

# Lecture FYS- STK3155/4155, August 26, 2024

$$C(\beta) = \text{MSE} = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \sum_{j=0}^{p-1} \beta_j x_i^j)^2$$

$$\hat{y}_i = \sum_{j=0}^{p-1} \beta_j x_i^j \quad (\text{Model})$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

↑  
optimal

$X$  orthogonal  $\quad X^T X = X X^T = I$

$$= \begin{bmatrix} 1 & & & \\ & 1 & & 0 \\ & & 1 & \\ & 0 & & 1 \end{bmatrix}$$

$$\hat{y} = X \hat{\beta} = X \underbrace{(X^T X)^{-1}}_I X^T y$$

$$= y$$

$$A = X (X^T X)^{-1} X^T$$

$A^2 = A$  (Projection matrix)

$$\tilde{y} = Ay$$

$\tilde{y}$  can be viewed as a projection of  $y$  in terms of the column vectors of  $X$

$$X = \begin{bmatrix} 1 & x_0 & x_0 & \dots & x_0^{p-1} \\ 1 & x_1 & & & \\ \vdots & \vdots & & & \\ 1 & x_{n-1} & \dots & x_{n-1}^{p-1} \end{bmatrix}$$

$$\frac{\partial C}{\partial \beta} = 0 = -\frac{2}{n} X^T (Y - X\beta)$$

Ex 1

$$\frac{\partial^2 C}{\partial \beta \partial \beta^T} = \frac{2}{n} \overbrace{X^T X}^{\text{Hessian}}$$

(MSE)

SVD of a matrix

Decompose Matrix  $X$

$$X = U \Sigma V^T$$

$$X \in \mathbb{R}^{n \times p}$$

$$U \in \mathbb{R}^{n \times n}$$

$$V \in \mathbb{R}^{p \times p}$$

$$U U^T = U^T U = I$$

$$V V^T = V^T V = I$$

$$U = [u_0 \ u_1 \ \dots \ u_{n-1}]$$

$$u_i^T u_j = S_{ij}$$

$$V = [v_0 \ v_1 \ \dots \ v_{p-1}]$$

$$v_i^T v_j = S_{ij}$$

$$\Sigma \in \mathbb{R}^{n \times p} = \begin{bmatrix} \sigma_0 & 0 & \dots & 0 \\ 0 & \sigma_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_{p-1} \\ 0 & 0 & \dots & 0 \end{bmatrix}$$

$\sigma_0 > \sigma_1 > \sigma_2 \dots > \sigma_{p-1} > 0$

Example  $3 \times 2$

$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$\Sigma^T \Sigma = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} \sigma_0^2 & 0 \\ 0 & \sigma_1^2 \end{bmatrix}$$

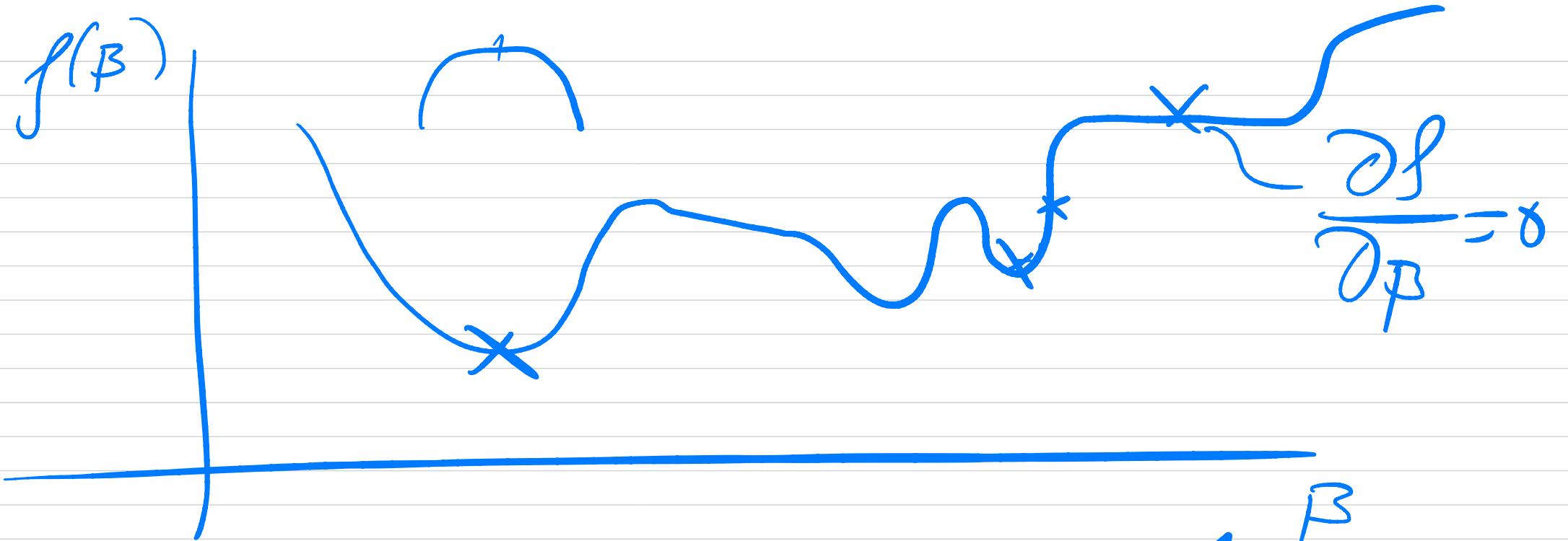
$$\Sigma \Sigma^T = \begin{bmatrix} 4 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix}$$

$$\cancel{\underset{P \times P}{XX^T}} = V \Sigma^T \underset{P \times n}{U^T} \underset{1 \times m}{U} \underset{n \times P}{\Sigma} V^T$$

$$= V \Sigma^T \Sigma V^T$$

$$X^T V = V \Sigma^T \Sigma \underset{V^T V}{\underbrace{V^T V}}$$

$$X^T X w_i = w_i^T \Sigma^2 \xrightarrow{\text{II } P \times P} \text{convex optimization}$$



$$\hat{y} = X\beta = X(X^T X)^{-1} X^T y$$

with SVD

$$X = U \Sigma V^T$$

$$\tilde{y}_{OLS} = \underbrace{U \Sigma V^T}_{m \times p} \left( \underbrace{V \Sigma^T \Sigma V^T}_{p \times p} \right)^{-\frac{1}{2}} \underbrace{V \Sigma^T U^T y}_{p \times m}$$

$$\overline{V \Sigma^T \Sigma V^T}$$

A, B are  
square invertible  
matrices

$$\frac{1}{AB} = A^{-1} B^{-1}$$

$$\tilde{y} = \left( \sum_{j=0}^{p-1} u_j u_j^T \right) y$$

if  $\forall \tau_i > 0$

$$\tilde{y} = \underbrace{\sum_{j=0}^{n-1} \gamma_j u_j u_j^T y}_1$$

$$y' = y$$

why interesting?  
our next method:  
Ridge Regression.

Trick in case you cannot  
invert  $\mathbf{X}^T \mathbf{X}$

$$\mathbf{X}^T \mathbf{X} \Rightarrow \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_{p \times p}$$

$$\hat{\beta}_{\text{Ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

$\lambda > 0$

And then cost function

$$C_{\text{Ridge}} = \frac{1}{n} \sum_{i=0}^{n-1} \left( y_i - \sum_{j=0}^{p-1} \beta_j x_i^j \right)^2$$

$$+ \lambda \sum_{j=0}^{p-1} \beta_j^2$$

Hyperparameter

constraint

$$\sum_{j=0}^{p-1} \beta_j^2 < t < \infty$$

$$CLASSO = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \sum_{j=0}^{p-1} \beta_j x_i^j)^2 + \lambda \sum_{j=0}^{p-1} |\beta_j|$$

$\overbrace{\quad\quad\quad}^{\lambda \|\beta\|_1}$

$$\frac{\partial C_{Ridge}}{\partial \beta} = 0 = \frac{-2}{n} X^T (y - X\beta) - 2\lambda \beta$$

$$-X^T(Y - X\beta) - \underbrace{\lambda}_{\lambda \rightarrow \lambda} \beta$$

$$\beta_{\text{Ridge}} = \left( X^T X + \lambda I \right)^{-1} X^T Y$$

## Lasso

$$\frac{d|F|}{d\beta} = \begin{cases} +1 & \beta > 0 \\ -1 & \beta < 0 \end{cases}$$

$$\hat{y}_{\text{Ridge}} = X \underbrace{U \Sigma V^T}_{\sim \Sigma^{-1}} (X^T X + \lambda I)^{-1} X^T y$$

SVD

$$= \left[ \sum_{j=0}^{p-1} \left( \frac{\sigma_j^2}{\sigma_j^2 + \lambda} \right) u_j y_j^T \right]$$

$$\sigma_0 > \sigma_1 - \dots > \sigma_{p-1} > 0$$

with  $\Gamma_j$  small and  
 $\Gamma_i$ , we can suppose  
(reduce) the role of a  
specific component  
 $u_j \Rightarrow$   
reduction of degrees  
of freedom.

