

# Statistics of LASSO estimator

Now we look at a setting in which  $d$  is large and

$$Y_i = \langle \vartheta^*, X_i \rangle + \varepsilon_i$$

where  $\vartheta^*$  is a sparse vector, meaning that the number of non zero components is "small" compared to  $d$ . In this situation we define the LASSO estimator

$$\tilde{\vartheta}_{\lambda} \in \arg \min_{\vartheta \in \mathbb{R}^d} \frac{1}{2n} \|X\vartheta - Y\|^2 + \lambda \|\vartheta\|_1.$$

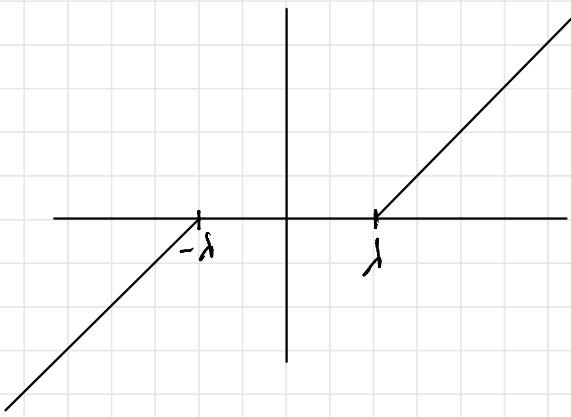
The optimality condition yields

$$0 \in \frac{1}{n} X^T (X\vartheta - Y) + \lambda \partial \|\cdot\|_1(\vartheta)$$

$$\Leftrightarrow \vartheta - \frac{1}{n} X^T (X\vartheta - Y) \in \vartheta + \lambda \partial \|\cdot\|_1(\vartheta) \\ = (I + \lambda \partial \|\cdot\|_1)(\vartheta)$$

$$\Leftrightarrow \vartheta = \text{prox}_{\lambda \|\cdot\|_1} \left( \vartheta - \frac{1}{n} X^T (X\vartheta - Y) \right)$$

$$\Leftrightarrow \vartheta = \text{soft}_{\lambda} \left( \vartheta - \frac{1}{n} X^T (X\vartheta - Y) \right)$$



when the component  $i$   
of  $\hat{\beta} - \frac{1}{n} X^T (X\hat{\beta} - y)$   
has absolute value  
less than  $\lambda$ ,  $\hat{\beta}_i = 0$ .

So, the  $l_1$  regularization term promotes sparsity.

Let's start with the statistical analysis.

$$\tilde{\theta}_{L_1} \in \operatorname{argmin}_{\theta \in \mathbb{R}^d} \frac{1}{2n} \|X\theta - Y\|^2 + \lambda \|\theta\|_1.$$

$$Y = X\theta^* + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n).$$

LEMMA 1

$$\|X\tilde{\theta}_{L_1} - X\theta^*\|^2 \leq 2n\lambda (\|\theta^*\|_1 - \|\tilde{\theta}_{L_1}\|_1) + 2\|\theta^* - \tilde{\theta}_{L_1}\|_1 \|X^T \varepsilon\|_\infty$$

Proof:

$$1. \quad \|X\tilde{\theta}_{L_1} - Y\|^2 + 2n\lambda \|\tilde{\theta}_{L_1}\|_1 \leq \|X\theta^* - Y\|^2 + 2n\lambda \|\theta^*\|_1$$

$$2. \quad \|X\tilde{\theta}_{L_1} - Y\|^2 = \|X\tilde{\theta}_{L_1} - X\theta^*\|^2 + \|X\theta^* - Y\|^2 + 2\langle X(\tilde{\theta}_{L_1} - \theta^*), -\varepsilon \rangle$$

$$\begin{aligned} \|X\tilde{\theta}_{L_1} - X\theta^*\|^2 &\leq \cancel{\|\varepsilon\|^2} + 2n\lambda (\|\theta^*\|_1 - \|\tilde{\theta}_{L_1}\|_1) - \cancel{\|\varepsilon\|^2} + 2\langle \tilde{\theta}_{L_1} - \theta^*, X^T \varepsilon \rangle \\ &\leq 2n\lambda (\|\theta^*\|_1 - \|\tilde{\theta}_{L_1}\|_1) + 2\|\tilde{\theta}_{L_1} - \theta^*\|_1 \|X^T \varepsilon\|_\infty. \end{aligned}$$

Now from LEMMA 1 we derive

$$\begin{aligned} \|X\tilde{\theta}_{L_1} - X\theta^*\|^2 &\leq 2n\lambda (\|\theta^*\|_1 - \|\tilde{\theta}_{L_1}\|_1) + 2(\|\theta^*\|_1 + \|\tilde{\theta}_{L_1}\|_1) \|X^T \varepsilon\|_\infty \\ &= 2\|\theta^*\|_1 (\|X^T \varepsilon\|_\infty + n\lambda) + 2\|\tilde{\theta}_{L_1}\|_1 (\|X^T \varepsilon\|_\infty - n\lambda) \end{aligned}$$

The purpose now is to make

$$\|X^T \varepsilon\|_\infty \leq n\lambda \quad \text{in high probability}$$

so to remove the term  $2\|\hat{\sigma}\|_1 (\|X^T \varepsilon\|_\infty - n\lambda)$  in the bound above.

Let's analyze the random vector  $X^T \varepsilon$ .

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_n^T \end{bmatrix} \in \mathbb{R}^{n \times d} \quad X^T = [x_1 \dots x_n] \in \mathbb{R}^{d \times n}$$

$$X = (X_{ij})_{\substack{1 \leq i \leq n \\ 1 \leq j \leq d}}$$

$X_{\cdot,j}$  are the columns of  $X$   
and  $X_{\cdot,i}^T$  are the rows of  $X^T$

$$\|X^T \varepsilon\|_\infty = \max_{1 \leq j \leq d} |\langle X_{\cdot,j}, \varepsilon \rangle|$$

$$\langle X_{\cdot,j}, \varepsilon \rangle = \sum_{i=1}^n X_{ij} \varepsilon_i \quad \text{it is a normal r.v.}$$

with zero mean and variance

$$\text{Var}(\langle X_{\cdot,j}, \varepsilon \rangle) = \sum_{i=1}^n X_{ij}^2 \text{Var}(\varepsilon_i) = \|X_{\cdot,j}\|^2 \sigma^2$$

$$\underline{\langle X_{\cdot,j}, \varepsilon \rangle \sim N(0, \|X_{\cdot,j}\|^2 \sigma^2)}$$

Reminder

If  $Z \sim N(0, \sigma^2)$

$$P(|Z| > t) \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right)$$

We look at  $IP(\{\|X^T \varepsilon\|_\infty > dn\})$

$$\|X^T \varepsilon\|_\infty = \max_{1 \leq i \leq d} |\langle X_{\cdot, i}, \varepsilon \rangle|$$

$$\{\|X^T \varepsilon\|_\infty > dn\} = \bigcup_{d=1}^d \{|\langle X_{\cdot, i}, \varepsilon \rangle| > dn\}$$

$$\begin{aligned} IP(\{\|X^T \varepsilon\|_\infty > dn\}) &\leq \sum_{d=1}^d IP(|\langle X_{\cdot, i}, \varepsilon \rangle| > nd) \\ &\leq 2 \sum_{d=1}^d \exp\left(-\frac{(nd)^2}{2 \|X_{\cdot, i}\|^2 \sigma^2}\right) \end{aligned}$$

Hyp:  $\|X_{\cdot, i}\| \leq M \sqrt{n}$  ( $X_{\cdot, i} \in \mathbb{R}^n$ )

$$\leq 2d \exp\left(-\frac{nd^2}{2M^2 \sigma^2}\right)$$

$$2d \exp\left(-\frac{nd^2}{2\beta^2\sigma^2}\right) \geq \delta \Leftrightarrow \frac{nd^2}{2M^2\sigma^2} = \log \frac{2d}{\delta}$$

$$\Leftrightarrow d = M\sigma \sqrt{\frac{2 \log(2d/\delta)}{n}}$$

With this choice we have

$$\mathbb{P}(\|X^T \varepsilon\|_\infty > dn) \leq \delta$$

and hence

$$\mathbb{P}(\|X^T \varepsilon\|_\infty \leq dn) \leq 1 - \mathbb{P}(\|X^T \varepsilon\|_\infty > dn) \leq 1 - \delta.$$

3. Therefore with probability  $\geq 1 - \delta$  we have

$$\|X^T \varepsilon\|_\infty \leq dn$$

and hence

$$\|X \hat{\vartheta}_{\ell_1} - X \vartheta^*\|^2 \leq \|\vartheta^*\|_1 4n\lambda.$$

In the end we have

$$\text{MSE}[X \hat{\vartheta}_{\ell_1}] \leq \|\vartheta^*\|_1 4M\sigma \sqrt{\frac{2 \log(2d/\delta)}{n}}$$

with probability greater than  $1 - \delta$ .

In the regime of  $d \geq n$ , say  $n = d \cdot q$  with  $q < 1$   
 we have

$$\text{MSE}[X \hat{\vartheta}_{\ell_2}] \leq \|\vartheta^*\|_1 4M\sigma \sqrt{\frac{2 \log 2d/\delta}{q \cdot d}}$$

Therefore if  $\|\vartheta^*\|_1 = O(d^{\alpha - \log^{\gamma} d})$  with  $\alpha < 1/2$   
 $\in \mathbb{R}^d$  and  $\gamma > 0$

we have that

$$\text{MSE}[X \hat{\vartheta}_{\ell_2}] \rightarrow 0 \text{ as } d \rightarrow \infty \text{ in high probability.}$$

For instance if  $\|\vartheta^*\|_{\infty} \leq B$ , then  $\|\vartheta^*\|_1 \leq B \underbrace{\|\vartheta^*\|_0}_{k = \#\text{supp } \vartheta^*}$

and hence we have

$$\text{MSE}[X \hat{\vartheta}_{\ell_2}] \leq 4MBk\sigma \sqrt{\frac{2 \log(2d/\delta)}{n}}$$

## Fast rates for Lasso estimator

Def  $X$  has incoherence  $k \in \mathbb{N}^*$  iff  $\left\| \frac{X^T X}{n} - I_d \right\|_{\infty} \leq \frac{1}{32k}$

Def  $\vartheta \in \mathbb{R}^d$ ,  $S \subset \{1, \dots, d\}$

$$\vartheta_S \in \mathbb{R}^d \quad (\vartheta_S)_j = \begin{cases} 0 & \text{if } j \notin S \\ \vartheta_j & \text{if } j \in S. \end{cases}$$

Lemma 2 Suppose that  $k \leq d$ ,  $S \subset \{1, \dots, d\}$   $|S| \leq k$ .

and that  $\vartheta \in \mathbb{R}^d$  is s.t.  $\|\vartheta_{S^c}\|_1 \leq 3 \|\vartheta_S\|_1$

Then  $\|\vartheta\|_2^2 \leq \frac{2}{n} \|X\vartheta\|_2^2$

Proof:  $\vartheta = \vartheta_S + \vartheta_{S^c}$

$$\frac{\|X\vartheta\|_2^2}{n} = \frac{\|X\vartheta_S\|_2^2}{n} + \frac{\|X\vartheta_{S^c}\|_2^2}{n} + \frac{2\langle X\vartheta_S, X\vartheta_{S^c} \rangle}{n}$$

$$\frac{\|X\vartheta_S\|_2^2}{n} = \left\langle \frac{X^T X}{n} \vartheta_S, \vartheta_S \right\rangle = \left\langle \left( \frac{X^T X}{n} - I_d \right) \vartheta_S, \vartheta_S \right\rangle + \|\vartheta_S\|_2^2$$

$$\left| \left\langle \left( \frac{X^T X}{n} - I_d \right) \vartheta_S, \vartheta_S \right\rangle \right| \leq \left\| \frac{X^T X}{n} - I_d \right\|_{\infty} \|\vartheta_S\|_2^2$$

Therefore  $\frac{\|X\vartheta_S\|_2^2}{n} \geq \|\vartheta_S\|_2^2 - \frac{\|\vartheta_S\|_1^2}{32k}$



Similarly

$$2. \frac{\|X\partial_{sc}\|^2}{n} \geq \|\partial_{sc}\|^2 - \frac{\|\partial_{sc}\|_1^2}{32k} \geq \|\partial_{sc}\|^2 - \frac{9\|\partial_s\|_1^2}{32}$$

$$\begin{aligned} 3. 2 \left| \left\langle \frac{X^T X}{n} \partial_s, \partial_{sc} \right\rangle \right| &= 2 \left\langle \left( \frac{X^T X}{n} - I_d \right) \partial_s, \partial_{sc} \right\rangle + 2 \underbrace{\langle \partial_s, \partial_{sc} \rangle}_0 \\ &\leq \frac{2}{32k} \|\partial_s\|_1 \|\partial_{sc}\|_1 \\ &\leq \frac{6}{32k} \|\partial_s\|_1^2 \end{aligned}$$

In the end we have

$$\begin{aligned} \frac{\|X\partial\|^2}{n} &\geq \|\partial_s\|^2 + \|\partial_{sc}\|^2 - \frac{\|\partial_s\|_1^2}{32k} - \frac{9\|\partial_s\|_1^2}{32k} - \frac{6\|\partial_s\|_1^2}{32k} \\ &= \|\partial_s\|^2 + \|\partial_{sc}\|^2 - \frac{16}{32k} \|\partial_s\|_1^2 \end{aligned}$$

$$\|\partial_s\|_1 \leq \sqrt{15} \|\partial_s\|_2 \leq \sqrt{k} \|\partial_s\|_2$$

$$\begin{aligned} \frac{\|X\partial\|^2}{n} &\geq \|\partial_s\|^2 + \|\partial_{sc}\|^2 - \frac{16k}{32k} \|\partial_s\|_2^2 \\ &= \|\partial_s\|^2 + \|\partial_{sc}\|^2 - \frac{\|\partial_s\|_2^2}{2} \\ &\geq \frac{\|\partial_s\|^2 + \|\partial_{sc}\|^2}{2} = \frac{\|\partial\|^2}{2}. \quad \square \end{aligned}$$

## Theorem

Suppose that  $X$  has incoherence  $\kappa \in \mathbb{N}^*$  and that

$$\|\text{supp } \theta^*\| \leq \kappa. \text{ Set } d = 4\sigma \sqrt{\frac{\log(2d/\delta)}{n}}$$

Then with probability  $\geq 1 - \delta$  we have

$$\text{MSE} [X \hat{\theta}_{\ell_1}] \leq 9 \cdot 2^5 \kappa \sigma^2 \frac{\log(2d/\delta)}{n}$$

Proof: It follows from LEMMA 1 that

$$\|X \hat{\theta}_{\ell_1} - X \theta^*\|^2 \leq 2nd (\|\theta^*\|_1 - \|\hat{\theta}_{\ell_1}\|_1) + 2 \|\hat{\theta}_{\ell_1} - \theta^*\|_1 \|X^T \varepsilon\|_\infty$$

Now we make  $\|X^T \varepsilon\|_\infty \leq \frac{nd}{2}$  with high probability

$$\begin{aligned} \text{IP}(\|X^T \varepsilon\|_\infty \geq t) &\leq \sum_{\ell=1}^d \text{IP}(|\langle X_{\cdot, \ell}, \varepsilon \rangle| \geq t) \\ &\leq 2 \sum_{\ell=1}^d \exp\left(-\frac{t^2}{2 \|X_{\cdot, \ell}\|^2 \sigma^2}\right) \end{aligned}$$

From the incoherence we have

$$\left(\frac{X^T X}{n} - \text{Id}\right)_{ij} = \frac{\langle X_{\cdot, i}, X_{\cdot, j} \rangle}{n} - \delta_{ij}$$

and hence  $\left| \frac{\|X_{\cdot i}\|^2}{n} - 1 \right| \leq \frac{1}{32k}$

which implies that  $\frac{\|X_{\cdot i}\|^2}{n} \leq 1 + \frac{1}{32k} \leq 2$ .

and hence  $\|X_{\cdot i}\|^2 \leq 2n$  ( $\|X_{\cdot i}\| \leq \sqrt{2n}$ )

Therefore  $\mathbb{P}(\|X^T \varepsilon\|_\infty > t) \leq 2d \exp\left(-\frac{t^2}{4n\sigma^2}\right)$

if we pick  $t = \frac{dn}{2}$  we have

$$\mathbb{P}\left(\|X^T \varepsilon\|_\infty > \frac{dn}{2}\right) \leq 2d \exp\left(-\frac{d^2 n}{16\sigma^2}\right)$$

$$2d \exp\left(-\frac{d^2 n}{16\sigma^2}\right) = \delta \Leftrightarrow \frac{d^2 n}{16\sigma^2} = \log \frac{2d}{\delta}$$

$$\Leftrightarrow d = 4\sigma \sqrt{\frac{\log(2d/\delta)}{n}}$$

Thus, with that choice we have

$$\mathbb{P}\left(\|X^T \varepsilon\|_\infty > \frac{dn}{2}\right) \leq \delta \text{ and hence}$$

$$\mathbb{P} \left( \|X^T \varepsilon\|_\infty \leq \frac{d\eta}{2} \right) \geq 1 - \delta.$$

Therefore with  $p. \geq 1 - \delta$  we have

$$\|X \hat{\vartheta}_{\mathcal{L}_1} - X \vartheta^*\|^2 \leq 2n\lambda (\| \vartheta^* \|_1 - \| \hat{\vartheta}_{\mathcal{L}_1} \|_1) + n\lambda \| \hat{\vartheta}_{\mathcal{L}_1} - \vartheta^* \|_1$$

Now set  $S = \text{supp } \vartheta^*$ . Then, we have

$$(\hat{\vartheta}_{\mathcal{L}_1} - \vartheta^*)_{S^c} = \hat{\vartheta}_{\mathcal{L}_1, S^c}, \text{ therefore}$$

$$\begin{aligned} \|X \hat{\vartheta}_{\mathcal{L}_1} - X \vartheta^*\|^2 &\leq -n\lambda \| \hat{\vartheta}_{\mathcal{L}_1, S^c} \|_1 + 2n\lambda (\| \vartheta^* \|_1 - \| \hat{\vartheta}_{\mathcal{L}_1, S} \|_1) \\ &\quad + n\lambda \| \hat{\vartheta}_{\mathcal{L}_1, S} - \vartheta^* \|_1 \\ &\leq -n\lambda \| \hat{\vartheta}_{\mathcal{L}_1, S^c} \|_1 + 3n\lambda \| \hat{\vartheta}_{\mathcal{L}_1, S} - \vartheta^* \|_1. \end{aligned}$$

If we define  $\vartheta = \hat{\vartheta}_{\mathcal{L}_1} - \vartheta^*$ , we have.

$$\|X \hat{\vartheta}_{\mathcal{L}_1} - X \vartheta^*\|^2 \leq -n\lambda \| \vartheta_{S^c} \|_1 + 3n\lambda \| \vartheta_S \|_1.$$

and in particular :  $\| \vartheta_{S^c} \|_1 \leq 3 \| \vartheta_S \|_1$

Therefore, by Lemma 2,

$$\| \vartheta_S \|_1 \leq \sqrt{|S|} \| \vartheta_S \|_2 \leq \sqrt{K} \| \vartheta \|_2 \leq \sqrt{\frac{2K}{n}} \| X \vartheta \|_2.$$

In the end, we have with  $p \geq 1 - \delta$

$$\|X \tilde{\theta}_{\mathcal{L}_1} - X \theta^*\|_2^2 \leq 3nd \|\tilde{\theta}_S\|_1 \leq 3n\lambda \sqrt{\frac{2k}{n}} \|\cancel{X(\tilde{\theta}_{\mathcal{L}_1} - \theta^*)}\|$$

and hence

$$\|X \tilde{\theta}_{\mathcal{L}_1} - X \theta^*\|^2 \leq 18nd^2k$$

that is

$$\frac{1}{n} \|X \tilde{\theta}_{\mathcal{L}_1} - X \theta^*\|^2 \leq 18 \cdot k \cdot 16 \sigma^2 \frac{\log(2d/\delta)}{n} \quad \square$$