

2. Bias and variance of ridge regression.

Ridge regression:

$$\hat{\beta}_\tau = \arg \min_{\beta} \|X\beta - y\|_2^2 + \tau \|\beta\|_2^2, \quad \tau \geq 0$$

True model $y = X\beta^* + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$.

Prove that

$$E[\hat{\beta}_\tau] = S_\tau^{-1} S \beta^*, \quad \text{Cov}[\hat{\beta}_\tau] = S_\tau^{-1} S S_\tau^{-1} \sigma^2$$

where

$$S = X^T X, \quad S_\tau = X^T X + \tau \mathbb{I}.$$

By taking the derivative of the loss function with respect to β , one can obtain

$$\frac{\partial \text{Loss}}{\partial \beta} = 2 [X^T (y - X\beta) + \tau \beta] \stackrel{!}{=} 0,$$

for which we can solve for the zero point to in turn obtain the optimal $\hat{\beta}_\tau$ as per

$$\hat{\beta}_\tau = (X^T X + \tau \mathbb{I})^{-1} X^T y.$$

Then we have the expectation value as

$$\begin{aligned} E[\hat{\beta}_\tau] &= E[(X^T X + \tau \mathbb{I})^{-1} X^T y] \quad y = X\beta^* + \varepsilon \\ &= E[(X^T X + \tau \mathbb{I})^{-1} X^T (X\beta^* + \varepsilon)] \\ &= E[(X^T X + \tau \mathbb{I})^{-1} X^T X \beta^*] + E[(X^T X + \tau \mathbb{I})^{-1} X^T \varepsilon] \\ &= E[S_\tau^{-1} S \beta^*] + E[S_\tau^{-1} X^T] E[\varepsilon] \\ &= S_\tau^{-1} S E[\beta^*] \quad \quad \quad = 0 \text{ since } \varepsilon \sim N(0, \sigma^2) \\ &= S_\tau^{-1} S \beta^* \end{aligned}$$

And then the variance becomes

$$\begin{aligned}
 \text{Cov}[\hat{\beta}_c] &= E[(\hat{\beta}_c - E[\hat{\beta}_c])^2] \\
 &= E[(\hat{\beta}_c - E[\hat{\beta}_c])(\hat{\beta}_c - E[\hat{\beta}_c])^T] \\
 &= E[\hat{\beta}_c \hat{\beta}_c^T - \hat{\beta}_c E[\hat{\beta}_c]^T - E[\hat{\beta}_c] \hat{\beta}_c^T + E[\hat{\beta}_c] E[\hat{\beta}_c]^T] \\
 &= E[\hat{\beta}_c \hat{\beta}_c^T] - E[\hat{\beta}_c] E[\hat{\beta}_c]^T - \cancel{E[\hat{\beta}_c] E[\hat{\beta}_c]^T} + \cancel{E[\hat{\beta}_c] E[\hat{\beta}_c]^T} \\
 &= E[(S_c^{-1} X^T y)(S_c^{-1} X^T y)^T] - (S_c^{-1} S \beta^*)(S_c^{-1} S \beta^*)^T \\
 &= E[(S_c^{-1} X^T (X \beta^* + \varepsilon))(S_c^{-1} X^T (X \beta^* + \varepsilon))^T] - (S_c^{-1} S \beta^*)(S_c^{-1} S \beta^*)^T
 \end{aligned}$$

$$\begin{aligned}
 & \left[(S_c^{-1} X^T (X \beta^* + \varepsilon))(S_c^{-1} X^T (X \beta^* + \varepsilon))^T = \right. \\
 & \quad = (S_c^{-1} X^T X \beta^*)(S_c^{-1} X^T X \beta^*)^T + S_c^{-1} X^T \varepsilon (S_c^{-1} X^T X \beta^*)^T \\
 & \quad \quad + S_c^{-1} X^T X \beta^* (S_c^{-1} X^T \varepsilon)^T + S_c^{-1} X^T \varepsilon (S_c^{-1} X^T \varepsilon)^T \left. \right] \\
 & = E[(S_c^{-1} S \beta^*)(S_c^{-1} S \beta^*)^T] + E[S_c^{-1} X^T] \underbrace{E[\varepsilon]}_{=0} E[S_c^{-1} S \beta^*]^T \\
 & \quad + \underbrace{E[S_c^{-1} S \beta^*]}_{=0} E[\varepsilon]^T E[S_c^{-1} X^T] + E[S_c^{-1} X^T \varepsilon (S_c^{-1} X^T \varepsilon)^T] \\
 & \quad - (S_c^{-1} S \beta^*)(S_c^{-1} S \beta^*)^T \\
 & = \cancel{(S_c^{-1} S \beta^*)(S_c^{-1} S \beta^*)^T} + E[S_c^{-1} X^T] E[\varepsilon \varepsilon^T] E[X (S_c^{-1})^T] \\
 & \quad - \cancel{(S_c^{-1} S \beta^*)(S_c^{-1} S \beta^*)} \\
 & = S_c^{-1} X^T \underbrace{E[\varepsilon \varepsilon^T]}_{= \sigma^2} X (S_c^{-1})^T \\
 & \quad \text{Covariance matrix symmetric} \Rightarrow (S_c^{-1})^T = S_c^{-1} \\
 & = S_c^{-1} X^T X S_c^{-1} \sigma^2 \\
 & = S_c^{-1} S S_c^{-1} \sigma^2
 \end{aligned}$$

