

Lecture September 11

Interpretation of OLS & Ridge

$$\hat{y}^{\text{OLS}} = \hat{X}\hat{\beta} = \underbrace{X(X^T X)^{-1} X^T}_A Y = AY$$

$$A^2 = (X(X^T X)^{-1} X^T)^2 = X(X^T X)^{-1} X \\ = A \in \mathbb{R}^{n \times n}$$

\hat{y}^{OLS} is an orthogonal projection of y onto the space spanned by the columns of X

With $\hat{\beta}$ we have the residuals

$$\hat{\epsilon} = y - \hat{y}^{\text{OLS}} = y - X\hat{\beta} \\ = (I - X(X^T X)^{-1} X^T) y$$

The residuals $I = A + Q$
are the projections $Q = I - A$
of y onto the orthogonal complement of the space spanned by the columns

of X

SVD $X \in \mathbb{R}^{n \times p} (\mathbb{C}^{n \times p})$

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

$$U \in \mathbb{R}^{n \times n} \quad U^T U = U U^T = I$$
$$\Sigma \in \mathbb{R}^{n \times p}$$

$$\sigma_i \geq \sigma_{i+1} \leftarrow \sigma_1 \geq \sigma_2 \geq \sigma_3 \dots \geq \sigma_p \geq 0$$

$$V \in \mathbb{R}^{p \times p} \quad V V^T = V^T V = I$$

Example

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

$$\Sigma \Sigma^T = \begin{matrix} 3 \times 2 \\ \begin{bmatrix} 2 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \end{matrix} \begin{matrix} 2 \times 3 \\ \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \end{matrix}$$

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

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

standard case $n > p$
 $(n \gg p)$

$$\hat{y}^{OLS} = \underline{\underline{x}} \hat{\beta} = x(x^T x)^{-1} x^T y$$

$$= \underbrace{u \Sigma v^T}_{\downarrow} \left(v \Sigma^T \underbrace{u^T u}_{\stackrel{=I}{\Sigma}} \Sigma v^T \right)^{-1}$$

$$\qquad \qquad \qquad \times v \Sigma^T u^T y$$

$$v \Sigma^T v^T = \Sigma^2 \in \mathbb{R}^{P \times P}$$

$$= u \Sigma v^T \frac{I}{\Sigma^2} v \Sigma^T u^T y$$

$$= u \Sigma \frac{1}{\Sigma^2} \Sigma^T u^T y$$

$$u \in \mathbb{R}^{n \times n} \sim \mathbb{R}^{P \times P}$$

$$\Sigma \Sigma^T \in \mathbb{R}^{n \times n}$$

$$u \frac{\Sigma \Sigma^T}{\Sigma^2} u^T = \boxed{uu^T - I}$$

$$= \sum_{j=0}^{p-1} u_j u_j^T$$

$$\Rightarrow \hat{y}^{OLS} = \sum_{j=0}^{p-1} u_j u_j^T y$$

$\downarrow = 0$

$u^T g$ are the coordinates
of g wrt to the
orthonormal u

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

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

$$= \begin{bmatrix} u_{00} & u_{01} & \dots & u_{0m-1} \\ u_{10} & u_{11} & \dots & u_{1m-1} \\ u_{20} & u_{21} & \dots & u_{2m-1} \\ \vdots & \vdots & \ddots & \vdots \\ u_{m-10} & u_{m-11} & \dots & u_{m-1m-1} \end{bmatrix}$$

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

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

$$U = \begin{bmatrix} P & n-p \\ \vdots & \vdots \\ P & 0 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} P & 0 \\ 0 & 0 \end{bmatrix}$$

$$V^T = \begin{bmatrix} P \\ P \end{bmatrix}$$

$$X^T X \propto \text{cov}[x]$$

$$\text{var}(\hat{\beta}) \propto (X^T X)^{-1}$$

$$X^T X = V \Sigma^T U^T \alpha \Sigma U = V \Sigma^2 V^T$$

$(\mathbf{x}^T \mathbf{x}) \mathbf{v} = \mathbf{v} \Sigma^2 \Rightarrow$ eigenvectors
of $\mathbf{x}^T \mathbf{x}$ are \mathbf{v} with eigen
values Σ^2

$$\mathbf{x} \mathbf{x}^T = \mathbf{u} \Sigma \mathbf{v}^T \mathbf{v} \Sigma^T \mathbf{u}^T$$

$$= \mathbf{u} (\Sigma \Sigma^T) \mathbf{u}^T \Rightarrow \\ (\mathbf{x} \mathbf{x}^T) \mathbf{u} = \mathbf{u} (\Sigma \Sigma^T)$$

eigenvectors of $\mathbf{x} \mathbf{x}^T$ are
 \mathbf{u} with eigenvalues
 $\Sigma \Sigma^T$

$$\mathbf{u} = \text{evec}(\mathbf{x} \mathbf{x}^T)$$

$$\mathbf{v} = \text{evec}(\mathbf{x}^T \mathbf{x})$$

$$\Sigma^2 = \text{eval}(\mathbf{x} \mathbf{x}^T) = \text{eval}(\mathbf{x}^T \mathbf{x})$$

(non-zero)

Ridge

$$\hat{\beta}^{\text{Ridge}} = \underbrace{(\mathbf{x}^T \mathbf{x} + \lambda \mathbb{I})^{-1}}_{\in \mathbb{R}^{P \times P}} \mathbf{x}^T \mathbf{y}$$

$$\Rightarrow \hat{\mathbf{y}}^{\text{Ridge}} = \mathbf{x} \hat{\beta}^{\text{Ridge}}$$

$$\begin{aligned}
 &= u \Sigma v^T \left(v \Sigma^T \underbrace{u^T u}_{\Sigma} \Sigma v^T + \lambda \Sigma \right)^{-1} \\
 &\quad \times v \Sigma^T u^T y \\
 &= \left(\sum_{j=0}^{p-1} u_j \frac{\tau_j^2}{\tau_j^2 + \lambda} u_j^T \right) y
 \end{aligned}$$

$\lambda \geq 0$ then we have

$$\frac{\tau_j^2}{\tau_j^2 + \lambda} \leq 1$$

Ridge shrinks the coordinates of y by a factor $\frac{\tau_j^2}{\tau_j^2 + \lambda}$

A greater amount of shrinkage is applied to the coordinates of basis vector with smaller values of τ_j^2 .

Example

$$X^T X = X X^T = \mathbb{1} \quad p = n$$

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

$$\Rightarrow \hat{y}^{\text{Ridge}} = \frac{1}{1+\lambda} \hat{y}^{\text{OLS}}$$

$$X^T X = V \Sigma^2 V^T = C[x] \cdot n$$

$$(X^T X) v = \lambda v^2$$

eigenvector

of the covariance matrix
and the singular values

$$\frac{\lambda_j^2}{n} \quad \lambda_j \geq \lambda_{j+1}$$

The first component
 v_1 has the largest
singular value λ_1^2

With Ridge we are
effectively shrinking
those components with
small covariance / variance
(small eigenvalues λ_{ii}^2)

Σ_p has eigenvalue σ_p^2

Ridge shrinks these values by $\frac{\sigma_j^2}{\sigma_j^2 + \lambda}$

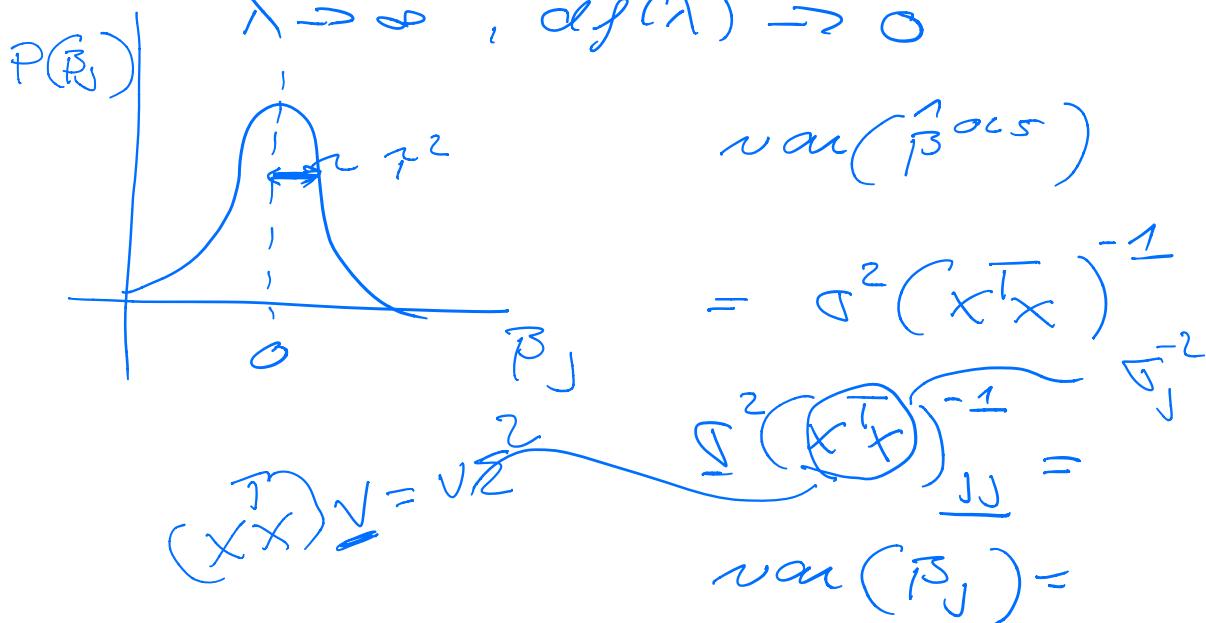
degrees of freedom

$$df(\lambda) = \text{Tr} \left[\bar{x} (\bar{x}^T \bar{x} + \lambda I)^{-1} \right]$$

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

$\lambda = 0$ then $df(\lambda) = p$

$\lambda \rightarrow \infty$, $df(\lambda) \rightarrow 0$



$$\text{var}(\hat{\beta}_{\text{Ridge}}) = \sigma^2 \left[\bar{x}^T \bar{x} + \lambda I \right]^{-1} \bar{x}^T \bar{x}$$

$$\times \left\{ E[x^T x + \lambda I]^{-1} \right\}^T$$

Ridge if $\text{var}(\hat{\beta}_j)$ is large or small $\hat{\sigma}_j^2$ from $x^T x$, with λ in Ridge we can shrink $\text{var}(\hat{\beta}_j^{\text{Ridge}})$ to smaller values,