

## OLS

We assume:  $y = x\beta + \varepsilon$  with  $E(\varepsilon) = 0$ ,  $\text{Var}(\varepsilon) = \sigma^2 I$

OLS estimator:  $\beta^* = (x^T x)^{-1} x^T y = H y$

$$\tilde{\beta} = Cy = (H+D)y$$

q1).  $E(\tilde{\beta}) = E[(H+D)y] = (H+D)E(y)$  and  $E(y) = x\beta + \underbrace{E(\varepsilon)}_0$

$$\text{So } E(\tilde{\beta}) = (H+D)x\beta$$

Since  $\tilde{\beta}$  is unbiased:

$$E(\tilde{\beta}) = \beta \quad \text{i.e.} \quad \underbrace{Hx\beta + Dx\beta}_{\substack{= \\ I\beta}} = \beta$$

$$\text{i.e. } Dx\beta = 0$$

$$\text{i.e. } Dx = 0$$

We use the assumption that  $x \in \text{Ker}(D)$

• We have  $\text{Var}(\beta_{OLS}^*) = \sigma^2 H H^T$

Let's compute  $\text{Var}(\tilde{\beta})$

$$\text{Var}(\tilde{\beta}) = \text{Var}(Cy)$$

$$= \text{Var}(x\beta + \varepsilon)$$

$$= C^T \text{Var}(x\beta + \varepsilon) C$$

$$= \text{Var}(\varepsilon) = \sigma^2 I$$

$$= \sigma^2 (H H^T + \underbrace{H D^T + D H^T}_{=0} + D D^T)$$

$\underbrace{H D^T + D H^T}_{=0}$  since  $x \in \text{Ker}(D)$ ,  $x^T \in \text{Ker}(D^T)$

$$\text{Var}(\tilde{\beta}) = \underbrace{\sigma^2 H H^T}_{\text{Var}(\beta^*)} + \underbrace{\sigma^2 D D^T}_{>0 \text{ since } D \text{ non-zero}} > \text{Var}(\beta^*)$$



## Ridge regression

$$\beta_{\text{ridge}}^* = \underset{\beta}{\operatorname{argmin}} \underbrace{(y_c - x_c \beta)^T (y_c - x_c \beta)}_{f(\beta)} + \lambda \|\beta\|_2^2$$

q2) We have the first order optimality condition

$$\frac{\partial f(\beta_{\text{ridge}}^*)}{\partial \beta} = -2x_c^T (y - x\beta) + 2\lambda \beta = 0$$

i.e

$$\beta_{\text{ridge}}^* = (x^T x + \lambda I)^{-1} x^T y$$

Let's compute the bias:

$$E(\beta_{\text{ridge}}^*) = E[(\lambda I + x^T x)^{-1} x^T x \beta] + E[(\lambda I + x^T x)^{-1} x^T \varepsilon]$$

$$E(\beta_{\text{ridge}}^*) = (\lambda I + x^T x)^{-1} x^T x \beta = 0$$

$$\neq \beta \quad \text{for } \lambda \neq 0$$

The estimator of ridge regression is biased.

## SVD decomposition

q3) We have  $x_c = UDV^T$

$$\beta_{\text{ridge}}^* = (x^T x + \lambda I)^{-1} x^T y$$

$$= (VD^T DV^T + \lambda VV^T)^{-1} VD^T U^T y$$

$$= V (D^2 + \lambda)^{-1} \underbrace{V^T V}_{I} D^T U^T y$$

$$= V (D^2 + \lambda)^{-1} D U^T y$$

since  $y = x\beta = UDV^T \beta$ :

$$\beta_{\text{ridge}}^* = V (D^2 + \lambda)^{-1} D^T D V^T \beta$$

$$\beta_{\text{ridge}}^* = V \operatorname{diag} \left( \frac{d_i^2}{d_i^2 + \lambda} \right) V^T \beta$$



## Variance

q4) We have:

$$\begin{aligned}\text{Var}(\beta^*_{\text{ridge}}) &= E[(\beta^*_{\text{ridge}} - E[\beta^*_{\text{ridge}}])(\beta^*_{\text{ridge}} - E[\beta^*_{\text{ridge}}])^T] \\ &= E\left[\left((x^T x + \lambda I)^{-1} x^T \varepsilon\right) \left((x^T x + \lambda I)^{-1} x^T \varepsilon\right)^T\right] \\ &= (x^T x + \lambda I)^{-1} x^T \underbrace{E[\varepsilon \varepsilon^T]}_{\sigma^2 I} x (x^T x + \lambda I)^{-1}\end{aligned}$$

$$\text{Var}(\beta^*_{\text{ridge}}) = \sigma^2 V (D^T D + \lambda I)^{-1} D^T D (D^T D + \lambda I)^{-1} V^T$$

$$\text{Var}(\beta^*_{\text{ridge}}) = \sigma^2 V \text{diag}\left(\frac{d_i^2}{(d_i^2 + \lambda)^2}\right) V^T$$

$$\text{And } \text{Var}(\beta^*_{\text{OLS}}) = \sigma^2 (x^T x)^{-1} = \sigma^2 V \text{diag}\left(\frac{1}{d_i^2}\right) V^T$$

$$\begin{aligned}\text{Thus } \underbrace{\text{Var}(\beta^*_{\text{ridge}}) - \text{Var}(\beta^*_{\text{OLS}})}_{< 0} &= \sigma^2 V \text{diag}\left(\frac{d_i^2}{(d_i^2 + \lambda)^2} - \frac{1}{d_i^2}\right) V^T \\ &= -\frac{2\lambda d_i^2 + \lambda^2}{d_i^2(d_i^2 + \lambda)} < 0\end{aligned}$$

And so:

$$\text{Var}(\beta^*_{\text{OLS}}) > \text{Var}(\beta^*_{\text{ridge}})$$

q5) Using what we found for the variance and bias for the ridge:

$$E(\beta^*_{\text{ridge}}) = V^T \text{diag}\left(\frac{d_i^2}{d_i^2 + \lambda}\right) V^T \beta \xrightarrow{\lambda \rightarrow \infty} 0$$

$$\text{Var}(\beta^*_{\text{ridge}}) = \sigma^2 V \text{diag}\left(\frac{d_i^2}{(d_i^2 + \lambda)^2}\right) V^T \xrightarrow{\lambda \rightarrow \infty} 0$$

q6) IF  $x_c^T x_c = I_d$ :  $(x^T x)^{-1} = I_d$

$$\beta^*_{\text{ridge}} = \underbrace{(x^T x + \lambda I)}_{I}^{-1} x^T y = \frac{1}{1+\lambda} \underbrace{I}_{(x^T x)^{-1}} x^T y = \frac{1}{1+\lambda} (x^T x)^{-1} x^T y$$

$$\beta^*_{\text{ridge}} = \frac{1}{1+\lambda} \beta^*_{\text{OLS}} y$$



## Elastic Net

$$\beta_{\text{Elastic}}^* = \underset{\beta}{\operatorname{argmin}} \underbrace{(y_c - x_c \beta)^T (y_c - x_c \beta)}_{f''(\beta)} + \lambda_2 \| \beta \|_2^2 + \lambda_1 \| \beta \|_1$$

q7) We have:

$$\frac{\partial f(\beta_{\text{Elastic}}^*)}{\partial \beta} = -2x_c^T (y - x \beta) + 2\lambda_2 \beta + \lambda_1 (1 - \alpha) (\pm 1) = 0$$

$$2\beta (\underbrace{x_c^T x_c}_{\lambda_2} + \underbrace{\lambda_2}_{\lambda_2}) = 2 \underbrace{x_c^T y}_{\lambda_1} \pm \underbrace{\lambda_1 (1 - \alpha)}_{\lambda_1}$$

$$(x^T x)^{-1} x^T y = \beta_{\text{OLS}}^* \quad \text{since } (x^T x)^{-1} = \text{Id}$$

$$\beta_{\text{Elastic}}^* = \frac{\beta_{\text{OLS}}^* \pm \frac{\lambda_1}{2}}{1 + \lambda_2}$$