

CS446: Machine Learning, Fall 2017, Homework 2

Name: Lanxiao Bai (lbai5)

Worked individually

Problem 1

Solution: Since we have model

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \varepsilon$$

by letting

$$\frac{\partial \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2}{\partial \mathbf{w}} = 0$$

,

we can get the least square solution

$$\mathbf{w}_{\text{LS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Problem 2

Solution: In order to check if the least square estimation is biased, we can calculate the bias of the estimation

$$\begin{aligned} \text{bias}(\mathbf{w}_{\text{LS}}) &= \mathbb{E}[\mathbf{w}_{\text{LS}}] - \mathbf{w} \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] - \mathbf{w} \\ &= \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\mathbf{w} + \varepsilon)] - \mathbf{w} \\ &= \mathbb{E}[(\mathbf{X}^{-1}(\mathbf{X}^T)^{-1} \mathbf{X}^T \mathbf{X}\mathbf{w} + (\mathbf{X}^{-1}(\mathbf{X}^T)^{-1} \mathbf{X}^T \mathbf{X}\mathbf{w}\varepsilon)] - \mathbf{w} \\ &= \mathbb{E}[\mathbf{w}] + (\mathbf{X}^{-1}(\mathbf{X}^T)^{-1} \mathbf{X}^T \mathbb{E}[\varepsilon]) - \mathbf{w} \\ &= \mathbf{w} + \mathbf{0} - \mathbf{w} = \mathbf{0} \end{aligned}$$

So we see that the least square estimator is unbiased.

Problem 3

Solution:

$$\begin{aligned}
 \text{Var}(\mathbf{w}_{\text{LS}}) &= \text{Var}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) \\
 &= \text{Var}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \mathbf{w} + \varepsilon)) \\
 &= \text{Var}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w}) + \text{Var}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \varepsilon) \\
 &= \text{Var}(\mathbf{w}) + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\varepsilon) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\
 &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 I (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\
 &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \\
 &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}
 \end{aligned}$$

Problem 4

Solution: Now that we have the objective function

$$\arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

similarly, we let

$$\frac{\partial (\|\mathbf{y} - \mathbf{X} \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2)}{\partial \mathbf{w}} = 0$$

and get

$$\mathbf{w}_{\text{ridge}} = (\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Problem 5

Solution: Similarly, we calculate the bias of $\mathbf{w}_{\text{ridge}}$,

$$\begin{aligned}
 \text{bias}(\mathbf{w}_{\text{ridge}}) &= \mathbb{E}[\mathbf{w}_{\text{ridge}}] - \mathbf{w} \\
 &= \mathbb{E}[(\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] - \mathbf{w} \\
 &= \mathbb{E}[(\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X} \mathbf{w} + \varepsilon)] - \mathbf{w} \\
 &= \mathbb{E}[(\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w}] - \mathbf{w} \\
 &= (\mathbf{I}_D + \lambda (\mathbf{X}^T \mathbf{X})^{-1}) \mathbf{w}
 \end{aligned}$$

Sta (2006)

Hence, the estimator is biased.

Problem 6

Solution:

$$\begin{aligned}
 \text{Var}(\mathbf{w}_{\text{ridge}}) &= \text{Var}((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}) \\
 &= \text{Var}((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \mathbf{w}) + \sigma^2 ((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) ((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T)^T \\
 &= \sigma^2 ((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) ((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T)^T \\
 &= \sigma^2 (\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} ((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1})^T
 \end{aligned}$$

Problem 7

Solution:

$$\begin{aligned}
 \text{tr}(\text{Var}(\mathbf{w}_{\text{LS}})) &= \text{tr}(\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}) \\
 &= \sigma^2 \text{tr}((\mathbf{X}^T \mathbf{X})^{-1}) \\
 &= \sigma^2 \text{tr}((\mathbf{V} \mathbf{S}^2 \mathbf{V}^T)^{-1}) \\
 &= \sigma^2 \text{tr}(\mathbf{V}(\mathbf{S}^2)^{-1} \mathbf{V}^{-1}) \\
 &= \sigma^2 \text{tr}((\mathbf{S}^2)^{-1}) \\
 &= \sigma^2 \sum_{i=1}^n \frac{1}{s_i^2}
 \end{aligned}$$

and

$$\begin{aligned}
 \text{tr}(\text{Var}(\mathbf{w}_{\text{ridge}})) &= \text{tr}(\sigma^2(\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} ((\lambda \mathbf{I}_D + \mathbf{X}^T \mathbf{X})^{-1})^T) \\
 &= \sigma^2 \text{tr}(\mathbf{V}(\lambda \mathbf{I}_D + \mathbf{S}^2)^{-1} \mathbf{V}^T (\mathbf{V} \mathbf{S}^2 \mathbf{V}^T) (\mathbf{V}(\lambda \mathbf{I}_D + \mathbf{S}^2)^{-1} \mathbf{V}^T)^T) \\
 &= \sigma^2 \text{tr}(\mathbf{V}(\lambda \mathbf{I}_D + \mathbf{S}^2)^{-1} \mathbf{V}^T (\mathbf{V} \mathbf{S}^2 \mathbf{V}^T) \mathbf{V}^T ((\lambda \mathbf{I}_D + \mathbf{S}^2)^{-1})^T \mathbf{V}) \\
 &= \sigma^2 \text{tr}((\lambda \mathbf{I}_D + \mathbf{S}^2)^{-1} \mathbf{S}^2 (\lambda \mathbf{I}_D + \mathbf{S}^2)^{-1}) \\
 &= \sigma^2 \text{tr}(((\lambda \mathbf{I}_D + \mathbf{S}^2)^{-1})^2 \mathbf{S}^2) \\
 &= \sigma^2 \sum_{i=1}^n \frac{s_i^2}{\lambda + s_i^2}
 \end{aligned} \tag{1}$$

Problem 8

Solution: We see that the bias of \mathbf{w}_{LS} is smaller than that of $\mathbf{w}_{\text{ridge}}$ but the variance of \mathbf{w}_{LS} is larger than that of $\mathbf{w}_{\text{ridge}}$.

References

(2006). Regularization: Ridge regression and the lasso.

URL <http://statweb.stanford.edu/~tibs/sta305files/Rudyregularization.pdf>