

# SDS 387 Linear Models

Fall 2025

Lecture 18 - Tue, Oct 30, 2025

Instructor: Prof. Ale Rinaldo

- Recall that if  $R(\beta) = \mathbb{E}[(y - x^T \beta)^2]$  is the prediction risk associated to  $\beta \in \mathbb{R}^d$  then  
then  
$$R(\beta) = \underbrace{\| \beta - \beta^* \|_{\Sigma}^2}_{\text{estimation}} + \underbrace{\mathbb{E}[(y - \mathbb{E}[y|x])^2]}_{\sigma^2 \text{ variance}} + \underbrace{\mathbb{E}[(\mathbb{E}[y|x] - x^T \beta^*)^2]}_{\eta^2 \text{ non linearity}}$$

where  $\beta^* = \mathbb{E}[xx^T]^{-1} \mathbb{E}[y \cdot x]$   
 $\downarrow$   
projection parameter

- If  $\beta = \beta^*$  then  $R(\beta^*) = \sigma^2 + \eta^2$

$$= \inf_{\beta \in \mathbb{R}^d} R(\beta)$$

①

because  $\|\beta - \beta^*\|^2 \leq 0$  if  $\beta = \beta^*$

and  $\inf_{\beta \in \mathbb{R}^d} \mathbb{E} \left[ \left( \mathbb{E}[Y|X] - X^T \beta \right)^2 \right] = \eta^2$

$$= \mathbb{E} \left[ \left( \mathbb{E}[Y|X] - X^T \beta^* \right)^2 \right]$$

↓

$$0 \leq R(\beta) - R(\beta^*) \quad \text{Excess risk of } \beta$$

- Of course if linearity holds, i.e.  $\mathbb{E}[\varepsilon] = 0$   
 $Y = X^T \beta^* + \varepsilon$   $\varepsilon \perp X$

then  $\eta = 0$  and

$$R(\beta^*) = \sigma^2 = \mathbb{E} \left[ (Y - X^T \beta^*)^2 \right]$$
$$= \mathbb{E} [\varepsilon^2] = \text{var}[\varepsilon]$$

- Back to last lecture:  $(X_i, Y_i)_{i=1, \dots, n} \stackrel{\text{i.i.d.}}{\sim} P_{X,Y}$

$$Y = \begin{bmatrix} Y_1 \\ \vdots \\ Y_n \end{bmatrix} \in \mathbb{R}^n$$

$$\Phi = \begin{bmatrix} \Phi_1^T \\ \vdots \\ \Phi_n^T \end{bmatrix}$$

$n \times d$

↓

design matrix or feature matrix (2)

$$\Phi_n = \varphi(x_n) \in \mathbb{R}^d \quad \text{with feature}$$

The plug-in estimator of  $\beta^*$  (either the projection parameter or the linear parameter if the model is linear)

$$\text{is } \hat{\beta} = \underset{\beta \in \mathbb{R}^d}{\operatorname{argmin}} \hat{R}(\beta) \quad \text{where}$$

$$\begin{aligned} \text{empirical risk} \quad \hat{R}(\beta) &= \hat{\mathbb{E}}_n [(Y_i - \Phi_i^\top \beta)^2] \\ &= \frac{1}{n} \sum_{i=1}^n (Y_i - \Phi_i^\top \beta)^2 \\ &= \frac{1}{n} \| Y - \Phi \beta \|^2 \end{aligned}$$

$\begin{matrix} n \times 1 & n \times d & d \times 1 \end{matrix}$

Then  $\hat{\beta} = (\Phi^\top \Phi)^{-1} \Phi^\top Y$

OLS estimator

provided that  $\Phi$  has full column rank ( $=d$ )

$$\hookrightarrow \Phi^\top \Phi \text{ invertible}$$

$\begin{matrix} d \times d \end{matrix}$

This requires  $n \geq d$  (check)!

PA/ The function  $\beta \in \mathbb{R}^d \rightarrow \hat{R}(\beta) \in \mathbb{R}$   
 is strictly convex because its Hessian  
 ( $d \times d$  matrix of mixed 2<sup>nd</sup> derivatives) is

$$\hat{\Sigma} = \frac{1}{n} \Phi^T \Phi \succ 0 \quad \text{by assumption for all } \beta \in \mathbb{R}^d$$

So  $\hat{\beta}$  is the only vector s.t.

$$\nabla \hat{R}(\hat{\beta}) = 0$$

$$\Updownarrow$$

$$-\frac{2}{n} \Phi^T (Y - \Phi \hat{\beta}) = 0$$

$$\Updownarrow$$

$$\Phi^T \Phi \hat{\beta} = \Phi^T Y$$

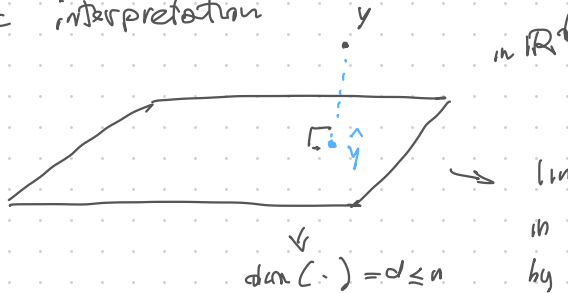
Normal equations

Using inverting of  $\Phi^T \Phi$  we get

$$\hat{\beta} = (\Phi^T \Phi)^{-1} \Phi^T Y$$

$\square$

• Geometric interpretation



linear subspace  
 in  $\mathbb{R}^n$  spanned  
 by columns of  $\Phi$ , i.e.  
 $\text{col}(\Phi)$   $\textcircled{4}$

where  $\hat{Y} = \Phi \hat{\beta} = \underbrace{\Phi (\Phi^T \Phi)^{-1} \Phi^T}_H Y$

$\downarrow$   
 fitted values

$= HY$

where  $H$ , the hat matrix, is the orthogonal projection onto  $C(\Phi)$ . ( $H$  is symmetric and  $H^2 = H \cdot H = H$ ).

Next  $e = Y - \hat{Y} = \underbrace{(I_n - H)}_{\text{also an orthogonal projection (symmetric + idempotent) onto } C(\Phi)^\perp} Y \in \mathbb{R}^n$

$\downarrow$   
 residuals

• Direct sum decomposition

$$Y = \hat{Y} + e \quad \text{where} \quad \langle \hat{Y}, e \rangle = 0$$

$$\frac{\|Y\|^2}{n} = \frac{\|\hat{Y}\|^2}{n} + \frac{\|e\|^2}{n}$$

$\updownarrow$

$$\frac{\sum y_i^2}{n} = \frac{\sum \hat{y}_i^2}{n} + \frac{\sum e_i^2}{n}$$

$\downarrow$

$\rightarrow$  least squares error

"proportion of energy" explained by the model

• Numerical considerations. How do you compute  $\hat{\beta}$ ?

1) Requires matrix inversion of  $\Phi^T \Phi$   
order  $d^3$  in fact order  $n \cdot d^2$

2) Gradient descent iterative procedure that  
starting from  $\beta_0 \in \mathbb{R}^d$  (i.e.  $\beta_0 = 0$  or  
any point in the orthogonal complement  
of  $\text{kernel}(\Phi)$ )

apply the following recursion:

$$t \geq 1 \quad \beta_t = \beta_{t-1} + \gamma \underbrace{\nabla \frac{\hat{R}}{2}}_{>0 \text{ stepsize or learning rate}}(\beta_{t-1})$$

$$= \beta_{t-1} - \gamma \left( \Phi^T (\Phi \beta_{t-1} - y) \right)$$

This has complexity  $n \cdot d$  and as  $t \rightarrow \infty$

$$\beta_t \rightarrow \hat{\beta}$$

↓  
See section  
5.2 of Bach's  
book

- What if  $\Phi^T \Phi$  is not invertible? Assume  $\text{rank}(\Phi) = n$  (for example,  $d > n$ ).

Then the normal equations

$$\Phi^T \Phi \beta = \Phi^T y$$

have infinitely many solutions, because if say  $\hat{\beta}$  solve the normal equations, so does  $\hat{\beta} + v$

where  $v \in \text{ker}(\Phi)$

$$\hookrightarrow \{x \in \mathbb{R}^d : \Phi x = 0\}$$

because

$$\begin{aligned} \Phi^T \Phi (\hat{\beta} + v) &= \Phi^T \Phi \hat{\beta} + \underbrace{\Phi^T \Phi v}_{=0} \\ &= \Phi^T \Phi \hat{\beta} \end{aligned}$$

- When  $\text{rank}(\Phi) = n$  then any solution  $\hat{\beta}$  to the normal equation is such that

$$\hat{y} = \Phi \hat{\beta} = y$$

i.e.  $\hat{R}(\hat{\beta}) = 0$  and the model

interpolate the data

- Among the infinitely many solutions to the normal equations there is a canonical one:

the min-norm solution, the one with smallest Euclidean norm. This is defined as

$$\hat{\beta}_{MN} = (\Phi^T \Phi)^+ \Phi^T Y$$

where for an  $A_{m \times n}$  its Moore-Penrose

pseudo inverse  $A^+$  is a  $n \times m$  matrix s.t.

i)  $AA^+A = A$  ( $AA^+$  maps columns of  $A$  to themselves, so it is an identity on  $C(A)$ )

ii)  $A^+AA^+ = A^+$

iii)  $AA^+$  is symmetric  
 $A^+A$

- Extra properties  $\text{kernel}(A^+) = \text{kernel}(A^T)$   
 $C(A^+) = C(A^T)$

$AA^+$  and  $A^+A$  are idempotent

↓  
orthogonal projection  
onto  $C(A)$

↓  
orthogonal projection  
onto  $C(A^T) \rightarrow$  row space of  $A$



• If  $A = U \Sigma V^T$

$m \times n$        $m \times m$      $k \times k$      $k \times n$

where  $\Sigma$  is diagonal with positive diagonal elements (the singular values) and  $k = \text{rank}(A)$

Then

$$A^+ = V \Sigma^{-1} U^T$$

• Back to interpolation (i.e.  $\text{rank}(\Phi) = n$ )

$$\hat{\beta}_{MN} = \Phi^+ Y = \text{argmin} \{ \|Y\| \text{ s.t. } \Phi \beta = Y \}$$

and gradient descent  $\rightarrow \hat{\beta}_{MN}$