambda > 0$.

## Remarks

We are typically interested in the case where $p > n$ and even $p >> n \Rightarrow$ Gram matrix $\Psi_n$ is degenerate, *i.e*

$$\min_{\substack{\delta \in \mathbb{R}^p \\ \delta \neq 0}} \frac{(\delta^\top \psi_n \delta)^{1/2}}{\sqrt{n}\|\delta\|_2} \equiv \min_{\substack{\delta \in \mathbb{R}^p \\ \delta \neq 0}} \frac{\|X\delta\|_2}{\sqrt{n}\|\delta\|_2} = 0$$

OLSE does not work in this case, since it requires positive definiteness of $\psi_n$, *i.e*

$$\min_{\substack{\delta \in \mathbb{R}^p \\ \delta \neq 0}} \frac{\|X\delta\|_2}{\sqrt{n}\|\delta\|_2} > 0$$

The Lasso require much weaker assumptions. Replace

- the min by the minimum over a restricted set of vectors
- the norm $\|\delta\|_2$ by the $\ell_2$ norm of only a part of $\delta$.

For $J_0 = J(\beta^*)$, one of the properties of the Lasso is that the residuals $\delta = \widehat{\beta^{\lambda,L}} - \beta^*$ satisfy, with probability close to 1,

$$\|\delta\|_{1,J_0^c} \leq 3\|\delta\|_{1,J_0}.$$

# RE($s, c_0$) Assumption

Restricted Eigenvalues
Sparsity
Constant ( 1 comme un si-desfous)

For an integer $s$ such that $1 \le s \le p$ and a positive number $c_0$, we assume that the following condition is satisfied:

$$\kappa(s, c_0) := \min_{\substack{J_0 \subseteq \{1,\dots,p\} \\ |J_0| \le s}} \min_{\substack{\delta \ne 0 \\ \|\delta\|_{1,J_0^c} \le c_0 \|\delta\|_{1,J_0}}} \frac{\|X\delta\|_2}{\sqrt{n}\|\delta\|_{2,J_0}} > 0.$$

plus petit vp de la matrice restreinte à l'espace $J_0$.

hypothèse faible
( ce qu'on recherche
on recherche )

# Theorem~1

[a]Consider model (1). Under $\mathcal{H}$ and RE($s$, 3) Assumptions, the estimator $\widehat{\beta^{\lambda,L}}$, for $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, with $A > 2\sqrt{2}$ is such that with probability greater than $1 - p^{1-\frac{A^2}{8}}$ ::

$$\|X(\widehat{\beta^{\lambda,L}} - \beta^*)\|_2^2 \leq \frac{16A^2}{\kappa(s,3)^2}\sigma^2 s \log(p). \qquad (2)$$

*petit*

$$\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1 \leq \frac{16A}{\kappa(s,3)^2}\sigma s \sqrt{\frac{\log(p)}{n}}, \qquad (3)$$

$$M(\widehat{\beta^{\lambda,L}}) \leq \frac{64\phi_{\max}}{\kappa(s,3)^2}, \qquad (4)$$

→ *valeur propre maximale*

*Inégalités de concentration*

where $\phi_{\max}$ denote the maximal eigenvalue of the Gram matrix $\psi_n$.

---

[a]Bickel, Ritov and Tsybakov,*Simultaneous analysis of Lasso and Dantzig selector* (2009, AOS)

To prove this theorem, we need previous results.

Let $V_j := \frac{1}{n} \sum_{i=1}^{n} \epsilon_i x_{i,j} = \frac{1}{n} X_j^\top \epsilon$, we define the set

$$\Omega = \bigcap_{j=1}^{p} \left\{ 2|V_j| \le \lambda \right\} = \bigcap_{j=1}^{p} \left\{ \left| \frac{1}{n} X_j^\top \epsilon \right| \le \lambda/2 \right\}$$

*espace où le temps amène de sélectionner des variables sur trop variables au bruit.*

*les variables n'ont pas été prises*

Then for $A > 2\sqrt{2}$ and $\lambda = A\sigma \sqrt{\frac{\log(p)}{n}}$, we have

$$\mathbb{P}(\Omega) > 1 - p^{1 - \frac{A^2}{8}},$$

# Proof of Lemma~1

Denote by $\eta = \sqrt{n}V_j = \frac{1}{\sqrt{n}} \sum_{i=1}^{\sqrt{n}} x_{i,j}\epsilon_i$, then

$$\mathbb{E}(\eta) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} x_{i,j}\mathbb{E}(\epsilon_i) = 0$$

$$\mathbb{V}ar(\eta) = \frac{1}{n} \sum_{i=1}^{n} x_{i,j}^2 \mathbb{V}ar(\epsilon_i) = \frac{\sigma^2}{n} \sum_{i=1}^{n} x_{i,j}^2 = \frac{\sigma^2}{n} X_j^\top X_j = \sigma^2 \quad \text{by Assumption } \mathcal{H}$$

> *La éléments diagonaux*

Then $\frac{\eta}{\sigma} \sim \mathcal{N}(0,1)$ and

> *union des complémentaires*

$$\mathbb{P}(\Omega^c) = \mathbb{P}\left(\bigcup_{j=1}^{p} \{2|V_j| > \lambda\}\right) = \mathbb{P}\left(\bigcup_{j=1}^{p} \{2|\eta| > \lambda \sqrt{n}\}\right) \leq \sum_{j=1}^{p} \mathbb{P}\left(2|\eta| > \lambda \sqrt{n}\right)$$

$$\leq p\mathbb{P}\left(2|\eta| > \lambda \sqrt{n}\right) = p\mathbb{P}\left(|\frac{\eta}{\sigma}| > \frac{\lambda \sqrt{n}}{2\sigma}\right) \leq pe^{-\frac{\lambda^2 n}{8\sigma^2}} = pe^{-A^2 \frac{\log(p)}{8}} = p^{1-A^2/8}$$

> $Z \sim \mathcal{N}(0,1)$
> $\mathbb{P}(|Z| > t) \leq e^{-\frac{t^2}{2}}$

*Donc $\Omega$ est vrai avec grande probabilité.*

Consider model (1). Under $\mathcal{H}$ Assuptions, The estimator $\widehat{\beta^{\lambda,L}}$ for $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, with $A > 2\sqrt{2}$ is such that with probability greater than $1 - p^{1-\frac{A^2}{8}}$ :

$$\|X(\widehat{\beta^{\lambda,L}} - \beta^*)\|_n^2 + \lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1 \leq 4\lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0}, \qquad (5)$$

where $J_0 = J(\beta^*)$.

*The proof will be given last*

$\delta = \hat{\beta}^L - \beta^*$

Lemme 2 : $\|X\delta\|_n^2 + \lambda\|\delta\|_1 \leq 4\lambda\|\delta\|_{1,J_0}$

With probability $> 1 - p^{1 - \frac{A^2}{8}}$, we have

$$\lambda \|\widehat{\beta}^{\lambda, L} - \beta^*\|_1 \leq \overbrace{\|X(\widehat{\beta}^{\lambda, L} - \beta^*)\|_n^2}^{> 0} + \lambda \|\widehat{\beta}^{\lambda, L} - \beta^*\|_1$$

$$\Rightarrow \quad \lambda \|\widehat{\beta}^{\lambda, L} - \beta^*\|_1 \leq 4\lambda \|\widehat{\beta}^{\lambda, L} - \beta^*\|_{1, J_0} \qquad \text{Lemma 2}$$

$$\Rightarrow \quad \|\widehat{\beta}^{\lambda, L} - \beta^*\|_1 \leq 4 \|\widehat{\beta}^{\lambda, L} - \beta^*\|_{1, J_0} \qquad /\lambda$$

$$\Rightarrow \quad \|\widehat{\beta}^{\lambda, L} - \beta^*\|_{1, J_0} + \|\widehat{\beta}^{\lambda, L} - \beta^*\|_{1, J_0^c} \leq 4 \|\widehat{\beta}^{\lambda, L} - \beta^*\|_{1, J_0}$$

*somme sur $J_0$*        *somme sur $J_0^c$*

---

**Then**

$$\|\widehat{\beta}^{\lambda, L} - \beta^*\|_{1, J_0^c} \leq 3 \|\widehat{\beta}^{\lambda, L} - \beta^*\|_{1, J_0}. \tag{6}$$

$\|x\|_{1, s} \leq 3 \|x\|_{1, s^c}$

# KKT Conditions

Consider (1). Under $\mathcal{H}$ assumption, any Lasso solution $\widehat{\beta^L}$ satisfies the following necessary and sufficient condition (KKT)

$$\begin{cases} \frac{1}{n} X_j^\top (Y - X\widehat{\beta^{\lambda,L}}) = \lambda \, \text{sign}(\widehat{\beta_j^{\lambda,L}}) & \text{if} \quad \widehat{\beta_j^{\lambda,L}} \neq 0 \\ \frac{1}{n} |X_j^\top (Y - X\widehat{\beta^{\lambda,L}})| \leq \lambda & \text{if} \quad \widehat{\beta_j^{\lambda,L}} = 0 \end{cases}$$

# KKT Conditions - Proof

$$\widehat{\beta}^{\lambda,L} = \arg\min_{\beta \in \mathbb{R}^p} \left\{ \underbrace{\|Y - X\beta\|_n^2 + 2\lambda\|\beta\|_1}_{\text{on cherche la dérivée}} \right\} := \arg\min_{\beta \in \mathbb{R}^p} f(\beta),$$

The lasso solution is not differentiable at any point where $\beta_j$ is equal to zero. ($\rightarrow$ use the concept of a *subdifferential*).

- The subdifferential of the $\ell_1$ penalty at point $\beta$ is the vector with components

$$\partial(|\cdot|)(\beta_j) = \begin{cases} \text{sign}(\beta_j), & \text{if} \quad \beta_j \neq 0 \\ v_j \in [-1, 1], & \text{if} \quad \beta_j = 0 \end{cases}$$

- We know that the gradient of the first terms at point $\beta$ is

$$-\frac{2}{n} X^\top (Y - X\beta).$$

Any Lasso solution $\widehat{\beta}^L$ satisfies the following necessary and sufficient condition (KKT)

$$\vec{0} \in \partial(f)(\widehat{\beta}^L) \Leftrightarrow \text{KKT Conditions}$$

[a]Consider model (1). Under $\mathcal{H}$ and RE($s, 3$) Assumptions, the estimator $\widehat{\beta}^{\lambda,L}$, for $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, with $A > 2\sqrt{2}$ is such that with probability greater than $1 - p^{1-\frac{A^2}{8}}$ ::

$$\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_2^2 \leq \frac{16A^2}{\kappa(s,3)^2}\sigma^2 s \log(p).$$

$$\|\widehat{\beta}^{\lambda,L} - \beta^*\|_1 \leq \frac{16A}{\kappa(s,3)^2}\sigma s \sqrt{\frac{\log(p)}{n}},$$

$$M(\widehat{\beta}^{\lambda,L}) \leq \frac{64\phi_{\max}}{\kappa(s,3)^2},$$

where $\phi_{\max}$ denote the maximal eigenvalue of the Gram matrix $\psi_n$.

---

[a]Bickel, Ritov and Tsybakov,*Simultaneous analysis of Lasso and Dantzig selector* (2009, AOS)

# Proof of Theorem~1 : the prediction bound

By Lemma~2, for $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, with $A > 2\sqrt{2}$ is s.t. w.h.p.

$$\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n^2 + \lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_1 \leq \underbrace{4\lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{1,J_0}}_{=\delta}$$

$$\sum_{j\in J_0}|\delta_j| \leq \sqrt{\Sigma \cdot 1}\sqrt{\Sigma \delta_j^2}$$

By Cauchy Schwarz, and for $M(\beta^*) \leq s$:

$$\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n^2 + \lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_1 \leq 4\sqrt{s}\lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{2,J_0} \qquad (7)$$

$$\kappa \leq \frac{\|X\delta\|_2}{\sqrt{n}\|\delta\|_{2,J_0}}$$

By RE(s,3) assumption, for $\delta = \widehat{\beta}^{\lambda,L} - \beta^*$ and $\kappa := \kappa(s,3)$

$$(\Leftrightarrow) \; \kappa\|\delta\|_{2,J_0} \leq \frac{\|X\delta\|_2}{\sqrt{n}}$$

$$\kappa\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{2,J_0} \leq \frac{1}{\sqrt{n}}\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_2 = \|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n. \qquad (8)$$

Then using (7) and (8)

$$\begin{aligned}
\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n^2 &\leq \|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n^2 + \lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_1 \leq 4\sqrt{s}\lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{2,J_0} \\
&\leq \frac{4\sqrt{s}\lambda}{\kappa}\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n.
\end{aligned}$$

By simplifying by $\|X\widehat{\beta}^{\lambda,L} - \beta^*\|_n$ and squaring, we have

$$\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n^2 \leq \frac{16s\lambda^2}{\kappa^2}. \qquad (9)$$

For $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, we get (2).

# Theorem~1

[a]Consider model (1). Under $\mathcal{H}$ and RE$(s, 3)$ Assumptions, the estimator $\widehat{\beta}^{\lambda, L}$, for $\lambda = A\sigma \sqrt{\frac{\log(p)}{n}}$, with $A > 2\sqrt{2}$ is such that with probability greater than $1 - p^{1 - \frac{A^2}{8}}$ ::

$$\|X(\widehat{\beta}^{\lambda, L} - \beta^*)\|_2^2 \leq \frac{16A^2}{\kappa(s, 3)^2} \sigma^2 s \log(p).$$

$$\|\widehat{\beta}^{\lambda, L} - \beta^*\|_1 \leq \frac{16A}{\kappa(s, 3)^2} \sigma s \sqrt{\frac{\log(p)}{n}},$$

$$M(\widehat{\beta}^{\lambda, L}) \leq \frac{64\phi_{\max}}{\kappa(s, 3)^2},$$

where $\phi_{\max}$ denote the maximal eigenvalue of the Gram matrix $\psi_n$.

---

[a]Bickel, Ritov and Tsybakov,*Simultaneous analysis of Lasso and Dantzig selector* (2009, AOS)

## Proof of Theorem~1: the estimation bound

Under RE(s,3), for $\delta = \widehat{\beta}^{\lambda,L} - \beta^*$ and $\kappa := \kappa(s, 3)$,

$$\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{2,J_0} \leq \frac{1}{\kappa}\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n. \tag{10}$$

Using the previous prediction bound (9)   $\|X\delta\|_n \leq \frac{4\sqrt{s}\lambda}{\kappa}$

$$\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{2,J_0} \leq \frac{1}{\kappa}\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n \leq \frac{4\sqrt{s}\lambda}{\kappa^2}. \tag{11}$$

As $\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{1,J_0^c} \leq 3\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{1,J_0}$, a consequence of Lemma~2

$$
\begin{aligned}
\|\widehat{\beta}^{\lambda,L} - \beta^*\|_1 &= \|\widehat{\beta}^{\lambda,L} - \beta^*\|_{1,J_0} + \|\widehat{\beta}^{\lambda,L} - \beta^*\|_{1,J_0^c} \leq 4\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{1,J_0} \\
&\leq 4\sqrt{s}\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{2,J_0} \quad \text{Cauchy Sch. inequ.} \\
&\leq \frac{16s\lambda}{\kappa^2} \quad \text{by eq.(11)}
\end{aligned}
$$

for $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, we get (3).

# Theorem~1

[a]Consider model (1). Under $\mathcal{H}$ and RE($s$, 3) Assumptions, the estimator $\widehat{\beta}^{\lambda,L}$, for $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, with $A > 2\sqrt{2}$ is such that with probability greater than $1 - p^{1-\frac{A^2}{8}}$ ::

$$\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_2^2 \leq \frac{16A^2}{\kappa(s,3)^2}\sigma^2 s \log(p).$$

$$\|\widehat{\beta}^{\lambda,L} - \beta^*\|_1 \leq \frac{16A}{\kappa(s,3)^2}\sigma s \sqrt{\frac{\log(p)}{n}},$$

$$M(\widehat{\beta}^{\lambda,L}) \leq \frac{64\phi_{\max}}{\kappa(s,3)^2},$$

where $\phi_{\max}$ denote the maximal eigenvalue of the Gram matrix $\psi_n$.

---

[a]Bickel, Ritov and Tsybakov,*Simultaneous analysis of Lasso and Dantzig selector* (2009, AOS)

# Proof of Theorem~1: the estimation support bound

On $\Omega$, w.h.p *with high proba*

$$\left|\frac{1}{n}X_j^\top \epsilon\right| \leq \lambda/2, \quad \forall j = 1, \cdots, p \tag{12}$$

Moreover, the Lasso estimator $\widehat{\beta}^{\lambda,L}$ satisfies the KKT condition:

$$\begin{cases} \frac{1}{n}X_j^\top(Y - X\widehat{\beta}^{\lambda,L}) = \lambda\,\text{sign}(\widehat{\beta}_j^{\lambda,L}) & \text{if} \quad \widehat{\beta}_j^{\lambda,L} \neq 0 \\ \frac{1}{n}|X_j^\top(Y - X\widehat{\beta}^{\lambda,L})| \leq \lambda & \text{if} \quad \widehat{\beta}_j^{\lambda,L} = 0 \end{cases}$$

Then, we have $\quad \left|\dfrac{1}{n}X_j^\top(Y - X\widehat{\beta}^{\lambda,L})\right| = \lambda \quad \text{if} \quad \widehat{\beta}_j^{\lambda,L} \neq 0 \tag{13}$

Combining (12) and (13), and as $X\beta^* = Y - \epsilon$, we have

$$\begin{aligned} \left|\frac{1}{n}X_j^\top(X\beta^* - X\widehat{\beta}^{\lambda,L})\right| &= \left|\frac{1}{n}X_j^\top(Y - X\widehat{\beta}^{\lambda,L} - \epsilon)\right| \geq \left|\frac{1}{n}X_j^\top(Y - X\widehat{\beta}^{\lambda,L})\right| - \left|\frac{1}{n}X_j^\top \epsilon\right| \\ &\geq \lambda - \lambda/2 = \lambda/2 \quad \text{if} \quad \widehat{\beta}_j^{\lambda,L} \neq 0 \end{aligned}$$

$$\left| \frac{1}{n} X_j^\top (X\beta^* - X\widehat{\beta^{\lambda,L}}) \right| \geq \lambda/2 \quad \text{if} \quad \widehat{\beta}_j^{\lambda,L} \neq 0 \tag{14}$$

Denote by $\widehat{J} = J(\widehat{\beta^{\lambda,L}})$ and $M(\widehat{\beta^{\lambda,L}})$ the cardinal of $\widehat{J}$, then

$$\sum_{j=1}^{p} \left( \frac{1}{n} X_j^\top (X\beta^* - X\widehat{\beta^{\lambda,L}}) \right)^2 \geq \sum_{j \in \widehat{J}} \left( \frac{1}{n} X_j^\top (X\beta^* - X\widehat{\beta^{\lambda,L}}) \right)^2 \geq \sum_{j \in \widehat{J}} (\lambda/2)^2 \text{ by eq.(14)}$$

$$\underset{\text{minoré}}{} \quad = \quad M(\widehat{\beta^{\lambda,L}})\lambda^2/4 \tag{15}$$

First note that $\frac{\tilde{X}\tilde{X}^\top}{n}$ and $\frac{XX^\top}{n}$ have same maximal eigenvalue $\phi_{\max}$. Then,

$$\sum_{j=1}^{p} \left( \frac{1}{n} X_j^\top (X\beta^* - X\widehat{\beta^{\lambda,L}}) \right)^2 \quad = \quad \frac{1}{n} \left( X\beta^* - X\widehat{\beta^{\lambda,L}} \right)^\top \frac{XX^\top}{n} \left( X\beta^* - X\widehat{\beta^{\lambda,L}} \right)$$

$$\underset{\text{majoré}}{} \quad \leq \quad \frac{1}{n} \left( X\beta^* - X\widehat{\beta^{\lambda,L}} \right)^\top \phi_{\max} \left( X\beta^* - X\widehat{\beta^{\lambda,L}} \right)$$

$$= \quad \phi_{\max} \frac{1}{n} \| X\beta^* - X\widehat{\beta^{\lambda,L}} \|_2^2 \tag{16}$$

Combining (15) and (16), it comes

$$M(\widehat{\beta}^{\lambda,L}) \quad \leq \quad \frac{4\phi_{\max}}{\lambda^2}\|X\beta^* - X\widehat{\beta}^{\lambda,L}\|_n^2 \tag{17}$$

Using prediction bound (9),

$$M(\widehat{\beta}^{\lambda,L}) \quad \leq \quad \frac{4\phi_{\max}}{\lambda^2}\frac{16s\lambda^2}{\kappa^2} = \frac{64\phi_{\max}s}{\kappa^2}. \tag{18}$$

# Lemma~2

Consider model (1). Under $\mathcal{H}$ Assuptions, The estimator $\widehat{\beta}^{\lambda,L}$ for $\lambda = A\sigma\sqrt{\frac{\log(p)}{n}}$, with $A > 2\sqrt{2}$ is such that with probability greater than $1 - p^{1-\frac{A^2}{8}}$:

$$\|X(\widehat{\beta}^{\lambda,L} - \beta^*)\|_n^2 + \lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_1 \leq 4\lambda\|\widehat{\beta}^{\lambda,L} - \beta^*\|_{1,J_0}, \quad (19)$$

where $J_0 = J(\beta^*)$.

## Proof of Lemme~2

By definition of $\widehat{\beta^{\lambda,L}}$, we have <u>for all</u> $\beta \in \mathbb{R}^p$

$$\|Y - X\widehat{\beta^{\lambda,L}}\|_n^2 + 2\lambda\|\widehat{\beta^{\lambda,L}}\|_1 \leq \|Y - X\beta\|_n^2 + 2\lambda\|\beta\|_1$$

As $Y = X\beta^* + \epsilon$, it comes

$$\|X\beta^* - X\widehat{\beta^{\lambda,L}} + \epsilon\|_n^2 + 2\lambda\|\widehat{\beta^{\lambda,L}}\|_1 \leq \|X\beta^* - X\beta + \epsilon\|_n^2 + 2\lambda\|\beta\|_1$$

$\|\epsilon\|_n^2 +$ $\|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 + \frac{2}{n}\langle X\beta^* - X\widehat{\beta^{\lambda,L}}, \epsilon\rangle + 2\lambda\|\widehat{\beta^{\lambda,L}}\|_1 \leq \|X\beta^* - X\beta\|_n^2 + \frac{2}{n}\langle X\beta^* - X\beta, \epsilon\rangle + 2\lambda\|\beta\|_1$ $+ \|\epsilon\|_n^2$

<u>In particular</u> for $\beta = \beta^*$ and $V_j := \frac{1}{n}\sum_{i=1}^n \epsilon_i x_{i,j}$, we have

$$\|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 + \frac{2}{n}\langle X\beta^* - X\widehat{\beta^{\lambda,L}}, \epsilon\rangle + 2\lambda\|\widehat{\beta^{\lambda,L}}\|_1 \leq 2\lambda\|\beta^*\|_1$$

$$\Leftrightarrow \quad \|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 \leq 2\lambda\left(\|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1\right) + \frac{2}{n}\langle X\widehat{\beta^{\lambda,L}} - X\beta^*, \epsilon\rangle$$

$$\Leftrightarrow \quad \|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 \leq 2\lambda\left(\|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1\right) + \frac{2}{n}\sum_{j=1}^p\sum_{i=1}^n \epsilon_i x_{i,j}(\widehat{\beta_j^{\lambda,L}} - \beta_j^*)$$

$$\Leftrightarrow \quad \|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 \leq 2\lambda\left(\|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1\right) + 2\sum_{j=1}^p V_j(\widehat{\beta_j^{\lambda,L}} - \beta_j^*)$$

# Proof of Lemme~2

Then w.h.p $> 1 - p^{1 - \frac{A^2}{8}}$, on $\Omega = \bigcap_{j=1}^{p} \left\{ 2|V_j| \leq \lambda \right\}$

$$\|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 \leq 2\lambda \left( \|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1 \right) + 2 \sum_{j=1}^{p} V_j(\widehat{\beta_j^{\lambda,L}} - \beta_j^*)$$

$$\Leftrightarrow \quad \|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 \leq 2\lambda \left( \|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1 \right) + \lambda \sum_{j=1}^{p} (\widehat{\beta_j^{\lambda,L}} - \beta_j^*)$$

$$\Leftrightarrow \quad \|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 \leq 2\lambda \left( \|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1 \right) + \lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1$$

$$\Leftrightarrow \quad \|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 + \lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1 \leq 2\lambda \left( \|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1 \right) + \lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1 + \lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1$$

$$\Leftrightarrow \quad \|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 + \lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1 \leq 2\lambda(\|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1 + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_1) \tag{20}$$

Morover $p = J_0 \cup J_0^c$, where $J_0 = J(\beta^*) = \{j : \beta_j^* \neq 0\}$, we have

$$\begin{cases} \|\beta^*\|_1 = \|\beta^*\|_{1,J_0} + \|\beta^*\|_{1,J_0^c} = \|\beta^*\|_{1,J_0} \\ \|\widehat{\beta^{\lambda,L}}\|_1 = \|\widehat{\beta^{\lambda,L}}\|_{1,J_0} + \|\widehat{\beta^{\lambda,L}}\|_{1,J_0^c} \\ \|\widehat{\beta^{\lambda,L}} - \beta^*\|_1 = \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0^c} \end{cases}$$

# Proof of Lemme~2

As $\|a\| - \|b\| \leq \|a - b\|$, we have

$$\|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1 + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_1$$
$$= \|\beta^*\|_{1,J_0} - \|\widehat{\beta^{\lambda,L}}\|_{1,J_0} - \|\widehat{\beta^{\lambda,L}}\|_{1,J_0^c} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0^c}$$
$$= [\|\beta^*\|_{1,J_0} - \|\widehat{\beta^{\lambda,L}}\|_{1,J_0}] - \|\widehat{\beta^{\lambda,L}}\|_{1,J_0^c} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0^c}$$
$$\leq \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0} - \|\widehat{\beta^{\lambda,L}}\|_{1,J_0^c} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0^c}$$
$$\leq 2\|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0} + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0^c} - \|\widehat{\beta^{\lambda,L}}\|_{1,J_0^c}$$
$$\leq 2\|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0} + \|{\color{red}\rule{2cm}{0.3cm}} \beta^*\|_{1,J_0^c}$$
$$\leq 2\|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0}$$

Then w.h.p. $> 1 - p^{1 - \frac{A^2}{8}}$, on $\Omega$ and eq. (20):

$$\|X\beta^* - X\widehat{\beta^{\lambda,L}}\|_n^2 + \lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_1 \quad \leq 2\lambda(\|\beta^*\|_1 - \|\widehat{\beta^{\lambda,L}}\|_1 + \|\widehat{\beta^{\lambda,L}} - \beta^*\|_1)$$
$$\leq 4\lambda\|\widehat{\beta^{\lambda,L}} - \beta^*\|_{1,J_0}$$