

# SDS 387 Linear Models

Fall 2025

Lecture 19 - Tue, Nov 4, 2025

Instructor: Prof. Ale Rinaldo



OLS ESTIMATOR IN FIXED-DESIGN SETTING

Today we will assume

$$Y_i = \Phi_i^T \beta^* + \varepsilon_i \quad \text{where } \varepsilon_1, \dots, \varepsilon_n \sim \text{iid } (0, \sigma^2)$$

and

$\Phi_1, \dots, \Phi_n$  are fixed  
(deterministic)  
vectors in  $\mathbb{R}^d$

2 assumptions i) linearity

$$\mathbb{E}[Y_i] = \Phi_i^T \beta^*$$

(Aside, if the  $\Phi_i$ 's were random  
linearity means  $\mathbb{E}[Y_i | \Phi_i] = \Phi_i^T \beta^*$ )

ii) fixed covariates

- The fixed covariates assumption is unrealistic but you can assume that we are conditioning on the  $\Phi_i$ 's.

- When the model is linear and the covariates are random (i.e.  $\mathbb{E}[Y_i | \Phi_n] = \Phi_n^\top \beta^*$ ) the distribution of the covariates (hence the covariates themselves) is ancillary. So it is natural to condition on the  $\Phi_n$ 's. This is only true if linearity holds (otherwise we know that the projection parameter  $\beta^* = \mathbb{E}[\Phi \Phi^\top]^{-1} \mathbb{E}[Y \Phi]$  depends on the distribution of the covariates).
- See  
Buja et al. (2019)

- Remark: If  $\varepsilon_1, \dots, \varepsilon_n \stackrel{\text{iid}}{\sim} N(0, \sigma^2)$  then the likelihood of the data  $y_1, \dots, y_n$  is

$$\left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left\{-\frac{\sum_{i=1}^n \|y_i - \Phi_i \beta^*\|^2}{2\sigma^2}\right\}$$

and the OLS  $\hat{\beta}$  is the MLE of  $\beta^*$ .

- Now, for any  $\beta \in \mathbb{R}^d$ , the risk of  $\beta$  is

$$R(\beta) = \mathbb{E}_Y \left[ \frac{\|Y - \Phi \beta\|^2}{n} \right] = \mathbb{E}_{\varepsilon} \left[ \frac{\|\Phi (\beta^* - \beta) + \varepsilon\|^2}{n} \right]$$

vector in  $\mathbb{R}^n$   
of errors  
 $\varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}$

$$= (\beta^* - \beta)^\top \underbrace{\frac{1}{n} \Phi^\top \Phi}_{\text{H}} (\beta^* - \beta) + \mathbb{E}_{\varepsilon} \left[ \frac{\|\varepsilon\|^2}{n} \right]$$

$$= (\beta^* - \beta)^\top \sum_{i=1}^n (\beta^* - \beta) + \sigma^2$$

$$\frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^\top$$

$$= \|(\beta^* - \beta)\|_{\Sigma}^2 + \sigma^2$$



$$R(\beta^*)$$

The quantity  $R(\beta) - R(\beta^*) \geq 0$  is the excess risk.

- Remark: You can think of  $R(\beta)$  as  $\mathbb{E}_{Y \sim P} \left[ \frac{\|Y - \Phi\beta\|^2}{n} \right]$  if  $P$  were random (assuming linearity).

- So let's analyze  $R(\hat{\beta}) - R(\beta^*)$ , the excess risk of the OLS  $\hat{\beta}$ . Note the excess risk is a random variable! Let's compute its expectation:

- Remark: think of  $R(\hat{\beta})$  as  $\mathbb{E}_{Y_{\text{new}}} \left[ \frac{\|Y_{\text{new}} - \Phi\hat{\beta}\|^2}{n} \right]$  where  $Y_{\text{new}} \in \mathbb{R}^n$  is a new draw of data independent of the observations.

$$\mathbb{E}[R(\hat{\beta})] - R(\beta^*) = \mathbb{E}_{\hat{\beta}} \left[ (\beta^* - \hat{\beta})^T \underbrace{\Sigma}_{\|\beta^* - \hat{\beta}\|_{\Sigma}^2} (\beta^* - \hat{\beta}) \right]$$

$$\|\beta^* - \hat{\beta}\|_{\Sigma}^2 = \|\beta^* - \mathbb{E}[\hat{\beta}] + \mathbb{E}[\hat{\beta}] - \hat{\beta}\|_{\Sigma}^2 = \dots$$

= add (subtract)  $\mathbb{E}[\hat{\beta}]$

Exercise!

$$= \mathbb{E} \left[ \|\hat{\beta} - \mathbb{E}[\hat{\beta}]\|_{\Sigma}^2 \right] + \|\beta^* - \mathbb{E}[\hat{\beta}]\|_{\Sigma}^2$$

variance term for  $\hat{\beta}$

bias term

$$\mathbb{E} \left[ (\hat{\beta} - \mathbb{E}[\hat{\beta}])^T \Sigma (\hat{\beta} - \mathbb{E}[\hat{\beta}]) \right] = \mathbb{E} \left[ \text{tr} \left( \Sigma (\hat{\beta} - \mathbb{E}[\hat{\beta}]) (\hat{\beta} - \mathbb{E}[\hat{\beta}])^T \right) \right] \quad (3)$$

$$= \text{tr} \left( \underbrace{\mathbb{E}[(\hat{\beta} - \mathbb{E}[\hat{\beta}])(\hat{\beta} - \mathbb{E}[\hat{\beta}])^T]}_{\text{Var}[\hat{\beta}]} \right)$$

Next, let's evaluate the variance and bias of  $\hat{\beta}$ :

$$\text{bias : } \mathbb{E}[\hat{\beta}] = (\underline{\Phi}^T \underline{\Phi})^{-1} \underline{\Phi}^T \underbrace{\mathbb{E}[y]}_{\underline{\Phi} \beta^*} = \beta^* \quad \text{because } \underline{\Phi} \text{ is invertible}$$

↓  
no bias

This derivation is valid also when  $\underline{\Phi}$  is random:

$$\mathbb{E}[\hat{\beta}] = \mathbb{E}[\underbrace{\mathbb{E}[\hat{\beta} | \underline{\Phi}]}_{\beta^*}] = \beta^*$$

↓  
Linearity here is crucial.

$$\text{variance} \quad \text{Var}[\hat{\beta}] = \text{Var}[(\underline{\Phi}^T \underline{\Phi})^{-1} \underline{\Phi}^T y]$$

$$\begin{aligned} \text{Var}[Ay] &= (\underline{\Phi}^T \underline{\Phi})^{-1} \underline{\Phi}^T \underbrace{\text{Var}[y]}_{\sigma^2 I_n} \underline{\Phi} (\underline{\Phi}^T \underline{\Phi})^{-1} \\ &= \sigma^2 (\underline{\Phi}^T \underline{\Phi})^{-1} \end{aligned}$$

So we can now plug-in these expressions and get

$$\mathbb{E}[R(\hat{\beta})] - R(\beta^*) = \sigma^2 \frac{d}{n}$$

P/S We only need to evaluate

$$\begin{aligned}
\mathbb{E} \left[ \| \hat{\beta} - \underbrace{\mathbb{E}[\hat{\beta}]}_{\beta^*} \|_{\Sigma}^2 \right] &= \mathbb{E} \left[ \| \beta^* + (\Phi^\top \Phi)^{-1} \Phi^\top \varepsilon - \beta^* \|_{\Sigma}^2 \right] \\
(\Phi^\top \Phi)^{-1} \Phi^\top (\Phi \beta^* + \varepsilon) &= \mathbb{E} \left[ \| \sum_{i=1}^n \frac{\Phi^\top \varepsilon}{n} \|_{\Sigma}^2 \right] \\
&= \mathbb{E} \left[ \sigma^2 \frac{\Phi}{n} \sum_{i=1}^n \frac{\varepsilon_i \varepsilon_i^\top}{n} \frac{\Phi^\top \varepsilon}{n} \right] \\
&= \frac{1}{n} \mathbb{E} \left[ \underbrace{\varepsilon^\top \Phi (\Phi^\top \Phi)^{-1} \Phi^\top \varepsilon}_{H \text{ the hat matrix}} \right] \\
&= \frac{1}{n} \mathbb{E} [\varepsilon^\top H \varepsilon] \\
&= \frac{1}{n} \mathbb{E} [\text{tr}(H \varepsilon \varepsilon^\top)] \\
&= \frac{1}{n} \text{tr} \left( H \underbrace{\mathbb{E}[\varepsilon \varepsilon^\top]}_{\sigma^2 I_n} \right) \\
&= \frac{\sigma^2}{n} \text{tr}(H) \\
&= \frac{\sigma^2 d}{n} \quad \blacksquare
\end{aligned}$$

other proof:

$$\begin{aligned}
\mathbb{E} \left[ \underbrace{\| \hat{\beta} - \beta^* \|_{\Sigma}^2}_{\mathbb{E}[\hat{\beta}]} \right] &= \mathbb{E} \left[ \text{tr} \left( \sum (\hat{\beta} - \beta^*) (\hat{\beta} - \beta^*)^\top \right) \right] \\
&= \text{tr} \left( \sum \text{Var}[\hat{\beta}] \right) \\
&= \text{tr} \left( \sum \frac{\sigma^2}{n} \sum^{-1} \right) \\
&= \frac{\sigma^2}{n} \text{tr}(I_d) = \frac{\sigma^2 d}{n} \quad \blacksquare
\end{aligned}$$

- Remarks:
- i) The  $\frac{\sigma^2}{n}$  bound for the excess risk of  $\hat{\beta}$  (OLS) is optimal, in a minimax sense
  - ii) More refined analysis will give you high prob. bounds for excess risk.

iii) Recall that this result tells that

$$\begin{aligned} \mathbb{E} [R(\hat{\beta})] &= \mathbb{E}_{Y_{\text{new}}, Y} \left[ \frac{\|Y_{\text{new}} - \Phi \hat{\beta}\|^2}{n} \right] \\ &\downarrow \\ &= \sigma^2 \left( 1 + \frac{d}{n} \right) \end{aligned}$$

this is called the out-of-sample risk

What if we used the in-sample expected risk?

$$\mathbb{E} [\hat{R}(\hat{\beta})] = \mathbb{E}_Y \left[ \frac{\|Y - \Phi \hat{\beta}\|^2}{n} \right] = \sigma^2 \left( 1 - \frac{d}{n} \right)$$

*HW*

$\downarrow$

Wrong measure of risk! The risk is at least  $R(\beta^*) = \sigma^2$   
while  $\mathbb{E} [\hat{R}(\hat{\beta})] < \sigma^2$