

# SDS 387 Linear Models

Fall 2025

Lecture 25 - Tue, Dec 2, 2025

Instructor: Prof. Ale Rinaldo

## STATISTICAL INFERENCE IN ASSUMPTION-LEAN SETTINGS

Assume just that the pairs  $(Y, \Phi) \sim P_{Y, \Phi}$  on  $\mathbb{R} \times \mathbb{R}^d$  have each 2 moments.  $\mathbb{E}[Y^2] < \infty$  and  $\Sigma = \mathbb{E}[\Phi \Phi^T]$  is invertible

Then we saw that the projection parameter

$$\begin{aligned}\beta^* &= \arg\min_{\beta \in \mathbb{R}^d} \mathbb{E}[(Y - \Phi^T \beta)^2] \\ &= \arg\min_{\beta \in \mathbb{R}^d} \mathbb{E}[(\mathbb{E}[Y | \Phi] - \Phi^T \beta)^2] \\ &= \Sigma^{-1} \underbrace{\mathbb{E}[Y \cdot \Phi]}_{\Gamma \in \mathbb{R}^d}\end{aligned}$$

is well defined!

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

$$\sum_i \beta^* = \Gamma$$

- Using the theory of  $L_2$  projections, this implies that

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

$\forall a \in \mathbb{R}^d$

- So based on this

$$Y = \underbrace{\Phi^T \beta^* + (\mathbb{E}[Y | \Phi] - \Phi^T \beta^*)}_{\eta} + \underbrace{(Y - \mathbb{E}[Y | \Phi])}_{\varepsilon}$$

regression function

non-linearity  
(= 0 if the model  
is well specified)

intrinsic  
variability

$$= \Phi^T \beta^* + \delta$$

$\downarrow$   
 $\eta + \varepsilon$

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

variance term

because

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

by law of iterated expectation

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

$$\mathbb{E}[\eta \cdot \Phi(\omega)] = 0$$

$$j = 1, \dots, d$$

$\varepsilon$  is orthogonal to all r.v.'s of the form  $f(\Phi)$  any  $f$  s.t.

$$\text{Var}[f(\Phi)] < \infty$$

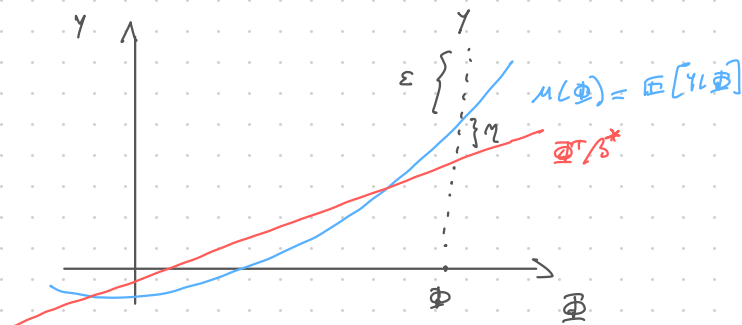
$$\hookrightarrow \mathbb{E}[\varepsilon \cdot \eta] = 0$$

- When the model is not well-specified the distribution of  $\Phi$  has to be taken into account because  $\beta^*$  depends on it. This is of course no longer the case when the model is well-specified (i.e.  $\mathbb{E}[Y|\Phi] = \Phi^T \beta^*$ )

$\hookrightarrow$  For a general discussion see Proposition 4.1 in Buja et al. (2019)

- Take home message: when the model is not linear we are facing an extra source of variability, namely the non-linearity!

$\hookrightarrow$  See Figure 1 in Buja et al. (2019)



- Remark: the issue is not just an increase in variance, but also the fact that the variance of  $\varepsilon$  and  $\eta$  given  $\Phi$  depends on  $\Phi \rightarrow$  variance not constant.

- Assume  $n$  iid pairs  $(Y_i, \Phi_i)$   $i=1, \dots, n$  from the unknown data generating distribution  $P_{Y, \Phi}$ .

We now show that we can estimate  $\beta^*$  using OLS

$$\hat{\beta} = \hat{\Sigma}^{-1} \hat{\Gamma} \quad \text{where}$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T \quad \hat{\Gamma} = \frac{1}{n} \sum_{i=1}^n Y_i \Phi_i$$

Then  $E[\hat{\beta}] \neq \beta^*$  (it is if the model is well-specified)

$$\text{Var}[\hat{\beta}] = E[\text{Var}[\hat{\beta} | \Phi_1, \dots, \Phi_n]] + \text{Var}[E[\hat{\beta} | \Phi_1, \dots, \Phi_n]]$$

= 0 when the model is well-specified (4)

- $\hat{\beta}$  is nonetheless a consistent estimator of  $\beta^*$

$$\hat{\beta} \xrightarrow{P} \beta^*$$

PA/ Recall that  $\hat{\beta} = \hat{\Sigma}^{-1} \cdot \hat{\Gamma}$  (assume that  $\hat{\Sigma}$  is invertible!)

$$\text{So WLLN } \hat{\Sigma} \xrightarrow{P} \Sigma = \mathbb{E}[\Phi \Phi^T]$$

$$\hat{\Gamma} \xrightarrow{P} \Gamma = \mathbb{E}[\Psi \cdot \Phi]$$

By CMT  $\hat{\Sigma}^{-1} \xrightarrow{P} \Sigma^{-1}$  so by Slutsky's

$$\hat{\beta} \xrightarrow{P} \Sigma^{-1} \Gamma = \beta^* \quad \square$$

Remark: this is much more complicated in high-dim settings  
where  $d = d(n)$ .

• CLT for  $\hat{\beta}$

- Recall: if  $\Phi$  is random and the model is linear

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

$$\hookrightarrow \Sigma = \mathbb{E}[\Phi \Phi^T]$$

- To establish a CLT for  $\hat{\beta}$  is an assumption heavy settings, let's consider these quantities:

$$\psi_i = \sum_{j=1}^n \Phi_j (Y_j - \Phi_j^T \beta^*) \quad \text{not computable!}$$

$i = 1, \dots, n$

Then

$$\frac{1}{n} \sum_{i=1}^n \psi_i = \Sigma^{-1} (\hat{\Gamma} - \hat{\Sigma} \beta^*)$$

Next,

$$\hat{\Sigma} (\hat{\beta} - \beta^*) = \hat{\Gamma} - \hat{\Sigma} \beta^*$$

$\Rightarrow$

$$\begin{aligned} \Sigma^{-1} \hat{\Sigma} \sqrt{n} (\hat{\beta} - \beta^*) &= \sqrt{n} \Sigma^{-1} (\hat{\Gamma} - \hat{\Sigma} \beta^*) \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_i \end{aligned}$$

• Now,  $E[\psi_i] = 0$  by the normal equations

$$\begin{aligned} \text{Var}[\psi_i] &= \Sigma^{-1} V \Sigma^{-1} \quad \text{where} \\ &\downarrow \text{by iid condition} \\ V &= \text{Var}[\Phi_n(Y_n - \Phi_n^T \beta^*)] \end{aligned}$$

the sandwich variance

• Remark if the model is well specified and

$$Y_n - \Phi_n^T \beta^* = \varepsilon_n \sim (0, \sigma^2) \quad \text{then}$$

$$\text{Var}[\psi_i] = \sigma^2 \Sigma^{-1}$$

• So by multivariate CLT

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_i \xrightarrow{d} N_d(0, \Sigma^{-1} V \Sigma^{-1})$$

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

But  $\Sigma^{-1} \hat{\Sigma} \sqrt{n}(\hat{\beta} - \beta^*) - \sqrt{n}(\hat{\beta} - \beta^*) = op(1)$

because  $(\Sigma^{-1} \hat{\Sigma} - I_d) \xrightarrow{p} 0$  by CMT

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

- Issue: we need to have a consistent estimator of the sandwich covariance. A natural estimator is the plug-in estimator:

$$\hat{\Sigma}^{-1} \hat{V} \hat{\Sigma}^{-1}$$

where  $\hat{V} = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T (y_i - \Phi_i^T \hat{\beta})^2$

- We need to show that

$$\hat{\Sigma}^{-1} \hat{V} \hat{\Sigma}^{-1} \xrightarrow{p} \Sigma^{-1} V \Sigma^{-1}$$

We already know that  $\hat{\Sigma}^{-1} \xrightarrow{p} \Sigma^{-1}$ . We need to show that

$$\hat{V} \xrightarrow{p} V$$

- We will first define

$$\tilde{V} = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T (y_i - \Phi_i^T \beta^*)^2 \quad (\text{not computable})$$

By WLLN  $\hat{V} \xrightarrow{P} V$

So all we need to do is to prove that

$$\hat{V} - \tilde{V} \xrightarrow{P} 0$$

We have

$$\begin{aligned} \hat{V} - \tilde{V} &= \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T \left[ (\Phi_i^T \hat{\beta})^2 - (\Phi_i^T \beta^*)^2 + 2 Y_i \Phi_i^T (\beta^* - \hat{\beta}) \right] \\ &= \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T \left[ \underbrace{(\Phi_i^T (\hat{\beta} - \beta^*))^2}_{\downarrow} + 2 (Y_i - \Phi_i^T \beta^*) \Phi_i^T (\beta^* - \hat{\beta}) \right] \end{aligned}$$

Next,

$$\|\hat{V} - \tilde{V}\|_{op} \leq \frac{1}{n} \sum_{i=1}^n \|\Phi_i\|^2 \left[ \begin{array}{c} \downarrow \\ \cdot \end{array} \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \|\Phi_i\|^2 (\Phi_i^T (\hat{\beta} - \beta^*))^2 + \frac{2}{n} \sum_{i=1}^n \|\Phi_i\| |Y_i - \Phi_i^T \beta^*| \|\Phi_i\| |\Phi_i^T (\beta^* - \hat{\beta})|$$

by Cauchy Schwarz

$$\begin{aligned} &\leq \underbrace{\frac{1}{n} \sum_{i=1}^n \|\Phi_i\|^2 (\Phi_i^T (\hat{\beta} - \beta^*))^2}_A + 2 \underbrace{\sqrt{\frac{1}{n} \sum_{i=1}^n \|\Phi_i\|^4 (Y_i - \Phi_i^T \beta^*)^2}}_B \underbrace{\sqrt{\frac{1}{n} \sum_{i=1}^n \|\Phi_i\|^2 |\Phi_i^T (\beta^* - \hat{\beta})|}}_A \\ &= A + 2 \sqrt{A} \sqrt{B} \end{aligned}$$

Next by Cauchy Schwarz

$$A \leq \underbrace{\left[ \frac{1}{n} \sum_{i=1}^n \|\Phi_i\|^4 \right]}_{\xrightarrow{P} \mathbb{E}[\|\Phi_n\|^4]} \|\hat{\beta} - \beta^*\| \xrightarrow{P} 0$$

$\xrightarrow{P} \mathbb{E}[\|\Phi_n\|^4]$  which we assume to be finite



As for B:

$$B \rightarrow \text{tr} \left( \text{Var} \left( \mathbb{E}_1 (y_1 - \mathbb{E}_1^\top \beta^*) \right) \right) = \text{tr}(V)$$

also finite

By Slutsky's

$$A + 2\sqrt{A} \sqrt{B} \xrightarrow{P} 0$$

$\hookrightarrow$

For large  $n$ :

$$\sqrt{n}(\hat{\beta} - \beta^*) \approx N_d(0, \hat{\Sigma}^{-1} \hat{V} \hat{\Sigma}^{-1})$$

- Remark: To carry out this program in high-dim settings is highly non-trivial.