

Statistical Methods - Homework #3

1. Consider the multiple linear regression model:

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad \varepsilon_i \sim (0, \sigma^2) \text{ independently for } i = 1, \dots, n, \quad (1)$$

where $\mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\beta} \in \mathbb{R}^p$ and $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top \in \mathbb{R}^{n \times p}$ is full rank. ($\text{rank}(X) = p$) Letting $\mathbf{y} = (y_1, \dots, y_n)^\top$ and $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$, (1) can equivalently be written as

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \mathbf{E}(\boldsymbol{\varepsilon}) = \mathbf{0}, \text{ and } \text{Var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n.$$

(A) Find the ridge estimator of $\boldsymbol{\beta}$, denoted $\hat{\boldsymbol{\beta}}_\lambda$, that minimizes

$$\ell(\boldsymbol{\beta}) = (\mathbf{y} - X\boldsymbol{\beta})^\top (\mathbf{y} - X\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta} = \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2 + \lambda \boldsymbol{\beta}^\top \boldsymbol{\beta}.$$

Solution: Differentiate $\ell(\boldsymbol{\beta})$ with respect to $\boldsymbol{\beta}$ and set to zero:

$$\frac{\partial \ell}{\partial \boldsymbol{\beta}} = -2X^\top (\mathbf{y} - X\boldsymbol{\beta}) + 2\lambda \boldsymbol{\beta} = \mathbf{0}.$$

Rearranging gives the normal equations for the ridge estimator:

$$(X^\top X + \lambda I_p) \boldsymbol{\beta} = X^\top \mathbf{y}.$$

Since $X^\top X + \lambda I_p$ is positive definite (hence invertible) for any $\lambda > 0$,

$$\boxed{\hat{\boldsymbol{\beta}}_\lambda = (X^\top X + \lambda I_p)^{-1} X^\top \mathbf{y}.} \quad (2)$$

(B) Find $\mathbf{E}(\hat{\boldsymbol{\beta}}_\lambda)$. What is the value of λ that makes $\mathbf{E}(\hat{\boldsymbol{\beta}}_\lambda) = \boldsymbol{\beta}$?

Solution: Substituting $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ into (2) and taking expectation:

$$\begin{aligned} \mathbf{E}(\hat{\boldsymbol{\beta}}_\lambda) &= (X^\top X + \lambda I_p)^{-1} X^\top \mathbf{E}(\mathbf{y}) \\ &= (X^\top X + \lambda I_p)^{-1} X^\top X \boldsymbol{\beta}. \end{aligned} \quad (3)$$

In general, $\mathbf{E}(\hat{\boldsymbol{\beta}}_\lambda) \neq \boldsymbol{\beta}$ unless $\lambda = 0$, in which case

$$(X^\top X + 0 \cdot I_p)^{-1} X^\top X = (X^\top X)^{-1} X^\top X = I_p,$$

so $\mathbf{E}(\hat{\boldsymbol{\beta}}_0) = \boldsymbol{\beta}$. Hence, $\hat{\boldsymbol{\beta}}_\lambda$ is biased for $\lambda > 0$, and the value $\boxed{\lambda = 0}$ makes the estimator unbiased.

(C) Find $\text{Var}(\hat{\beta}_\lambda)$. Also, show that

$$\text{Var}(\hat{\beta}_\lambda) \preceq \text{Var}(\hat{\beta}_0).$$

(Hint: You may use eigen-value decomposition on $X^\top X$: There exists a matrix $P \in \mathbb{R}^{p \times p}$ and a diagonal matrix $\Lambda \in \mathbb{R}^{p \times p}$ such that $P^\top P = PP^\top = I$ and $X^\top X = P\Lambda P^\top$).

Solution: Since $\hat{\beta}_\lambda = (X^\top X + \lambda I_p)^{-1} X^\top \mathbf{y}$ is a linear function of \mathbf{y} with $\text{Var}(\mathbf{y}) = \sigma^2 I_n$,

$$\begin{aligned} \text{Var}(\hat{\beta}_\lambda) &= (X^\top X + \lambda I_p)^{-1} X^\top \cdot \sigma^2 I_n \cdot X (X^\top X + \lambda I_p)^{-1} \\ &= \sigma^2 (X^\top X + \lambda I_p)^{-1} X^\top X (X^\top X + \lambda I_p)^{-1}. \end{aligned} \quad (4)$$

For the OLS estimator ($\lambda = 0$): $\text{Var}(\hat{\beta}_0) = \sigma^2 (X^\top X)^{-1}$.

To prove $\text{Var}(\hat{\beta}_\lambda) \preceq \text{Var}(\hat{\beta}_0)$, we show $D := \text{Var}(\hat{\beta}_0) - \text{Var}(\hat{\beta}_\lambda) \succeq 0$. Let $X^\top X = P\Lambda P^\top$ be the eigendecomposition, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ with $\lambda_j > 0$. Then $X^\top X + \lambda I_p = P(\Lambda + \lambda I_p)P^\top$, so

$$\begin{aligned} \text{Var}(\hat{\beta}_0) &= \sigma^2 P \Lambda^{-1} P^\top, \\ \text{Var}(\hat{\beta}_\lambda) &= \sigma^2 P (\Lambda + \lambda I_p)^{-1} \Lambda (\Lambda + \lambda I_p)^{-1} P^\top. \end{aligned}$$

Hence

$$D = \sigma^2 P \underbrace{[\Lambda^{-1} - (\Lambda + \lambda I_p)^{-1} \Lambda (\Lambda + \lambda I_p)^{-1}]}_{=: \Delta} P^\top.$$

Since Δ is diagonal, its j -th entry is

$$\Delta_{jj} = \frac{1}{\lambda_j} - \frac{\lambda_j}{(\lambda_j + \lambda)^2} = \frac{(\lambda_j + \lambda)^2 - \lambda_j^2}{\lambda_j(\lambda_j + \lambda)^2} = \frac{\lambda(\lambda + 2\lambda_j)}{\lambda_j(\lambda_j + \lambda)^2} \geq 0$$

for all $\lambda \geq 0$ and $\lambda_j > 0$. Thus $\Delta \succeq 0$, and therefore $D = \sigma^2 P \Delta P^\top \succeq 0$, i.e.,

$$\text{Var}(\hat{\beta}_\lambda) \preceq \text{Var}(\hat{\beta}_0). \quad \square$$

(D) Show the bias-variance decomposition of mean-squared error(MSE):

$$\mathbb{E} \left[(\hat{\beta}_\lambda - \beta)(\hat{\beta}_\lambda - \beta)^\top \right] = \text{Var}(\hat{\beta}_\lambda) + (\beta - \mathbb{E}[\hat{\beta}_\lambda])(\beta - \mathbb{E}[\hat{\beta}_\lambda])^\top.$$

Solution: Let $\mu := \mathbb{E}[\hat{\beta}_\lambda]$ and write $\hat{\beta}_\lambda - \beta = (\hat{\beta}_\lambda - \mu) + (\mu - \beta)$. Then

$$\begin{aligned} &\mathbb{E} \left[(\hat{\beta}_\lambda - \beta)(\hat{\beta}_\lambda - \beta)^\top \right] \\ &= \mathbb{E} \left[\left((\hat{\beta}_\lambda - \mu) + (\mu - \beta) \right) \left((\hat{\beta}_\lambda - \mu) + (\mu - \beta) \right)^\top \right] \\ &= \mathbb{E} \left[(\hat{\beta}_\lambda - \mu)(\hat{\beta}_\lambda - \mu)^\top \right] + \mathbb{E} \left[(\hat{\beta}_\lambda - \mu) \right] (\mu - \beta)^\top + (\mu - \beta) \mathbb{E} \left[(\hat{\beta}_\lambda - \mu)^\top \right] + (\mu - \beta)(\mu - \beta)^\top. \end{aligned}$$

Since $\mathbb{E}[\hat{\beta}_\lambda - \mu] = \mathbf{0}$, the two cross terms vanish. Therefore,

$$\mathbb{E} \left[(\hat{\beta}_\lambda - \beta)(\hat{\beta}_\lambda - \beta)^\top \right] = \underbrace{\mathbb{E} \left[(\hat{\beta}_\lambda - \mu)(\hat{\beta}_\lambda - \mu)^\top \right]}_{= \text{Var}(\hat{\beta}_\lambda)} + \underbrace{(\mu - \beta)(\mu - \beta)^\top}_{= (\beta - \mathbb{E}[\hat{\beta}_\lambda])(\beta - \mathbb{E}[\hat{\beta}_\lambda])^\top},$$

which gives the desired bias-variance decomposition. \square