

# High-dimensional Statistics

6<sup>th</sup> Lecture

A. DALALYAN

MASTER MVA

## 0. RECALL

We observe  $(Y_i)_{i=1, \dots, n}$  such that

$$(1) \quad y = f^* + \xi \quad ; \quad \xi \sim \mathcal{N}(0, \sigma^2 I_n)$$

where  $f^* = (f_1^*, \dots, f_n^*) \in \mathbb{R}^n$  is an unknown vector (called signal)

$\xi = (\xi_1, \dots, \xi_n) \in \mathbb{R}^n$  is a random noise

Dictionary:  $f^* = \Phi \cdot \beta^*$  where  $\Phi \in \mathbb{R}^{n \times p}$  is a known design matrix

Sparsity: the dimension  $p$  of  $\beta^*$  is large, possibly larger than  $n$ , but many coordinates of  $\beta^*$  vanish:

$$s = \|\beta^*\|_0 = \sum_{j=1}^p \mathbb{1}(|\beta_j^*| \neq 0) \ll p.$$

Lasso: 
$$\hat{\beta}^{\text{Lasso}} \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \|y - \Phi \beta\|_2^2 + C \|\beta\|_1 \right\}$$

THEOREM: If  $\Phi$  is orthogonal, i.e.,  $\frac{1}{n} \Phi^T \Phi = I_p$ , then  $\hat{\beta}^{\text{Lasso}}$  coincides with the soft thresholding procedure.

More precisely, for  $\lambda = \frac{C}{2n}$ , we have

$$\hat{\beta}_j^{\text{Lasso}} = \left( \left| \frac{1}{n} (\Phi^T y)_j \right| - \lambda \right)_+ \cdot \text{sign}((\Phi^T y)_j), \quad j=1, \dots, p.$$

## 1. Remarks

1) From now on, we will parameterize the Lasso by  $\lambda$  instead of  $C$ :

$$(2) \quad \hat{\beta}_\lambda^{\text{Lasso}} \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|y - \Phi \beta\|_2^2 + \lambda \|\beta\|_1 \right\}.$$

2) The condition  $\frac{1}{n} \Phi^T \Phi = I_p$  essentially means that the dictionary

3) Even if  $\Phi$  is not ortho-normal (ON), the Lasso is efficiently computable even for very large values of  $p$ .

Lemma:  $\hat{\beta}$  is a solution to (2) if and only if there exist  $(\hat{u}, \hat{v}, \hat{t}) \in \mathbb{R}^p \times \mathbb{R}^n \times \mathbb{R}$  such that  $(\hat{\beta}, \hat{u}, \hat{v}, \hat{t})$  is a solution of the second-order cone program (SOCP):

$$\left. \begin{aligned} \min \left\{ \frac{1}{2n} \times t + \lambda \sum_{j=1}^p u_j \right\} \\ \text{subject to : } & -u_j \leq \beta_j \leq u_j \quad \forall j=1, \dots, p \\ & y - \Phi\beta = v \\ & \|(v; t; \frac{1}{2})\|_2 \leq t + \frac{1}{2} \end{aligned} \right\} \text{(SOCP1)}$$

Proof.

1) We prove first that if  $(\hat{\beta}, \hat{u}, \hat{v}, \hat{t})$  is a solution of (SOCP-1) then  $\hat{\beta}$  is a solution of (2). Indeed, let  $\beta$  be any vector of  $\mathbb{R}^p$ . Define  $u_j = |\beta_j|$ ,  $j=1, \dots, p$ ;  $v = y - \Phi\beta$  and  $t = \|y - \Phi\beta\|_2^2$ . One easily checks that all the constraints of (SOCP1) are fulfilled for  $(\beta, u, v, t)$  defined in such a way. Therefore,  $(\beta, u, v, t)$  is a feasible solution and

$$(3) \quad \frac{1}{2n} \times \hat{t} + \lambda \sum_{j=1}^p \hat{u}_j \leq \frac{1}{2n} \times t + \lambda \sum_{j=1}^p u_j = \frac{1}{2n} \|y - \Phi\beta\|_2^2 + \lambda \|\beta\|_1.$$

On the other hand, since  $(\hat{\beta}, \hat{u}, \hat{v}, \hat{t})$  is a solution of (SOCP1) we have [notice that the third constraint is equivalent to  $\|v\|_2^2 \leq t$ ]:  $\|y - \Phi\hat{\beta}\|_2^2 \leq \hat{t}^2$  and  $\|\hat{\beta}\|_1 \leq \|\hat{u}\|_1$ .

Combined with (3), this implies that

$$\frac{1}{2n} \|y - \Phi\hat{\beta}\|_2^2 + \lambda \|\hat{\beta}\|_1 \leq \frac{1}{2n} \|y - \Phi\beta\|_2^2 + \lambda \|\beta\|_1.$$

Since this holds true for every  $\beta \in \mathbb{R}^p$ ,  $\hat{\beta}$  is a solution of (2).



2) Let  $\hat{\beta}$  be a solution of (2) and define

$$\hat{u}_j = |\hat{\beta}_j|, j=1, \dots, p; \quad \hat{v} = y - X\hat{\beta} \quad \text{and} \quad \hat{t} = \|\hat{v}\|_2^2.$$

One easily checks that  $(\hat{\beta}, \hat{u}, \hat{v}, \hat{t})$  verifies the constraints of (SOCP1). Furthermore, if  $(\beta, u, v, t)$  is another vector satisfying the constraints of (SOCP1), we have

$$(4) \quad \frac{1}{2n} \|y - \Phi\beta\|_2^2 + \lambda \|\beta\|_1 \leq \frac{1}{2n} \|y - \Phi\hat{\beta}\|_2^2 + \lambda \|\hat{\beta}\|_1$$

The left-hand side of (4) is equal to

$$\frac{1}{2n} \hat{t} + \lambda \sum_{j=1}^p \hat{u}_j$$

whereas the right-hand side of (4) is bounded as follows:

$$\frac{1}{2n} \|y - \Phi\beta\|_2^2 + \lambda \|\beta\|_1 \leq \frac{1}{2n} \|v\|_2^2 + \lambda \sum_{j=1}^p u_j \leq \frac{1}{2n} t + \lambda \sum_{j=1}^p u_j.$$

This implies that

$$\frac{1}{2n} \hat{t} + \lambda \sum_{j=1}^p \hat{u}_j \leq \frac{1}{2n} t + \lambda \sum_{j=1}^p u_j$$

for every  $(\beta, u, v, t)$  satisfying the constraints. ■

### 3. Prediction loss of the Lasso: the "slow" rates.

$$y = \Phi\beta^* + \varepsilon \quad ; \quad \hat{\beta}^{\text{Lasso}} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{2n} \|y - \Phi\beta\|_2^2 + \lambda \|\beta\|_1 \right\}$$

We wish to evaluate the prediction loss

$$l_n(\hat{f}, f^*) = \frac{1}{n} \|\hat{f} - f^*\|_2^2 = \frac{1}{n} \|\Phi(\hat{\beta} - \beta^*)\|_2^2 \triangleq l_n(\hat{\beta}, \beta^*).$$

Note: we divide by  $n$  since  $\|\hat{f} - f^*\|_2^2$  is a sum containing  $n$  terms.

Condition C1: For  $\forall j \in \{1, \dots, p\}$ , we have  $\|\Phi^j\|_2^2 = \sum_{i=1}^n \varphi_{ij}^2 \leq n$ .

Lemma For every  $\delta \in (0, 1)$ , the event  $\mathcal{B} = \{ \|\Phi^T \xi\|_\infty \leq \sigma \sqrt{2n \ln(p/\delta)} \}$  satisfies  $\mathbb{P}(\mathcal{B}) \geq 1 - \delta$ .

Proof. We will prove that  $\mathbb{P}(\mathcal{B}^c) \leq \delta$ . Set  $x = \sqrt{2 \ln(p/\delta)}$  and  $\xi_j = (\Phi^j)^T \xi / (\sigma \cdot \|\Phi^j\|_2)$ ;  $j = 1, \dots, p$ . Since Gaussian distribution is stable by affine transformations, we have  $\xi_j \sim \mathcal{N}(0, 1) \forall j$ .

Therefore,

$$\begin{aligned} \mathbb{P}(\mathcal{B}^c) &= \mathbb{P}\left(\max_j |(\Phi^j)^T \xi| > \sigma \sqrt{2n \ln(p/\delta)}\right) \\ &= \mathbb{P}\left(\bigcup_{j=1}^p \{ |(\Phi^j)^T \xi| > \sigma \sqrt{n} \cdot x \}\right) \\ &\leq \sum_{j=1}^p \mathbb{P}\left(|(\Phi^j)^T \xi| > \sigma \sqrt{n} \cdot x\right) \\ &\leq \sum_{j=1}^p \mathbb{P}\left(|(\Phi^j)^T \xi| > \sigma \|\Phi^j\|_2 \cdot x\right) \\ &= \sum_{j=1}^p \mathbb{P}\left(|\xi_j| > x\right) = 2p \mathbb{P}\left(\xi_1 > x\right) \end{aligned}$$

where for the last equality we used the symmetry of the Gaussian distribution. To complete the proof we will use the Gaussian tail bound:  $\mathbb{P}(\xi_1 > x) \leq \frac{1}{2} \exp(-x^2/2)$ ,  $\forall x > 0$ .

This leads to

$$\mathbb{P}(\mathcal{B}^c) \leq 2p \times \frac{1}{p} \times \exp\left(-\frac{1}{2} (\sqrt{2 \ln(p/\delta)})^2\right) = \delta. \quad \blacksquare$$

For completeness, we provide the proof of the Gaussian tail bound.

$$\mathbb{P}(\xi_1 > x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-u^2/2} du.$$

Set  $G(x) = \mathbb{P}(\xi_1 > x) - \frac{1}{2} e^{-x^2/2}$ . One easily checks that

$$G'(x) = -\frac{1}{\sqrt{2\pi}} e^{-x^2/2} + \frac{x e^{-x^2/2}}{2} = \frac{e^{-x^2/2}}{2} \left(x - \sqrt{\frac{2}{\pi}}\right).$$

This implies that  $G$  is decreasing on  $]0, \sqrt{2/\pi}[$  and increasing on  $]\sqrt{2/\pi}, +\infty[$ . Therefore,  $G(x) \leq \max(G(0), G(+\infty)) = 0 \quad \blacksquare$



From now on, we will always work on the event  $\mathcal{B}$ .

Theorem: If  $\Phi$  satisfies condition C1, then on  $\mathcal{B}$ , it holds that

(5) 
$$l_n(\hat{\beta}^L, \beta^*) \leq \inf_{\beta \in \mathbb{R}^p} \left\{ l_n(\beta, \beta^*) + 4\lambda \|\beta\|_1 \right\}$$
 provided that  $\lambda$  is chosen  $\geq \sigma \sqrt{\frac{2}{n} \ln(p/\delta)}$ .

Proof: Let  $\beta \in \mathbb{R}^p$  be any vector. Then

$$\frac{1}{2n} \|y - \Phi \hat{\beta}\|_2^2 + \lambda \|\hat{\beta}\|_1 \leq \frac{1}{2n} \|y - \Phi \beta\|_2^2 + \lambda \|\beta\|_1.$$

$$\frac{1}{2n} \|\Phi(\beta^* - \hat{\beta}) + \xi\|_2^2 + \lambda \|\hat{\beta}\|_1 \leq \frac{1}{2n} \|\Phi(\beta^* - \beta) + \xi\|_2^2 + \lambda \|\beta\|_1$$

$$\frac{1}{n} \|\Phi(\beta^* - \hat{\beta})\|_2^2 + 2\lambda \|\hat{\beta}\|_1 \leq \frac{2}{n} \xi^T \Phi(\hat{\beta} - \beta) + 2\lambda \|\beta\|_1 + \frac{1}{n} \|\Phi(\beta^* - \beta)\|_2^2$$

$$l_n(\hat{\beta}, \beta^*) \leq l_n(\beta, \beta^*) + 4\lambda \|\beta\|_1 + \underbrace{\frac{2}{n} \xi^T \Phi(\hat{\beta} - \beta) - 2\lambda (\|\beta\|_1 + \|\hat{\beta}\|_1)}_{\leq 0 \text{ on } \mathcal{B}}$$

En effet

$$\frac{2}{n} \xi^T \Phi(\hat{\beta} - \beta) \leq \frac{2}{n} (\Phi^T \xi)^T (\hat{\beta} - \beta) \leq \frac{2}{n} \|\Phi^T \xi\|_\infty \cdot \|\hat{\beta} - \beta\|_1$$

$$\leq \frac{2}{n} \sigma \sqrt{2n \ln(p/\delta)} \times \|\hat{\beta} - \beta\|_1$$

$$\leq 2\lambda \|\hat{\beta} - \beta\|_1 \leq 2\lambda (\|\hat{\beta}\|_1 + \|\beta\|_1) \quad \square$$

### Commentaires

① On dit que (5) est une Oracle Inequality with "slow" rate.

If an oracle tells us what is the best sparse approximation  $\bar{\beta}$  of  $\beta^*$ , then can use it as estimator of  $\beta^*$  and get the loss  $l_n(\bar{\beta}, \beta^*)$ .

Since we do not know  $\bar{\beta}$ , we use  $\hat{\beta}^{\text{Lasso}}$  and get the loss bounded

$$\text{by } l_n(\bar{\beta}, \beta^*) + \underbrace{\text{Const} \times \sqrt{\frac{\ln(p)}{n}} \times \|\bar{\beta}\|_1}_{\text{this is called "slow" rate.}}$$