

Svetlana Telnova 2 group

1. We have two absolutely identical preliminary standardized regressors x and \bar{x} . The dependent variable y is centered.

In the ridge regression one minimizes the loss function

$$\text{loss}(\hat{\beta}) = (y - \hat{y})^T(y - \hat{y}) + \lambda \hat{\beta}^T \hat{\beta}, \quad \hat{y} = \hat{\beta}_1 x + \hat{\beta}_2 \bar{x}.$$

- (a) Find the optimal $\hat{\beta}_1$ and $\hat{\beta}_2$ for fixed λ .
- (b) What happens to the estimates when $\lambda \rightarrow \infty$?
- (c) What happens to the sum $\hat{\beta}_1 + \hat{\beta}_2$ when $\lambda \rightarrow 0$?

$$\text{loss}(\hat{\beta}) = (y - \hat{y})^T(y - \hat{y}) + \lambda \hat{\beta}^T \hat{\beta} =$$

$$\min (y - \hat{y})^T(y - \hat{y}) + \lambda \hat{\beta}^T \hat{\beta} \quad \hat{y} = \hat{\beta}_1 x + \hat{\beta}_2 \bar{x}.$$

$$\frac{\partial \text{loss}(\hat{\beta})}{\partial \hat{\beta}} = \frac{\partial}{\partial \hat{\beta}} \lambda \hat{\beta}^T \hat{\beta} = 2 \lambda \hat{\beta}$$

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}$$

$$X \hat{\beta} = \begin{bmatrix} x_1 & \dots & x_n \\ x_2 & \dots & x_n \\ \vdots & \ddots & \vdots \\ x_n & \dots & x_n \end{bmatrix} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} = \begin{bmatrix} x_1 \hat{\beta}_1 + x_2 \hat{\beta}_2 \\ x_2 \hat{\beta}_1 + x_3 \hat{\beta}_2 \\ \vdots \\ x_n \hat{\beta}_1 + x_{n+1} \hat{\beta}_2 \end{bmatrix} \quad \text{which is equivalent}$$

$$\hat{y} = \hat{\beta}_1 x + \hat{\beta}_2 \bar{x} \rightarrow \text{can represent } \hat{y} = X \hat{\beta}.$$

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

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

$$= (y^T y - 2 \hat{\beta}^T X^T y + \hat{\beta}^T X^T X \hat{\beta})$$

$$\frac{\partial \text{loss}}{\partial \hat{\beta}} = \frac{\partial (y^T y - 2 \hat{\beta}^T X^T y + \hat{\beta}^T X^T X \hat{\beta} + \lambda \hat{\beta}^T \hat{\beta})}{\partial \hat{\beta}} = 0.$$

$$- 2 X^T y + 2 X^T X \hat{\beta} + 2 \lambda \hat{\beta} = 0.$$

$$(X^T X + 2\lambda) \hat{\beta} = X^T y$$

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

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

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = (X^T X + \lambda)^{-1} X^T y$$

$$\begin{aligned} \text{loss}(\hat{\beta}) &= (y - \hat{y})^T (y - \hat{y}) + \lambda \hat{\beta}^T \hat{\beta} = (y - (\hat{\beta}_1 + \hat{\beta}_2) x)^T (y - (\hat{\beta}_1 + \hat{\beta}_2) x) + \\ &+ \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2) = \left\| (y - (\hat{\beta}_1 + \hat{\beta}_2) x) \right\|^2 + \lambda (\hat{\beta}_1^2 + \hat{\beta}_2^2) \\ \left\{ \begin{array}{l} \frac{\partial \text{loss}}{\partial \beta_1} = 2(-x^T)(y - (\hat{\beta}_1 + \hat{\beta}_2)x) + 2\lambda \hat{\beta}_1 = 0 \\ \frac{\partial \text{loss}}{\partial \beta_2} = 2(-x^T)(y - (\hat{\beta}_1 + \hat{\beta}_2)x) + 2\lambda \hat{\beta}_2 = 0 \end{array} \right. \end{aligned}$$

The equations are symmetric $\Rightarrow \hat{\beta}^* = \hat{\beta}_1 = \hat{\beta}_2$.

$$-2x^T(y - 2\hat{\beta}^*x) + 2\lambda \hat{\beta}^* = 0$$

$$-2x^Ty + 4x^Tx\hat{\beta}^* + 2\lambda \hat{\beta}^* = 0.$$

$$(2x^Tx + \lambda) \hat{\beta}^* = x^Ty$$

$$\hat{\beta}^* = (2x^Tx + \lambda)^{-1}x^Ty = \hat{\beta}_1 = \hat{\beta}_2, \text{ - optimal } \hat{\beta}_1 \text{ and } \hat{\beta}_2 \text{ for fixed } \lambda.$$

b) $\lim_{\lambda \rightarrow \infty} (2x^Tx + \lambda)^{-1} \cdot x^Ty = 0 \Rightarrow (\hat{\beta}_1^*, \hat{\beta}_2^*) = (0, 0)$

c) $\lim_{\lambda \rightarrow 0} 2(2x^Tx + \lambda)^{-1} \cdot x^Ty = \hat{\beta}^* \Rightarrow \text{The sum of estimations converges to the optimal estimator } \hat{\beta}^*$

2. Consider the model $y = X\beta + u$ where β is non-random, $\mathbb{E}(u | X) = 0$, the matrix X of size $n \times k$ has rank $X = k$, but $\text{Var}(u | X) = \sigma^2 W$ with $W \neq I$. Let $\hat{\beta}$ be the standard OLS estimator of β .

- (a) Find $\mathbb{E}(\hat{\beta} | X), \mathbb{E}(\hat{\beta})$.
- (b) Find $\text{Var}(\hat{\beta} | X)$.
- (c) How do you think, will the standard confidence interval for β be valid in this case?
- (d) Find $\text{Cov}(y, \hat{\beta} | X)$.

$$y = X\beta + u \quad \mathbb{E}(u | X) = 0. \quad \begin{matrix} X \\ n \times k \end{matrix}$$

$$\text{Var}(u | X) = \sigma^2 W \quad W \neq I$$

$$\mathbb{E}(X\beta | X)$$

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

$$\mathbb{E}(\hat{\beta} | X) = \mathbb{E}((X^T X)^{-1} X^T y | X) = (X^T X)^{-1} X^T \mathbb{E}(X\beta + u | X) =$$

$$= (X^T X)^{-1} X^T (X^T X) \underbrace{\beta + 0}_{\beta} = \beta$$

$$\mathbb{E}(\hat{\beta}) = \mathbb{E}((X^T X)^{-1} X^T y) = \mathbb{E}((X^T X)^{-1} X^T (X\beta + u)) = (X^T X)^{-1} X^T \mathbb{E}(X\beta + u) =$$

$$= \underbrace{X^T X^{-1} X^T X}_{\text{const}} \beta + \mathbb{E}(u) = \beta + \mathbb{E}(u).$$

$$B) \text{Var}(\hat{\beta} | X) = \text{Var}((X^T X)^{-1} X^T y | X) = (X^T X)^{-1} X^T \text{Var}(u | X) \cdot \underbrace{[(X^T X)^{-1} X^T]^T}_{\text{const}}$$

$$= (X^T X)^{-1} X^T \sigma^2 W X \cdot (X^T X)^{-1} = \sigma^2 (X^T X)^{-1} X^T W X \cdot (X^T X)^{-1}$$

$$\textcircled{1} \quad \text{Var}(y | X) = \text{Var}(X\beta + u | X) = \text{Var}(u | X) \cdot X^T = \sigma^2 W$$

c) No, the CI will be invalid since it is based on assumption that $\text{Var}(u | X) = \sigma^2 I$. In our case $\text{Var}(u | X) = \sigma^2 W$ ($W \neq I$) indicates heteroscedasticity or correlation among errors. Thus, Gauss Markov conditions are violated, the standard estimators and CI are invalid.

$$d) \text{Cov}(y; \hat{\beta} | X) = \text{Cov}(y; (X^T X)^{-1} X^T y | X) = \text{Cov}(y, y | X) [(X^T X)^{-1} X^T]^T =$$

$$= \text{Var}(y | X) \cdot X^T (X^T X)^{-1} = \sigma^2 W \cdot X^T (X^T X)^{-1}$$

$$\text{Var}(y | X) = \text{Var}(X\beta + u | X) = \text{Var}(u | X) = \sigma^2 W$$

3. Consider the matrix

$$X = \begin{pmatrix} 2 & 1 \\ -1 & 2 \\ 1 & 1 \end{pmatrix}.$$

- (a) Find the matrix $X^T X$ and diagonalize it.
- (b) Find the SVD of X .
- (c) Find the best approximation to X with rank equal to 1.

Remark: in the principal component analysis the variables in the matrix X should be standardized. If you can't do this by bare hands, feel free to use python, but provide code!

a) $X^T X = \begin{pmatrix} 2 & -1 & 1 \\ 1 & 2 & 1 \end{pmatrix} \times \begin{pmatrix} 2 & 1 \\ -1 & 2 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 6 & 1 \\ 1 & 6 \end{pmatrix}$

$$\begin{vmatrix} 6-\lambda & 1 \\ 1 & 6-\lambda \end{vmatrix} = 0$$

$$(6-\lambda)^2 - 1 = 0$$

$$(6-\lambda-1)(6-\lambda+1) = 0$$

$$\lambda = 5 \quad \lambda = 7.$$

① $\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$. $x_1 + x_2 = 0$ $\begin{pmatrix} x_1 \\ -x_1 \end{pmatrix} = x_1 \begin{pmatrix} 1 \\ -1 \end{pmatrix}$

② $\begin{pmatrix} -1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$. $-x_1 + x_2 = 0$ $x_1 = x_2$ $\begin{pmatrix} x_1 \\ x_1 \end{pmatrix} = x_1 \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$$A = \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$

$$A^{-1} = \left(\begin{array}{cc|cc} 1 & 1 & 1 & 0 \\ -1 & 1 & 0 & 1 \end{array} \right) \rightarrow \left(\begin{array}{cc|cc} 1 & 1 & 1 & 0 \\ 0 & 2 & 1 & 1 \end{array} \right) \rightarrow \left(\begin{array}{cc|cc} 2 & 0 & 1 & -1 \\ 0 & 2 & 1 & 1 \end{array} \right) \rightarrow \left(\begin{array}{cc|cc} 1 & 0 & \frac{1}{2} & -\frac{1}{2} \\ 0 & 1 & \frac{1}{2} & \frac{1}{2} \end{array} \right)$$

$$A \cdot A^{-1} = \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \times \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$D = \begin{pmatrix} 7 & 0 \\ 0 & 5 \end{pmatrix} \quad \begin{pmatrix} 2 & 1 \\ -1 & 2 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \times \begin{pmatrix} 7 & 0 \\ 0 & 5 \end{pmatrix} \times \begin{pmatrix} \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$$

Rectangular diagonal matrix:

$$B) \quad X \cdot X^T = \begin{pmatrix} \sqrt{7} & 0 \\ 0 & \sqrt{5} \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 2 & 1 \\ -1 & 2 \\ 1 & 1 \end{pmatrix} \times \begin{pmatrix} 2 & -1 & 1 \\ 1 & 2 & 1 \end{pmatrix} = \begin{pmatrix} 5 & 0 & 3 \\ 0 & 5 & 1 \\ 3 & 1 & 2 \end{pmatrix}$$

$$\begin{vmatrix} 5-\lambda & 0 & 3 \\ 0 & 5-\lambda & 1 \\ 3 & 1 & 2-\lambda \end{vmatrix} = 0$$

$$(5-\lambda)^2(2-\lambda) + 0 + 0 - (9(5-\lambda) + (5-\lambda)) = 0.$$

$$(5-\lambda)^2(2-\lambda) - 10(5-\lambda) = 0.$$

$$(5-\lambda)((5-\lambda)(2-\lambda) - 10) = 0.$$

$$10 - 2\lambda - 5\lambda + \lambda^2 - 10 = 0$$

$$① \quad \begin{pmatrix} 5 & 0 & 3 \\ 0 & 5 & 1 \\ 3 & 1 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \quad \lambda = 5, \quad \lambda^2 - 7\lambda = 0$$

$$x_1 = -\frac{3x_3}{5}, \quad \lambda = 0, \quad \lambda = 7,$$

$$5x_1 + 3x_3 = 0.$$

$$x_1 = -\frac{3x_3}{5}$$

$$v_3 = \begin{pmatrix} -\frac{3}{5} \\ -\frac{1}{5} \\ 1 \end{pmatrix}$$

$$5x_2 + x_3 = 0 \Rightarrow x_2 = -\frac{x_3}{5}$$

$$② \quad \begin{pmatrix} 0 & 0 & 3 \\ 0 & 0 & 1 \\ 3 & 1 & -3 \\ -2 & 0 & 3 \\ 0 & -2 & 1 \\ 3 & 1 & -5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$3x_3 = 0 \Rightarrow x_3 = 0, \quad v_2 = \begin{pmatrix} -\frac{1}{3} \\ 1 \\ 0 \end{pmatrix}$$

$$3x_1 + x_2 - 3x_3 = 0, \quad 3x_1 + x_2 = 0, \quad x_1 = -\frac{x_2}{3}$$

$$-2x_1 + 3x_3 = 0, \quad x_1 = \frac{3x_3}{2}, \quad v_1 = \begin{pmatrix} \frac{3}{2} \\ \frac{1}{2} \\ 1 \end{pmatrix} \rightarrow \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}$$

$$-2x_2 + x_3 = 0, \quad x_2 = \frac{x_3}{2}, \quad v_3 = \begin{pmatrix} \frac{3}{\sqrt{14}} \\ \frac{1}{\sqrt{14}} \\ \frac{1}{2} \end{pmatrix}$$

Normalizing vectors:

$$|v_1| = \sqrt{\frac{1}{9} + 1} = \frac{\sqrt{10}}{3}$$

$$|v_3| = \sqrt{\frac{9}{25} + \frac{1}{25} + 1} = \frac{\sqrt{35}}{5}$$

$$v_1' = \begin{pmatrix} \frac{3}{\sqrt{14}} \\ \frac{1}{\sqrt{14}} \\ \frac{1}{2} \end{pmatrix}$$

$$v_2 = \begin{pmatrix} -\frac{1}{\sqrt{10}} \\ \frac{3}{\sqrt{10}} \\ 0 \end{pmatrix}$$

$$v_1 = \begin{pmatrix} -\frac{3}{\sqrt{35}} \\ -\frac{1}{\sqrt{35}} \\ \frac{5}{\sqrt{35}} \end{pmatrix}$$

Normalizing : $\begin{pmatrix} 1 & 1 \\ -1 & 1 \\ \frac{\sqrt{2}}{2} & \frac{\sqrt{2}}{2} \end{pmatrix} \rightarrow \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}$
 length length

$$X = \begin{pmatrix} \frac{3}{\sqrt{4}} & -\frac{1}{\sqrt{10}} & -\frac{3}{\sqrt{35}} \\ \frac{1}{\sqrt{4}} & \frac{3}{\sqrt{10}} & -\frac{1}{\sqrt{35}} \\ 0 & 0 & \frac{5}{\sqrt{35}} \end{pmatrix}_{[3 \times 3]} \times \begin{pmatrix} \sqrt{7} & 0 \\ 0 & \sqrt{5} \\ 0 & 0 \end{pmatrix}_{n \times m} \times \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}_{2 \times 2} = \begin{pmatrix} 2 & 1 \\ -1 & 2 \\ 1 & 1 \end{pmatrix} \Rightarrow$$

SVD decomposition

$$\begin{pmatrix} \frac{3}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{3}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ 0 & 0 \end{pmatrix} \times \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix} = \begin{pmatrix} 2 & 1 \\ -1 & 2 \\ 1 & 1 \end{pmatrix}$$

c) $A_k = U_k S_k V_k^T$

$$A = [A_{k \times 1}] \times [n \times k] \times [k \times k] \times [k \times m]$$

U_k - keep the first k columns

S_k - keep only the top k singular val

V_k^T - keep the first k rows

$$\begin{pmatrix} \frac{3}{\sqrt{4}} \\ \frac{1}{\sqrt{4}} \\ \frac{2}{\sqrt{4}} \\ \frac{1}{\sqrt{4}} \end{pmatrix}_{[3 \times 1]} \times \begin{pmatrix} \sqrt{7} \\ 1 \end{pmatrix}_{[1 \times 2]} \times \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}_{[1 \times 2]} = \begin{pmatrix} \frac{3}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \\ \frac{2}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{pmatrix}_{[3 \times 1]} \times \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{pmatrix}_{[1 \times 2]} = \begin{pmatrix} \frac{3}{2} & \frac{3}{2} \\ \frac{1}{2} & \frac{1}{2} \\ 1 & 1 \end{pmatrix}_{[1 \times 2]}$$

\Rightarrow the matrix approximation of X with rank 1 is $(1; 1)$

4. The columns of X are standardized. You know the SVD of the matrix $X = UDV^T$. The diagonal elements of D are positive and ordered from highest to lowest, $d_{11} > d_{22} > \dots > 0$.

Let's maximize $\|Xw\|^2$ by choosing an optimal vector w subject to $\|w\|^2 = 1$.

- (d) Write the Lagrangian function for this problem.
- (e) Find the first order conditions. Differential is your friend!
- (f) Find the optimal w in terms of columns of V .

Hint: one may interpret the FOC in terms of eigenvalues and eigenvectors!

$$d) \begin{cases} \|Xw\|^2 \rightarrow \max_w \\ \text{s.t. } \|w\|^2 = 1 \end{cases}$$

$$\begin{aligned} \|Xw\|^2 &= (Xw)^T (Xw) = w^T X^T X w \\ \|w\|^2 &= w^T w \end{aligned} \Rightarrow \begin{cases} \max_w w^T X^T X w \\ \text{s.t. } w^T w = 1. \end{cases}$$

$$L = \|Xw\|^2 - \lambda (\|w\|^2 - 1) = w^T X^T X w - \lambda (w^T w - 1)$$

e) Denote $X^T X$ as A $(w_1 \ w_2 \dots \ w_n) \mid$

$$\begin{aligned} \text{F.O.C. } \begin{cases} \frac{\partial L}{\partial w} = 2 X^T X w - 2\lambda w = 0 \Rightarrow X^T X \underbrace{w}_{\substack{\text{eigenvector} \\ \text{matrix}}} = \underbrace{\lambda w}_{\substack{\text{eigenvalue}}} \quad (1) \\ \frac{\partial L}{\partial \lambda} = -w^T w + 1 = 0 \Rightarrow w^T w = 1 \quad (2) \end{cases} \end{aligned}$$

f) From (1) w is an eigenvector of $X^T X \Rightarrow$

$$\|Xw\|^2 = w^T X^T X w = w^T \lambda w = \lambda w^T w = \lambda \rightarrow \max_w \Rightarrow$$

the eigenvector of w should correspond to the maximum eigenvalue Γ is the normalized matrix of eigenvectors of $X^T X$ matrix. Thus, the optimal value of w corresponds to the maximum eigenvalue of matrix Γ