

# SDS 387

## Linear Models

Fall 2025

Lecture 24 - Thu, Nov 20, 2025

Instructor: Prof. Ale Rinaldo

### RANDOM DESIGN LINEAR REGRESSION

We still assume a well-specified linear regression model, but allow for the covariates (e.g. features) to be random. So now the model is

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

where  $\Phi_i$ ,  $i=1, \dots, n$ , are iid random vectors in  $\mathbb{R}^d$  from some distribution  $P_\Phi$  and the errors are s.t.

$$\varepsilon_1, \dots, \varepsilon_n | \Phi_1, \dots, \Phi_n \stackrel{\text{iid}}{\sim} (0, \sigma^2)$$

Note: the  $\Phi_i$ 's are ancillary for estimating

$\beta^*$  because the model is well specified

$\hookrightarrow$  it is ok to condition on the  $\Phi_i$ 's. ①

↳ See Buja et al. (2019) Stat Science

- So, we now observe  $n$  lot pairs

$$(Y_1, \Phi_1), \dots, (Y_n, \Phi_n) \text{ in } \mathbb{R} \times \mathbb{R}^d$$

↳ Remark: the first coordinate of each  $\Phi_i$  is a 1, to allow for an intercept.

- Now the risk function is defined as:

$$\beta \in \mathbb{R}^d \mapsto R(\beta) = \mathbb{E}_{\substack{Y, \Phi \\ \text{or } \varepsilon, \Phi}} [(Y - \Phi^\top \beta)^2]$$

↓  
Think of  $Y$  and  $\Phi$  as  $(Y_{\text{new}}, \Phi_{\text{new}})$  a new lot drawn from  $P_{Y, \Phi}$

- Prop. 3.9 in Bach's book. (expression for the predictive risk).

Let  $\Sigma = \mathbb{E} [\Phi \Phi^\top] \geq 0$

Then,  $\forall \beta \in \mathbb{R}^d$ ,

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

$$\begin{aligned}
 R(\beta) &= \mathbb{E}[(y - \Phi^T \beta)^2] = \mathbb{E}\left[\left(y - \Phi^T \beta^* + \Phi^T (\beta^* - \beta)\right)^2\right] \\
 &= \mathbb{E}\left[(y - \Phi^T \beta^*)^2\right] + \mathbb{E}\left[\left(\Phi^T (\beta^* - \beta)\right)^2\right] \\
 &\quad + 2 \underbrace{\mathbb{E}\left[(y - \Phi^T \beta^*)(\Phi^T (\beta^* - \beta))\right]}_{CV}
 \end{aligned}$$

- If the model is linear, as we assume, then

$$y - \Phi^T \beta^* = \varepsilon \text{ is such that}$$

$$\mathbb{E}[\varepsilon | \Phi] = 0 \text{ so that}$$

$$\begin{aligned}
 \mathbb{E}[CV] &= \mathbb{E}\left[\mathbb{E}[CV | \Phi]\right] \\
 &= \mathbb{E}_{\Phi} \left[ \mathbb{E}_{\varepsilon | \Phi} \left[ \varepsilon \Phi^T (\beta^* - \beta) \right] \right] \\
 &= \mathbb{E}_{\Phi} \left[ \Phi^T (\beta^* - \beta) \underbrace{\mathbb{E}_{\varepsilon | \Phi} [\varepsilon]}_{=0} \right] \\
 &= 0.
 \end{aligned}$$

- If is not linear, then our  $\beta^*$  is the projection parameter

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

projection  
parameter,  
the coefficients of  
the L<sub>2</sub> projection  
of  $Y$  (or of  $\mathbb{E}[Y|\Phi]$ )  
onto the linear span of  $\Phi$

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

$$= \Sigma^{-1} \mathbb{E}[Y \cdot \Phi] \quad \text{assuming } \Sigma > 0 \text{ and } \mathbb{E}[Y^2] < \infty$$

In this case

$\mathbb{E}[CV] = 0$  be the defining properties of L<sub>2</sub> projection,  
namely that  $Y - \Phi^T \beta^*$  is uncorrelated  
with any linear function of  $\Phi$ .

• Regardless,  $\mathbb{E}[CV] = 0$

↳

$$R(\beta) = \mathbb{E}[(\Phi^T(\beta^* - \beta))^2] + \mathbb{E}[(Y - \Phi^T \beta^*)^2]$$

$$= \|\beta - \beta^*\|_{\Sigma}^2 + \sigma^2 \quad \text{if the model is well-specified}$$

$$= \|\beta - \beta^*\|_{\Sigma}^2 + \mathbb{E}[(Y - \mathbb{E}[Y|\Phi])^2]$$

$$\xrightarrow{\text{non-linearity}} + \mathbb{E}[(\mathbb{E}[Y|\Phi] - \Phi^T \beta^*)^2]$$

$\hookrightarrow$  Thus  $R(\beta)$  decomposes as a sum of

$$\|\beta - \beta^*\|_1^2 \Sigma \quad \text{and an intrinsic, irreducible error: either } \sigma^2 \text{ or } \sigma^2 + \eta^2$$

■

$\hookrightarrow$  All we can do to minimize the risk is to minimize  $\|\beta - \beta^*\|_1^2 \Sigma$  because  $\sigma^2$  (or  $\sigma^2 + \eta^2$ ) do not depend on  $\beta$ .

The excess risk now is

$$R(\beta) = \begin{cases} \sigma^2 & \text{linear model} \\ \sigma^2 + \eta^2 & \text{mis-specified model.} \end{cases}$$

Assume we have data ( $n$  wk. pairs  $(y_i, \Phi_i)$ )  
 $i=1, \dots, n$ )

then we can compute the OLS estimator

$$\hat{\beta} = \sum_{i=1}^n \sum_{j=1}^n y_i \Phi_i \frac{\downarrow}{n}$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \Phi_i \Phi_i^\top = \hat{\Sigma}_n [\Phi \Phi^\top]$$

- the excess risk of  $\hat{\beta}$  is  $R(\hat{\beta}) = \|\hat{\beta} - \beta^*\|_1^2 \hat{\Sigma}$

$\downarrow$   
 random variable (5)

Prop 3.10 The expected excess risk of  $\hat{\beta}$  is:

$$\mathbb{E}[R(\hat{\beta})] = \frac{\sigma^2}{n} \mathbb{E}[\text{tr}(\Sigma \hat{\Sigma}^{-1})]$$

Remark  $\mathbb{E}[\text{tr}(\Sigma \hat{\Sigma}^{-1})] \geq d$  because

on the cone of PD matrices the map

$$A \mapsto \text{tr}(A^{-1}) \text{ is convex}$$

$$\hookrightarrow \mathbb{E}[\text{tr}(\Sigma \hat{\Sigma}^{-1})] \geq \text{tr}(\mathbb{E}[\hat{\Sigma}] \hat{\Sigma}^{-1}) = d$$

see page  
64 of  
Bach's book

PP/ Write  $\underline{\Phi}$  for the nxd matrix with

$$\text{rows } \underline{\Phi}_1^T, \dots, \underline{\Phi}_n^T.$$

$$\text{so } \hat{\Sigma}_i = \frac{1}{n} \underline{\Phi}^T \underline{\Phi}$$

$$\text{Similarly let } Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix} \in \mathbb{R}^n$$

So

$$\hat{\beta} = \hat{\Sigma}^{-1} \frac{\underline{\Phi}^T Y}{n} = \beta^* + \hat{\Sigma}^{-1} \frac{\underline{\Phi}^T \varepsilon}{n}$$

$$Y = \underline{\Phi} \beta^* + \varepsilon$$

Remark We are assuming throughout that  $\hat{\Sigma}$  is invertible!

$$\text{so } \mathbb{E} \left[ \| \hat{\beta} - \beta^* \|^2 \Sigma \right] = \mathbb{E} \left[ \| \hat{\Sigma}^{-1} \frac{\Phi^T \Sigma}{n} \|^2 \Sigma \right]$$

expected excess  
risk of  $\hat{\beta}$

$$= \mathbb{E} \left[ \varepsilon^T \frac{\Phi}{n} \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \frac{\Phi^T \Sigma}{n} \right]$$

$$= \mathbb{E} \left[ \text{tr} \left( \Sigma \left( \hat{\Sigma}^{-1} \frac{\Phi^T \Sigma}{n} \right) \left( \hat{\Sigma}^{-1} \frac{\Phi^T \Sigma}{n} \right)^T \right) \right]$$

$$= \mathbb{E}_{\Phi, \Sigma} \left[ \text{tr} \left( \Sigma \hat{\Sigma}^{-1} \frac{\Phi^T \Sigma}{n} \varepsilon^T \frac{\Phi}{n} \hat{\Sigma}^{-1} \right) \right]$$

$$= \mathbb{E}_{\Phi} \left[ \mathbb{E}_{\Sigma | \Phi} \left[ \cdot \right] \right]$$

$$= \mathbb{E}_{\Phi} \left[ \frac{1}{n} \text{tr} \left( \Sigma \hat{\Sigma}^{-1} \frac{\Phi^T}{n} \underbrace{\mathbb{E}_{\Sigma | \Phi} [\Sigma \varepsilon^T \Phi]}_{= 0^2 I_n} \Phi^T \hat{\Sigma}^{-1} \right) \right]$$

$$= \frac{0^2}{n} \mathbb{E}_{\Phi} \left[ \text{tr} \left( \Sigma \hat{\Sigma}^{-1} \frac{\Phi^T \Phi}{n} \hat{\Sigma}^{-1} \right) \right]$$

$$= \frac{0^2}{n} \mathbb{E}_{\Phi} \left[ \text{tr} (\Sigma \hat{\Sigma}^{-1}) \right]$$

②

- Theorem 1 by Mourtada (2022) (AoS 2024), 2157-2178)

a) Assume that  $d > n$  or the distribution of

$\hat{\Sigma}$  is not invertible  $\leftarrow$  the  $\Phi_i$ 's is degenerate (supported on a affine subspace of  $\mathbb{R}^d$ ) (7)

Then the minimax risk is infinity!

ii) if  $n \geq d$  and  $\hat{\Sigma}$  is invertible then

the minimax rule

$$\inf_{\tilde{\beta}} \sup_{\beta \in \mathbb{R}^d} \mathbb{E}_{\beta} [R(\tilde{\beta})] = \frac{\sigma^2}{n} \mathbb{E} [\text{tr}(\Sigma \tilde{\Sigma}^{-1})]$$

↓                      ↓  
 estimator            expectation w.r.t.

OLS is minimax optimal !!