

# SDS 387 Linear Models

Fall 2024

Lecture 25 - Tue, Dec 3, 2024

Instructor: Prof. Ale Rinaldo

- Last time: Assumption - lean inference :
  - ↪ see Statistical Science paper Models as approximations, part I
  - linear model is mis-specified
  - covers over our random
- White (1980) Consequences and Detection of mis-specified non-linear regression Models, JASA 76, 374-419-433
- $(\Phi, Y) \sim P_{\Phi, Y}$  on  $\mathbb{R}^{d+1}$  but no assumptions on the regression function  $x \in \mathbb{R}^d \mapsto \mathbb{E}[Y | \Phi = x]$  is made. We only assume 2<sup>nd</sup> moments for  $Y$  and  $\Phi$ .
- We can always write

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

↓  
signal  
regression function

↓  
noise, error  
where  $\mathbb{E}[\varepsilon | \Phi] = 0$   
 $\mathbb{E}[\varepsilon] = 0$

- We saw (and you should do it as an exercise) that, even if the model is not linear, the projection parameter

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

$$= \Sigma^{-1} \Gamma$$

where  $\Sigma = \mathbb{E}[\Phi \Phi^T]$  and  $\Gamma = \mathbb{E}[\Phi \cdot Y]$

assuming that  $\Sigma$  is invertible (and assuming  $\mathbb{E}[Y^2] < \infty$ )

Furthermore  $\beta^*$  satisfies

$$\Sigma \beta^* = \Gamma \quad \text{normal equations}$$

- $\beta^*$  is the focus of inference: vector of coefficients of the "best" approximation of  $Y$

↓  
measure of linear association  
btw  $Y$  and  $\Phi$

or  $\mathbb{E}[Y | \Phi]$  by linear functions  
of  $\Phi$ .

- Last time we saw a fundamental decomposition:

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

So  $Y = \Phi^T \beta^* + \varepsilon$   
 $\hookrightarrow \eta + \varepsilon$

Remark :  $E[\varepsilon^2] = E[\eta^2] + E[\varepsilon^2]$

i)  $\eta$  is orthogonal to the linear span  $\Phi$   $\rightarrow \{ \sum \Phi, \sum \in \mathbb{R}^d \}$

$$\left[ E[\eta \cdot \Phi(j)] = 0 \quad \forall j \right]$$

$\downarrow$   
 $j^{\text{th}}$  coordinate of  $\Phi$

ii)  $\varepsilon$  is orthogonal to all r.v.'s of the form  $f(\Phi)$  where  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  s.t.  
 $E[f^2(\Phi)] < \infty$   
 $\hookrightarrow$  as a result,  $E[\eta \cdot \varepsilon] = 0$

• Now the distribution of  $\Phi$  has to be taken into account, because  $\beta^*$  depends on it.

if  $Y = \Phi^T \beta^* + \varepsilon$   
 $\downarrow$   $\downarrow$  mean zero  
same  $\beta^*$

then  $\beta^*$  does not depend on the distribution of  $\Phi$

- Nonlinearity + random covariates  $\rightarrow$  extra uncertainty
- Assume  $n$  iid observations from  $P_{\Phi, Y}$ :

$$(\Phi_1, Y_1), \dots, (\Phi_n, Y_n) \stackrel{iid}{\sim} P_{\Phi, Y} \rightarrow \text{unknown}$$

$$\text{Let } \Phi_{n \times (d+1)} = \begin{bmatrix} 1 & \Phi_1^T \\ 1 & \Phi_2^T \\ \vdots & \vdots \\ 1 & \Phi_n^T \end{bmatrix} \text{ be the random design matrix}$$

and consider the OLS estimator:

$$\hat{\beta} = \hat{\Sigma}^{-1} \cdot \hat{\Gamma}$$

$$\text{where } \hat{\Sigma} = \frac{\Phi^T \Phi}{n} = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^T$$

plug-in estimator for  $\beta^*$

$$\hat{\Gamma} = \frac{1}{n} \sum_{i=1}^n Y_i \cdot \Phi_i$$

Remark:  $\mathbb{E}[\hat{\beta}] \neq \beta^*$

$$\text{Var}[\hat{\beta}] = \mathbb{E}[\text{Var}[\hat{\beta} | \Phi]] + \text{Var}[\mathbb{E}[\hat{\beta} | \Phi]]$$

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \downarrow$$

if  $\mathbb{E}[Y | \Phi] = \Phi \beta^*$  fixed

then  $\mathbb{E}[\hat{\beta} | \Phi] = \beta^*$

so  $\text{Var}[\mathbb{E}[\hat{\beta} | \Phi]] = 0$

• Consistency :

$$\hat{\beta} \xrightarrow{P} \beta^* \quad \text{as } n \rightarrow \infty \quad (\text{keeping } d \text{ fixed!})$$

PA/  $\hat{\beta} = \hat{\Sigma}^{-1} \cdot \hat{\Gamma}$

Now  $\hat{\Sigma} \xrightarrow{P} \Sigma$  by WLLN and  $\hat{\Sigma}^{-1} \xrightarrow{P} \Sigma^{-1}$  by CMT.

Next  $\hat{\Gamma} \xrightarrow{P} \Gamma$  by WLLN

So  $\hat{\beta} = \hat{\Sigma}^{-1} \hat{\Gamma} \xrightarrow{P} \Sigma^{-1} \Gamma = \beta^*$  by Slutsky's theorem.  $\Rightarrow$

Remark : When  $d$  grows with  $n$ , it is still not known how to eliminate the bias  $E[\hat{\beta}] - \beta^*$  efficiently.

• CLT for  $\hat{\beta}$

To establish a CLT for  $\hat{\beta}$ , define

$$\psi_i = \Sigma^{-1} \Phi_n(Y_i - \Phi_n^T \beta^*) \in \mathbb{R}^d$$

$$i = 1, \dots, n$$

Then :

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

Next,

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