

# FYS-STK3155/4155 lecture

## September 1, 2025

# FYS-STK3155/4155 lecture September 1, 2025

Analyzing Ordinary least squares  
and Ridge Regression.

$$X^T X = X X^T = \underline{1} = \begin{bmatrix} 1 & & \\ & 1 & 0 \\ 0 & & \ddots \\ & & & 1 \end{bmatrix}$$

$$\begin{aligned} \hat{y} &= X \underbrace{(X^T X)^{-1}}_{\underline{1}} X^T y \\ &= y \end{aligned}$$

SVD

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

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

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

$$UU^T = U^T U =$$

$$U = \begin{bmatrix} \vec{u}_0 & \vec{u}_1 & \dots & \vec{u}_{n-1} \end{bmatrix}$$

$\vec{u}_i^T \vec{u}_j = \delta_{ij}$

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

$$V^T V = V V^T = \mathbb{I}$$

$$V = \begin{bmatrix} \vec{v}_0 & \vec{v}_1 & \dots & \vec{v}_p \end{bmatrix}$$

$$\vec{v}_i^T \vec{v}_j = \delta_{i,j}$$

$$\begin{matrix} \vec{v}_0 > \vec{v}_1 > \dots > \vec{v}_p > 0 \end{matrix} \quad \vec{v}_i = \begin{bmatrix} A_0 \\ A_1 \\ \vdots \\ A_p \\ 0 \end{bmatrix}$$

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

$$\sigma_0 = 2 \quad \sigma_1 = 1$$

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

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

$$X^T X = \underbrace{V \Sigma^T}_{p \times m} \underbrace{U^T U}_{m \times m} \underbrace{\Sigma V^T}_{m \times p}$$

$$= V \underbrace{\Sigma^T \Sigma}_{\begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_p^2 \end{bmatrix}} V^T$$

$$([V, \Sigma^T \Sigma] = 0 \Rightarrow V \Sigma^T \Sigma = \Sigma^T \Sigma V$$

$$\boxed{X^T X = (\Sigma^T \Sigma) V V^T} \cdot V$$

$$V^T V = \underline{\underline{I}}$$

$$(X^T X) V = \Sigma^T \Sigma V$$

$$V = \begin{bmatrix} \vec{v}_0 & \vec{v}_1 & \dots & \vec{v}_{p-1} \end{bmatrix}$$

$$= \underbrace{(X^T X)}_A \vec{v}_i = \vec{v}_i^2 \vec{v}_i$$

$$\vec{v}_i^2 = \lambda_i \vec{v}_i$$

$$(X^T X) \vec{\beta}_n = \vec{\lambda}_n^2 \vec{\beta}_n$$

$$\lambda_i \geq 0$$

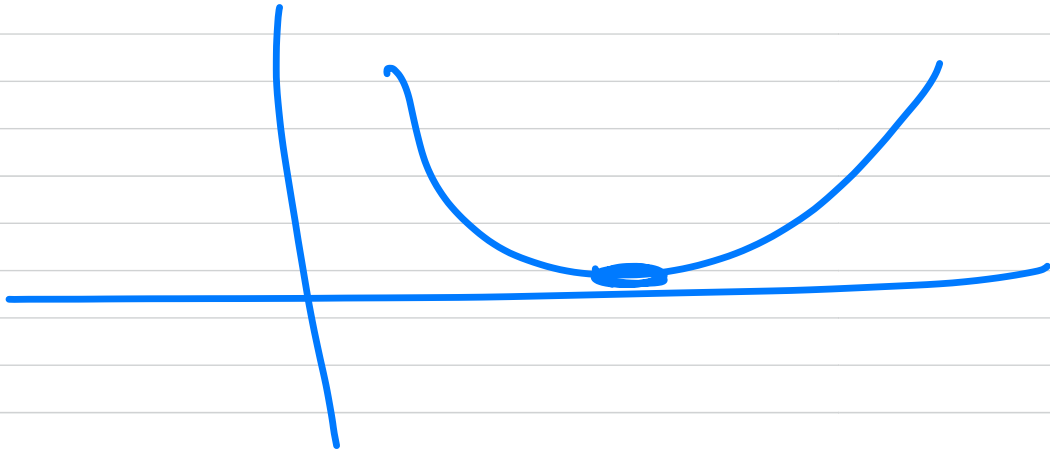
MSE

$$C(\theta) = \frac{1}{n} (\vec{X}\theta - y)^T (\vec{X}\theta - y)$$

$$\frac{\partial C}{\partial \theta} \quad \frac{\partial^2 C}{\partial \theta^2} = \frac{2}{n} X^T X$$



$$x^2 + 1 = f(x)$$

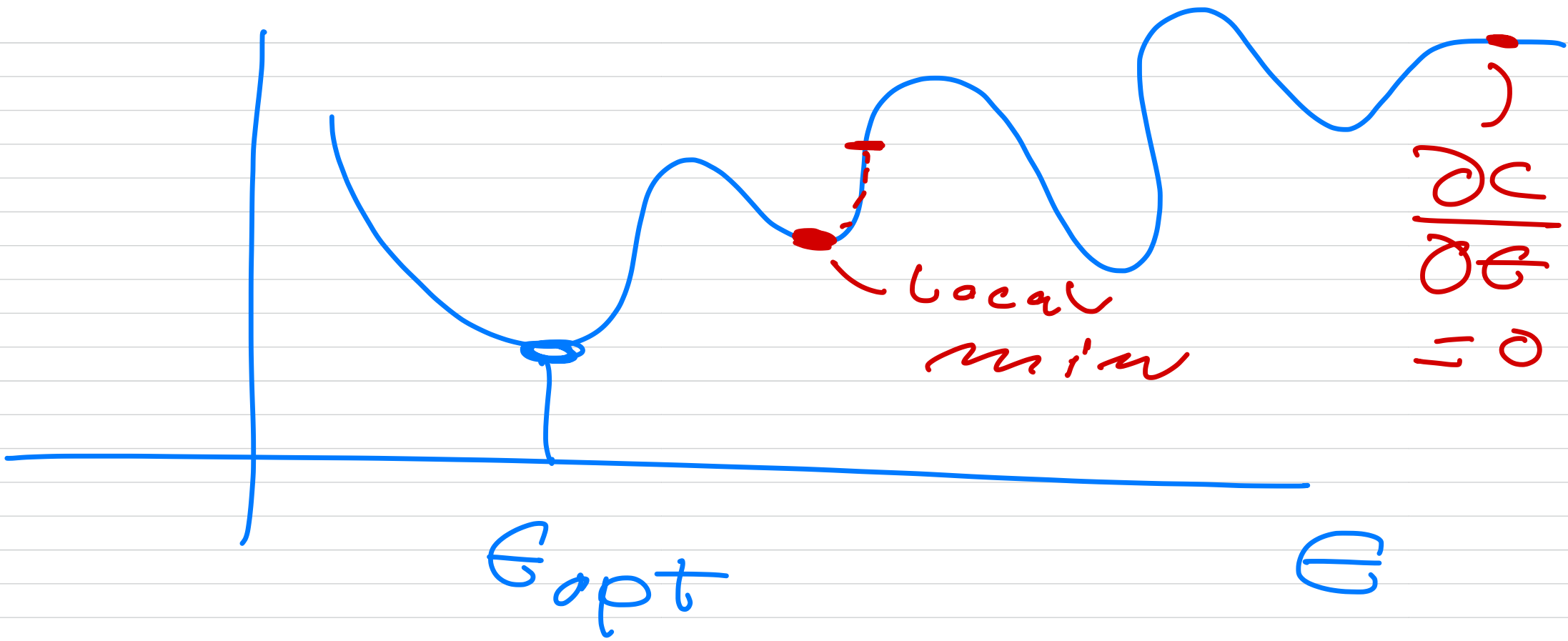


$$\frac{\partial^2 f}{\partial x^2} = 2 > 0$$

Ridge :

$$\frac{\partial^2 C_{\text{Ridge}}}{\partial \beta^2} = 2 \left( \frac{x^T x}{n} + \lambda \right)$$

$$\lambda \geq 0$$



SVD, OLS & Ridge

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

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

$$X (X^T X)^{-1} X^T y$$

$$= \left\{ U \Sigma V^T (V \Sigma^T \underbrace{A^T A}_{=I} \Sigma V^T)^{-1} \right\} X y$$
$$= U \Sigma V^T V \Sigma^T A^T A \Sigma V^T y$$

$$(V \Sigma^T \Sigma V^T)^{-1}$$

$A, B$  are square matrices  
and invertible

$$(AB)^{-1} = B^{-1}A^{-1}$$

$$\begin{aligned} \hat{y} &= (U \Sigma^T (\Sigma^T \Sigma)^{-1} \Sigma^T U^T) y \\ &= \left( \sum_{i=0}^{p-1} u_i u_i^T \right) y \end{aligned}$$

$$\hat{y}_{\text{Ridge}} = \left( \sum_{n=0}^{P-1} u_n u_n^T \frac{\sigma_n^2}{\sigma_n^2 + \lambda} \right) X^T Y$$

$E^T = [\theta_0, \theta_1, \dots, \theta_{p-1}]$

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

$$\sigma_i \geq 0$$

$$\lambda \rightarrow \infty \text{ (or large)}$$

Degrees of freedom through  $\hat{u}_i$  will be suppressed.

Taylor - expansion of  
 $C(E)$  around  $E^\uparrow$

$$E^\uparrow - E^{(n)}$$

$$\frac{dC}{dE} = g$$

at  
 $E = E^{(n)}$

$$C(E^\uparrow) = C(E^{(n)}) + g(E^{(n)}) (E^\uparrow - E^{(n)}) + \frac{1}{2} \frac{d^2 C}{dE^2} (E^\uparrow - E^{(n)})^2 + \dots$$

Truncate at  $\frac{d^2 C}{d\theta^2} \big|_{\theta = \theta^{(n)}}$

$$C(\theta^1) \approx C(\theta^{(n)}) +$$

$$\begin{aligned} & g(\theta^{(n)}) \left( \frac{dC}{d\theta} \bigg|_{\theta = \theta^{(n)}} (\theta^1 - \theta^{(n)}) \right. \\ & + \frac{1}{2} \left( \frac{d^2 C}{d\theta^2} \bigg|_{\theta = \theta^{(n)}} (\theta^1 - \theta^{(n)})^2 \right. \\ & \left. \left. H(\theta^{(n)}) \right) \right] \end{aligned}$$

$$= c(e^{(n)}) + g^{(n)} b^{(n)} + \frac{1}{2} (b^{(n)})^2 H^{(n)}$$

$$\frac{dc}{db} = g^{(n)} + b^{(n)} H^{(n)}$$

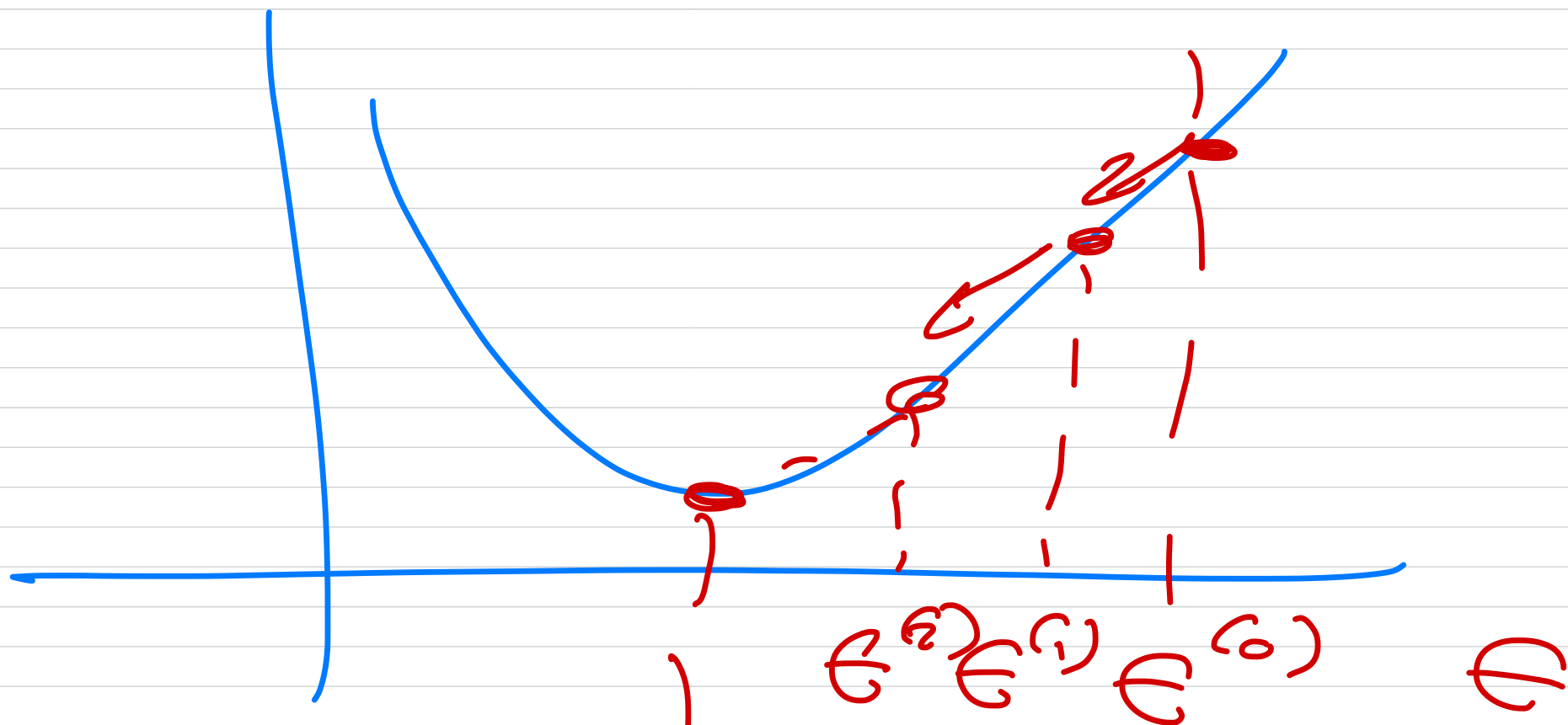
$$b^{(n)} = 0 \Rightarrow e^1 - e^{(n)} = (H^{(n)})^{-1} g^{(n)}$$



$$E^{\wedge} = E^{(n)} - (H^{(n)})^{-1} g^{(n)}$$

$$E^{\wedge} \rightarrow E^{(n+1)} = E^{(n)} - (H^{(n)})^{-1} g^{(n)}$$

recipe : start iteration  
with a guess  $E^{(0)}$   
keep iterating til  
 $|E^{(n+1)} - E^{(n)}| \leq \epsilon$



$$E^{(n+1)} = \Theta^{(n+1)} - (H^{(n)})^{-1} g^{(n)}$$

OLS

$$\nabla_{\theta} C(\theta) = \frac{2}{n} (X^T X \theta - X^T y)$$

$$\nabla_{\theta}^2 C(\theta) = \frac{2}{n} X^T X = H$$

$\begin{matrix} \uparrow \\ \text{vector} \end{matrix} \theta^{(n+1)} = \theta^{(n)} - \left( H^{-1}(\theta^{(n)}) \right) \nabla_{\theta} C(\theta^{(n)})$

$\begin{matrix} \uparrow \\ \text{independent} \\ \text{of } \theta \end{matrix}$

$H =$  Hessian matrix

$$H = \begin{bmatrix} \frac{\partial^2 C}{\partial \theta_0^2} & \frac{\partial^2 C}{\partial \theta_0 \partial \theta_1} & \dots & \frac{\partial^2 C}{\partial \theta_0 \partial \theta_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 C}{\partial \theta_{p-1} \partial \theta_0} & \dots & \dots & \frac{\partial^2 C}{\partial \theta_{p-1}^2} \end{bmatrix}$$

$$\theta^{(n+1)} = \theta^{(n)} - \left( H(\theta^{(n)}) \right)^{-1} \nabla_{\theta} C(\theta^{(n)})$$

$$E^{(n+1)} = E^{(n)} - \gamma \nabla_E (E^{(n)})$$

learning rate

constant  $\gamma$ : GD

$\gamma^{(n)}$ : { momentum  
ADAM  
RMSprop  
Adagrad  
: