

Lecture 12

Linear regression

Suppose we have an unknown vector  $\beta^* \in \mathbb{R}^d$ .

We observe linear observation of  $\beta^*$  of the form:

$$Y_i = \langle x_i, \beta^* \rangle + \varepsilon_i \quad i=1, \dots, n$$

$\downarrow \quad \downarrow \quad \downarrow$   
known known in  $\mathbb{R}^d$  noise

Assume  $\varepsilon_i$ 's are zero-mean and in  $\text{SubG}(\sigma^2)$

$$Y = X \beta^* + \varepsilon$$

$\uparrow \mathbb{R}^{n \times d}$

$$\begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} \cdots & x_1 & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & x_n & \cdots \end{bmatrix} \begin{bmatrix} \beta^* \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Fix design, we assume  $X$  is fixed.

[Another interesting regime is when  $x_i$ 's are random.]

Goal :

find  $\hat{\beta}$  such that  $\hat{\beta}$  is close to  $\beta^*$

what does close mean?

- small distance to  $\hat{\beta}$  and  $\beta^*$

say  $\|\hat{\beta} - \beta^*\|_2^2$  is small

- small "de-noising objective"

$Y = X\hat{\beta}$  is similar to  $Y = X\beta^*$

$$\frac{1}{n} \sum_{i=1}^n (\langle \alpha_i, \hat{\beta} \rangle - \langle \alpha_i, \beta^* \rangle)^2$$

$$= \frac{1}{n} \|X\hat{\beta} - X\beta^*\|_2^2$$

clearly, we do not have  $\beta^*$ . Thus, it is difficult to measure the quality of  $\hat{\beta}$ . What we usually do is to pick a "proxy" quantity for these measures and find  $\hat{\beta}$  that minimize them.

While a great deal of effort is dedicated to finding solutions. It is always important to look back and see the solution we have found via optimizing the proxy is indeed a good solution for the original objective as well.

Solution

$$\hat{\beta} \in \arg \min_{\beta} \|X\beta - Y\|_2^2 \quad *$$

$$= \arg \min_{\beta} \frac{1}{n} \sum_{i=1}^n (\langle x_i, \beta \rangle - y_i)^2$$

the gradient of  $\|X\beta - Y\|_2^2 = 0$

$$\|X\beta - Y\|_2^2 = (\beta^T X^T - Y^T) \cdot (X\beta - Y)$$

$$= \beta^T X^T X \beta - 2 \beta^T X^T Y + Y^T Y$$

$$\Rightarrow X^T X \beta = X^T Y$$

$$\Rightarrow \beta = \underbrace{(X^T X)}_{\text{pseudo inverse}}^T X^T Y$$

pseudo inverse

[not the focus of this lecture]

Given Sub-Gaussianity assumption on  $\varepsilon$ ,  
 what can we say about the error  
 of  $\hat{\beta}$ ?

Can we exploit any structure in  $X$ ?

such as low rank  $X$ ? or sparsity of  $\beta^*$ ?

since  $\hat{\beta}^*$  is an arg min in ~~\*~~ we have:

$$\|\hat{X}\hat{\beta} - Y\|_2^2 \leq \|X\beta^* - Y\|_2^2 = \|\varepsilon\|_2^2 \quad \textcircled{1}$$

on the other hand:

$$\|\hat{X}\hat{\beta} - Y\|_2^2 = \|\hat{X}\hat{\beta} - X\beta^* + \varepsilon\|_2^2$$

$$= \|\hat{X}\hat{\beta} - X\beta^*\|_2^2 - 2 \langle \varepsilon, \hat{X}\hat{\beta} - X\beta^* \rangle + \|\varepsilon\|_2^2$$

\textcircled{2}

$$\textcircled{1}, \textcircled{2} \Rightarrow \frac{1}{n} \|\hat{X}\hat{\beta} - X\beta^*\|_2^2 \leq \frac{2}{n} \langle \varepsilon, \hat{X}\hat{\beta} - X\beta^* \rangle$$

(basic inequality)

$$\Rightarrow \|X\hat{\beta} - X\beta^*\|_2 \leq 2 < \varepsilon, \frac{X\hat{\beta} - X\beta^*}{\|X\hat{\beta} - X\beta^*\|_2} >$$

$$\leq 2 \sup_{\beta} < \varepsilon, \frac{X\beta - X\beta^*}{\|X\beta - X\beta^*\|_2} > *$$

↳ does not on  $\hat{\beta}$  any more.

Let  $U = [u_1, \dots, u_r]$  be a matrix with orthonormal columns a basis for column space of  $X$

(where  $r$  is the rank of  $X^T X$ )

$\frac{X\beta - X\beta^*}{\|X\beta - X\beta^*\|_2}$  is a vector in column space of  $X$

Hence, it can be written in the basis  $u_i$ 's

$$\exists a : \frac{X\beta - X\beta^*}{\|X\beta - X\beta^*\|_2} = \frac{Ua}{\|Ua\|}$$

$$\|\lambda \hat{\beta} - \lambda \beta^*\|_2 \leq 2 \sup_{\|a\| \leq 1} \langle \varepsilon, a \rangle$$

$$\leq 2 \sup_{\|a\| \leq 1} \langle U^T \varepsilon, a \rangle$$

$$= 2 \|U^T \varepsilon\|$$

$$\Rightarrow \frac{1}{n} \|\lambda \hat{\beta} - \lambda \beta^*\|_2^2 \leq 4 \|U^T \varepsilon\|_2^2$$

$$U \in \mathbb{R}^{n \times r}$$

$$U^T \varepsilon = \begin{bmatrix} \underline{u_1} \\ \vdots \\ \underline{u_r} \end{bmatrix} \begin{bmatrix} \varepsilon \end{bmatrix}$$

$$\text{Let } v = U^T \varepsilon \Rightarrow v_i = \langle u_i, \varepsilon \rangle$$

$$\Rightarrow v_i = \sum_{j=1}^n u_{ij} \cdot \varepsilon_j$$

$$\Rightarrow v_i \in \text{Sub } G \left( \sum_{j=1}^n u_{ij}^2 \sigma^2 \right) \in \text{Sub } G(\sigma^2)$$

$$\mathbb{E}_{\varepsilon} \left[ \frac{1}{n} \| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* \|_2^2 \right] \leq \frac{4}{n} \mathbb{E}_{\varepsilon} [\| \mathbf{U}^\top \varepsilon \|_2^2]$$

$$= \frac{4}{n} \mathbb{E} \left\{ \sum_{i=1}^r v_i^2 \right\} \leq \frac{4}{n} \sum_{i=1}^r \mathbb{E}[v_i^2]$$

$$\leq \frac{4}{n} \sum_{i=1}^r (\sqrt{2} \sigma)^2 \leq \Theta \left( \frac{r \sigma^2}{n} \right)$$

depends on  $r$  not  $d$

$$r \leq \min \{ n, d \}$$

## Sparsity

Consider the case where all but  $k$  coordinate of  $\beta^*$  is zero.

$$B_0^d(k) := \{x \in \mathbb{R}^d : \|x\|_0 \leq k\}$$

↳ unit ball

we also pick  $\hat{\beta}$  from  $B_0^d(k)$

$$\hat{\beta} = \arg \min_{\substack{\beta \\ \beta \in B_0^d(k)}} \|X\beta - Y\|_2^2$$

Now we focus on bounding

$$\|\hat{X}\beta - X\beta^*\|_2^2.$$

Earlier, we have shown

$$\|\hat{X}\beta - X\beta^*\|_2 \leq 2 < \varepsilon, \frac{\|\hat{X}\beta - X\beta^*\|}{\|\hat{X}\beta - X\beta^*\|_2} >$$

$$\leq 2 \sup_{\beta \in B_0^d(K)} < \varepsilon, \frac{\|X(\beta - \beta^*)\|}{\|X(\beta - \beta^*)\|_2} >$$

Now  $\beta - \beta^*$  is a vector in  $\mathbb{R}^d$

with at most  $2k$  non-zero entries

Let  $S$  denote the set of indices

that are not zero.

We know  $|S| \leq 2k$

we can continue our bound by:

$$2 \sup_{\beta \in \beta_0^d(k)} \left\| X(\beta - \beta^*) \right\|_2 < \varepsilon, \frac{X(\beta - \beta^*)}{\|X(\beta - \beta^*)\|_2}$$

$$\leq 2 \max_{S \subseteq [d]} \sup_{\alpha_S} \left\| X \alpha_S \right\|_2 < \varepsilon, \frac{\|X \alpha_S\|_2}{\|X \alpha_S\|}$$

$|S| \leq 2k$       ↓

where  $\alpha_i$  is zero  
for all  $i \notin S$

Note that  $X \alpha_S$  lies in the column space of  $X_S$  restricted to columns that are in  $S$ .

$$\begin{bmatrix} 1, 2, 3 \\ \vdots \\ \alpha = \beta - \beta^* \end{bmatrix} \quad X \quad \left[ \begin{array}{c} \text{[yellow shaded]} \\ \text{[yellow shaded]} \\ \text{[yellow shaded]} \end{array} \right]$$

for example  
 $S = \{1, 2, 3\}$

Let  $U_S = [u_1, \dots, u_r]$  be a matrix where its columns form an orthonormal basis for the column space of  $X_S$ .

$$\text{hence } r \leq 2k$$

with a very similar argument as before

$$\|X\hat{\beta} - X\beta^*\|_2 \leq 2 \max_{\substack{S \subseteq [d] \\ |S| \leq 2k}} \|U_S^T \varepsilon\|_2$$

$$\Rightarrow E_S \left[ \frac{1}{n} \|X\hat{\beta} - X\beta^*\|_2^2 \right] \leq 2 E_S \left[ \max_{\substack{S \subseteq [d] \\ |S| \leq 2k}} \|U_S^T \varepsilon\|_2^2 \right] \leq \Theta \left( \frac{\sigma^2 k \log(d)}{n} \right)$$

next lecture ↗