

# SDS 387 Linear Models

Fall 2024

Lecture 24 - Tue, Nov 21, 2024

Instructor: Prof. Ale Rinaldo

- Last time: we assume a linear model with Gaussian covariates

$$Y_i = \Phi_i^\top \beta^* + \varepsilon_i \quad i=1, \dots, n$$

$$\Phi_1, \dots, \Phi_n \stackrel{i.i.d.}{\sim} N_d(0, I_d) \quad \text{and}$$

$$\varepsilon_1, \dots, \varepsilon_n \mid \Phi_1, \dots, \Phi_n \stackrel{i.i.d.}{\sim} (0, \sigma^2)$$

We are interested in the exact expression for the expected excess risk of  $\hat{\beta}$  OLS. If  $n \geq d+2$  then this value is

$$\sigma^2 \frac{d}{n-d-1}$$

if  $n = d$  or  $d+1$  this value is  $\infty$ .

Next if  $d > n$  we will consider the min-norm

least squares estimator

$$\hat{\beta}_{MN} = \Phi^+ Y = \underbrace{\Phi^T (\Phi \Phi^T)^{-1}}_{n \times n \text{ of full rank}} Y$$

$\downarrow$   
 $\Phi$  matrix with  
 $n \times d$  with row  $\perp$   
 $\Phi^T$

Last time we started analyzing the expected excess risk of  $\hat{\beta}_{MN}$ :

$$\mathbb{E} [\|\beta^* - \hat{\beta}_{MN}\|^2] = T_1 + T_2$$

$\downarrow$  bias term  
 $\searrow$  usual variance term

$$= \mathbb{E} [\underbrace{\|(\mathbb{I} - \Pi)\beta^*\|^2}_{\text{orthogonal projection onto null}(\Phi)}] + \sigma^2 \mathbb{E} [\text{tr}(\Phi \Phi^T)^{-1}]$$

$$= \frac{n}{d} \|\beta^*\|^2 + \sigma^2 \mathbb{E} [\text{tr}(\Phi \Phi^T)^{-1}]$$

$(\Phi \Phi^T)^{-1}$  is an inverse Wishart with scale parameter

$\Gamma_n$  and  $d$  degrees of freedom. The diagonal entries are inverse of  $\chi_{d-n-1}^2$  and have

expectations equal to

$$\begin{cases} \frac{1}{d-n-1} & d \geq n+2 \\ \infty & d = n \text{ or } d = n+1 \end{cases}$$

So putting together all these terms:

$$\mathbb{E} \left[ \left\| \beta^* - \hat{\beta}_{\text{OLS}} \right\|^2 \right] = \begin{cases} \sigma^2 \frac{d}{n-d-1} & d \leq n-2 \\ \|\beta^*\|^2 \left[ \frac{d-1}{d} \right] + \sigma^2 \frac{1}{d-n-1} & d \geq n+2 \\ \infty & \text{otherwise} \end{cases}$$

See Belkin, Hsu, Xu (2020)

Letting  $\gamma = \lim_n \frac{d_n}{n}$

$$= \begin{cases} \sigma^2 \frac{\gamma}{1-\gamma} & \gamma < 1 \\ \|\beta^*\|^2 \left( 1 - \frac{1}{\gamma} \right) + \sigma^2 \frac{1}{\gamma-1} & \gamma > 1 \\ \infty & \gamma = 1 \end{cases}$$

- Ridge regression has smaller risk profiles when it is optimally tuned  $d = \frac{\sigma^2 \gamma}{\|\beta^*\|^2}$   
tuning parameter for ridge

Solution: use cross-validation. See references.

# INFERENCE

$$y_n = \Phi_n^\top \beta^* + \varepsilon_n$$

$$\Phi_n \stackrel{iid}{\sim} P_\Phi \quad \varepsilon_n | \Phi_n \sim (0, \sigma^2) \text{ indep}$$

In this task, we are interested in estimating  $\beta^*$ , assuming a random design and well specified model.

If the model is well specified (i.e.  $\mathbb{E}[y_n | \Phi_n] = \Phi_n^\top \beta^*$ ) this is a simple extension of <sup>and</sup>  $\text{Var}(y_n - \Phi_n^\top \beta^*) = \sigma^2$  the analysis we did in the fixed-design case.

- Assume throughout that  $\frac{\Phi_n^\top \Phi_n}{n} = \frac{\sum_1^n \Phi_n \Phi_n^\top}{n}$  is invertible with prob. 1. Then  $\hat{\beta} \xrightarrow{P} \beta^*$ .

To see this, write

$$\hat{\beta} = \beta^* + \frac{\sum_1^n \Phi_n^\top \varepsilon_n}{n}$$

and notice that

$$i) \quad \frac{\sum_1^n \Phi_n \Phi_n^\top}{n} \xrightarrow{P} \sum_1 = \mathbb{E}[\Phi_n \Phi_n^\top] \quad \text{by WLLN}$$

$$\text{so } \frac{\sum_1^n \Phi_n \Phi_n^\top}{n}^{-1} \xrightarrow{P} \sum_1^{-1} \quad \text{by CMT}$$

$$ii) \quad \frac{\sum_1^n \Phi_n^\top \varepsilon_n}{n} \xrightarrow{P} \mathbb{E}[\Phi_n^\top \varepsilon_n] =$$

$$\text{by WLLN} \quad \mathbb{E}_{\Phi_n} [\Phi_n \mathbb{E}[\varepsilon_n | \Phi_n]] = 0$$

(4)

$$n \cdot n) \quad \sum_{i=1}^n \frac{\Phi_i^T \varepsilon_i}{n} \xrightarrow{P} \sum_{i=1}^n \cdot 0 = 0$$

by Slutsky

Also

$$\sqrt{n} (\hat{\beta} - \beta^*) \xrightarrow{d} N_d(0, \sigma^2 \Sigma^{-1})$$

This is because

$$\begin{aligned} \sqrt{n} (\hat{\beta} - \beta^*) &= \sum_{i=1}^n \frac{\Phi_i^T \varepsilon_i}{\sqrt{n}} \\ &\downarrow \\ &\xrightarrow{P} \Sigma^{-1} \end{aligned}$$

$$\text{Next} \quad \frac{\Phi^T \Sigma}{\sqrt{n}} = \sqrt{n} \sum_{i=1}^n \frac{\Phi_i \cdot \varepsilon_i}{n}$$

↳ normalized average of  $n$  random vectors

We saw that  $\mathbb{E}[\Phi_i \cdot \varepsilon_i] = 0$ . Also

$$\begin{aligned} \text{Var}[\Phi_i \cdot \varepsilon_i] &= \mathbb{E}[\varepsilon_i^2 \Phi_i \Phi_i^T] \\ &= \mathbb{E}_{\Phi_i} \left[ \underbrace{\mathbb{E}[\varepsilon_i^2 | \Phi_i]}_{\sigma^2} \Phi_i \Phi_i^T \right] \\ &= \sigma^2 \mathbb{E}[\Phi_i \Phi_i^T] = \sigma^2 \Sigma \end{aligned}$$

(5)

So by CLT

$$\sqrt{n} \frac{\Phi^T \varepsilon}{n} \xrightarrow{d} N_d(0, \sigma^2 \Sigma_1)$$

Finally by Slutsky's theorem:

$$\begin{aligned} \sqrt{n}(\hat{\beta} - \beta^*) &= \hat{\Sigma}^{-1} \sqrt{n} \frac{\Phi^T \varepsilon}{n} \xrightarrow{d} \Sigma_1^{-1} N_d(0, \sigma^2 \Sigma_1) \\ &= N_d(0, \sigma^2 \Sigma_1^{-1}) \end{aligned}$$

## INFERENCE IN ASSUMPTION-FREE SETTINGS (OR DISTRIBUTION-FREE)

Assumption-free means

- 1)  $\mathbb{E}[Y | \Phi] \neq \Phi^T \beta^*$
- 2) random design

More generally, it means we are not willing to make assumptions beyond minimal ones that are needed for the model to be well-defined (e.g. moment assumptions)

So, assume  $(\Phi, Y) \sim P_{\Phi, Y} \in \mathbb{R}^d \times \mathbb{R}$

$\downarrow$   $\rightarrow$  response

covariates/  
features/regressors

Let's only assume that  $\Phi$  and  $Y$  have 2 <sup>full rank</sup>  $\uparrow$  moments

$\mathbb{E}[Y^2] < \infty$  and  $\mathbb{E}[\Phi \Phi^T] = \sum_{d \times d} \mathbb{1}$

Then, trivially:

$$Y = \mathbb{E}[Y | \Phi] + \underbrace{Y - \mathbb{E}[Y | \Phi]}_{\varepsilon}$$

↓  
regression  
function

$$x \in \mathbb{R}^d \mapsto \mathbb{E}[Y | \Phi = x]$$

$\varepsilon \rightarrow$  by construction  
 $\mathbb{E}[\varepsilon | \Phi] = 0$   
so  $\mathbb{E}[\varepsilon] = 0$

• In general  $\text{Var}[\varepsilon | \Phi]$  may depend on  $\Phi$

• We saw that

$$\beta^* = \underset{\beta \in \mathbb{R}^d}{\text{argmin}} \mathbb{E}[(Y - \Phi^T \beta)^2] = \underset{\beta \in \mathbb{R}^d}{\text{argmin}} \mathbb{E}[(Y - \mathbb{E}[Y | \Phi])^2]$$

↓  
projection  
parameter

$$= \Sigma^{-1} \mathbb{E}[\Phi \cdot Y] \quad \text{unique!}$$

In particular  $\beta^*$  satisfies the normal equations

$$\Sigma \beta^* = \mathbb{E}[\Phi \cdot Y]$$

This implies that  $\forall a \in \mathbb{R}^d$

$$\mathbb{E}[(Y - \Phi^T \beta^*) \Phi^T a] = \mathbb{E}[(\mathbb{E}[Y | \Phi] - \Phi^T \beta^*) \Phi^T a] = 0$$

With this in mind, we have the following decomposition:

$$\begin{aligned}
 Y &= \Phi^T \beta^* + \underbrace{\left( \mathbb{E}[Y | \Phi] - \Phi^T \beta^* \right)}_{\text{non linearity}} + \underbrace{\left( Y - \mathbb{E}[Y | \Phi] \right)}_{\text{error}} \\
 &= \Phi^T \beta^* + \eta + \varepsilon \\
 &= \Phi^T \beta^* + \delta \\
 &\quad \downarrow \\
 &\quad \eta + \varepsilon
 \end{aligned}$$

Importantly  $\mathbb{E}[\delta^2] = \mathbb{E}[\eta^2] + \mathbb{E}[\varepsilon^2]$

Also:

i)  $\eta$  is orthogonal (in  $L_2$  sense) to the linear span of  $\Phi$ :

$$\mathbb{E}[\eta \cdot \Phi(j)] = 0 \quad j = 1, \dots, d$$

$$\Phi = \begin{bmatrix} \Phi(1) \\ \vdots \\ \Phi(d) \end{bmatrix}$$

ii)  $\varepsilon$  is orthogonal to all r.v.'s of the form  $f(\Phi)$  some  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  s.t.  $\text{Var}[f(\Phi)] < \infty$

In particular

$$\mathbb{E}[\varepsilon \cdot \eta] = 0$$

$\hookrightarrow$  function of  $\Phi$  (8)