

Homework2

Name: Chenqing Ji

Student ID:11911303

Statistical Learning for Data Science

Due time: March 20, 2023 (Monday) 12:00am

1 Proof

Proof the bias-variance trade-off in detail:
(hint: bias/variance decomposition of MSE)

Proof:

For Bias/variance decomposition of MSE, we can have:

$$\begin{aligned} E \left(y_0 - \hat{f}(x_0) \right)^2 &= E \left(f(x_0) + \varepsilon - \hat{f}(x_0) \right)^2 \\ &= E \left(f(x_0) - E(\hat{f}(x_0)) + E(\hat{f}(x_0)) - \hat{f}(x_0) + \varepsilon \right)^2 \\ &= E \left(\begin{aligned} &\left(f(x_0) - E(\hat{f}(x_0)) \right)^2 + \left(E(\hat{f}(x_0)) - \hat{f}(x_0) \right)^2 \\ &+ 2 \left(f(x_0) - E(\hat{f}(x_0)) \right) \left(E(\hat{f}(x_0)) - \hat{f}(x_0) \right) \\ &+ \varepsilon^2 + 2\varepsilon \left(f(x_0) - E(\hat{f}(x_0)) \right) \\ &+ 2\varepsilon \left(E(\hat{f}(x_0)) - \hat{f}(x_0) \right) \end{aligned} \right) \end{aligned} \quad (1)$$

So, in the equation (1), there exists 6 terms. As we know, $f(x_0)$ can be regarded as a constant in this expression, $\hat{f}(x_0)$ and ε can be regarded as variables in this expression. It's worth noting that although $\hat{f}(x_0)$ can be regarded as variable, $E(\hat{f}(x_0))$ should be regarded as a constant.

Therefore, some terms can be simplified in the equation (1) like that:

$$E \left(2 \left(f(x_0) - E(\hat{f}(x_0)) \right) \left(E(\hat{f}(x_0)) - \hat{f}(x_0) \right) \right) = 0 \quad (2)$$

$$E \left(2\varepsilon \left(E(\hat{f}(x_0)) - \hat{f}(x_0) \right) \right) = 0 \quad (3)$$

$$E \left(2\varepsilon \left(f(x_0) - E(\hat{f}(x_0)) \right) \right) = 0 \quad (4)$$

Note that the equations (3) and (4) are true because $E(\varepsilon)$ is 0. Then, the equation (1) can be simplified as:

$$E \left(y_0 - \hat{f}(x_0) \right)^2 = \left(f(x_0) - E(\hat{f}(x_0)) \right)^2 + E \left(\left(E(\hat{f}(x_0)) - \hat{f}(x_0) \right)^2 \right) + E(\varepsilon^2) \quad (5)$$

From the equation (5), the first term represent the **Square** of the **Bias** for $\hat{f}(x_0)$. The second term represent the variance of $\hat{f}(x_0)$ and the third term represent the variance of ε , which is also called **irreducible error**.

Therefore, due to the Bias/variance decomposition of MSE, the expected test MSE can be expressed in the equation (6) below, which is called "The Bias-Variance Trade-off".

$$\text{MSE} = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\varepsilon) \quad (6)$$

2 Proof

For simple linear regression, proof:

$$\begin{aligned} \text{Var}(\hat{\beta}_0) &= \text{SE}(\hat{\beta}_0)^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \\ \text{Var}(\hat{\beta}_1) &= \text{SE}(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

Proof:

For simple linear regression: $y_i = \beta_1 x_i + \beta_0 + \varepsilon$, ε is a random variable, so y_i is also a random variable.

$$\text{Var}(\hat{\beta}_1) = \text{Var}\left(\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) = \text{Var}\left(\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i)}{\sum_{i=1}^n (x_i - \bar{x})^2} - \bar{y} \frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}\right)$$

$$\text{Because } \sum_{i=1}^n (x_i - \bar{x}) = 0, \bar{y} \frac{\sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} = 0$$

$$\text{Var}(\hat{\beta}_1) = \text{Var}\left(\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i)}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \text{Var}(y_i)}{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)^2}, \text{Var}(y_i) = \text{Var}(\varepsilon) = \sigma^2$$

$$\text{So } \text{Var}(\hat{\beta}_1) = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \sigma^2}{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right)^2} = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \Rightarrow \text{Var}(\hat{\beta}_0) = \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) = \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) - 2\bar{x} \text{Cov}(\bar{y}, \hat{\beta}_1)$$

$$\text{Suppose } c_i = \frac{(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Cov}(\bar{y}, \hat{\beta}_1) = \text{Cov}\left(\frac{1}{n} \sum y_i, \sum c_j y_j\right) = \frac{1}{n} \sum c_j \sum \text{Cov}(y_i, y_j) = \frac{1}{n} \sum c_j \text{Cov}(y_j, y_j) = \frac{1}{n} \sum c_i = 0$$

$$\text{So, Var}(\hat{\beta}_0) = \text{Var}(\bar{y}) + \bar{x}^2 \text{Var}(\hat{\beta}_1) = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

3 Problem

3.7.1 Solution:

The null hypotheses are that: For "TV", in the presence of radio ads and newspaper ads, TV ads have no effect on sales. Similarly, for "radio", in the presence of TV and newspaper ads, radio ads have no effect on sales; For "newspaper", in the presence of TV and radio ads, newspaper ads have no effect on sales. Then, from the p-value we can get that: The low p-values of TV and radio suggest that the null hypotheses for TV and radio are false. However, the high p-value of newspaper suggests that the null hypothesis for newspaper is true.

3.7.2 Solution:

As we all know, KNN is a non-parametric algorithm used for both classification and regression tasks. The main difference between KNN classifier and KNN regression methods is the type of output that they produce: KNN classifier method produces a categorical output, while KNN regression method produces a continuous output.

3.7.3 Solution:

- (a) The relationship between the starting salary after graduation (in thousands of dollars) and related five predictors is shown in (7) and (8):

$$\text{Salary} = 50 + 20 * GPA + 0.07 * IQ + 35 * Gender + 0.01 * (GPA * IQ) - 10 * (GPA * Gender) \quad (7)$$

It will convert to:

$$Y = 50 + 20 * X_1 + 0.07 * X_2 + 35 * X_3 + 0.01 * X_4 - 10 * X_5 \quad (8)$$

Where $X_4 = X_1 * X_2$ and $X_5 = X_1 * X_3$ For male, the gender is 0 and for female the gender is 1. Therefore, the equation (8) for male and female will convert to the equation (9) and (10):

$$Y = 50 + 20 * X_1 + 0.07 * X_2 + 0.01 * X_4 \quad (9)$$

$$Y = 50 + 20 * X_1 + 0.07 * X_2 + 35 + 0.01 * X_4 - 10 * X_1 \quad (10)$$

Based on the equation (7)(8)(9)(10), we can get that: Once the GPA(X_1) is high enough, for a fixed value of IQ and GPA, males earn more on average than females. Therefore, the answer iii. is true.

- (b) If a female with IQ of 110 and a GPA of 4.0, her profile will be: $X_1 = 4$, $X_2 = 110$, $X_4 = X_1 * X_2 = 4 * 110$.

Her salary will be:

$$\begin{aligned} Y &= 50 + 20 * 4 + 0.07 * 110 + 35 + 0.01(4 * 110) - 10 * 4 \\ &= 137.1 \end{aligned} \quad (11)$$

Therefore, her salary will be **137.1** (in thousands of dollars).

- (c) **False.** Because before we judge it, we should carefully check the p-value of the coefficient for the GPA/IQ interaction term to determine whether the interaction term is statistically significant or not.

3.7.4 Solution:

- (a) Yes, we would expect that one to be lower than the other. That is, the cubic regression will have a lower training RSS than the linear regression because the cubic regression can make a closer fit with the data that has with a larger irreducible error.
- (b) Different from (a), we expect the cubic regression to have a higher test RSS than the linear regression because the training process will overfit the cubic regression so it will have more error (RSS) than the linear regression in the test process.
- (c) My answer to this problem is the same as the problem (a). That is: the cubic regression will have a lower training RSS than the linear regression. Because the cubic regression still has higher flexibility than linear regression no matter what the true relationship between X and Y is. Then, it will fit the training points closer than the linear regression and reduce the training error (RSS).
- (d) It is a shame that there is not enough evidence to tell which test RSS will be lower between the linear regression and the cubic regression because we don't know how far it is from linear for the true relationship between X and Y. If the true relationship between X and Y is closer to linear than cubic, the test RSS for linear regression could be lower than the related cubic regression. However, if the true relationship between X and Y is closer to the cubic than linear, the test RSS for the cubic regression could be lower than the related linear regression.

3.7.5 Solution:

For this problem, we have the known information about equation (12) and (13):

$$\hat{y}_i = x_i \hat{\beta} \quad (12)$$

$$\hat{\beta} = \left(\sum_{i=1}^n x_i y_i \right) / \left(\sum_{i'=1}^n x_{i'}^2 \right). \quad (13)$$

Based on that, we can get:

$$\hat{y}_i = x_i \cdot \left(\frac{\sum_{i=1}^n x_i y_i}{\sum_{i'=1}^n x_{i'}^2} \right) \quad (14)$$

After that, we just need to place y_i outside the fraction and bring the summation sign before the fraction to get the following formula:

$$\hat{y}_i = \sum_{i'=1}^n \frac{x_{i'} x_i}{\sum_{j=1}^n x_j^2} y_{i'} \quad (15)$$

From the above equation, we can easily see that the form of this equation is similar to the following equation (16):

$$\hat{y}_i = \sum_{i'=1}^n a_{i'} y_{i'} \quad (16)$$

Therefore, compare with the equation (15) and (16), we can get the $a_{i'}$, that is:

$$a_{i'} = \frac{x_{i'} x_i}{\sum_{j=1}^n x_j^2} \quad (17)$$

3.7.6 Solution:

For the simple linear regression we know:

$$y = \hat{\beta}_0 + \hat{\beta}_1 x \quad (18)$$

Then, from the equation (3.4) in the textbook we can get:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (19)$$

Then, we can change its form:

$$\hat{\beta}_0 + \hat{\beta}_1 \bar{x} - \bar{y} = 0 \quad (20)$$

Then, if the least squares line always passes through the point (\bar{x}, \bar{y}) , we can easily see that the equation (18) will become:

$$\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x} \quad (21)$$

Then, we get the right value of $\hat{\beta}_0$ from equation (19) and substitute that into equation (21), we can get:

$$\bar{y} = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 \bar{x} \quad (22)$$

Because the right side of the equation (22) is \bar{y} , and the left side of the equation (22) is also \bar{y} . So, the equation (22) is always right when the least squares line always passes through the point (\bar{x}, \bar{y}) .

In a conclusion, in the case of simple linear regression, the least squares line always passes through the point (\bar{x}, \bar{y}) .

3.7.7 Proof:

In the case of simple linear regression, R2 statistic is defined as:

$$R^2 = \frac{TSS - RSS}{TSS} \quad (23)$$

Then, TSS and RSS is defined as:

$$TSS = \sum (y_i - \bar{y})^2 \quad (24)$$

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (25)$$

Then, the correlation between X and Y is:

$$Cor(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \quad (26)$$

Next, we start to proof the conclusion that this problem gives us. The process is shown below:

Firstly, we substitute the definition of TSS and RSS into R^2 , that is:

$$R^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (27)$$

Then, we look at the numerator and suppose the value of it is Z, by using the square variance formula it will become:

$$\begin{aligned} Z &= \sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \sum_{i=1}^n [(y_i - \bar{y}) - (y_i - \hat{y}_i)] [(y_i - \bar{y}) + (y_i - \hat{y}_i)] \\ &= \sum_{i=1}^n (\hat{y}_i - \bar{y})(2y_i - \bar{y} - \hat{y}_i) \end{aligned} \quad (28)$$

Next, in the case of simple linear regression of Y onto X , we should know that:

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (29)$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (30)$$

So, we can substitute the expression for $\hat{\beta}_0$ into \hat{y}_i , that is:

$$\begin{aligned} \hat{y}_i &= \hat{\beta}_0 + \hat{\beta}_1 x_i \\ &= \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i \\ &= \bar{y} + \hat{\beta}_1 (x_i - \bar{x}) \end{aligned} \quad (31)$$

Therefore, we can use the above conclusion from equation (31) to simplify the partial expression of Z , that is:

$$\hat{y}_i - \bar{y} = \hat{\beta}_1 (x_i - \bar{x}) \quad (32)$$

$$\begin{aligned} 2y_i - \bar{y} - \hat{y}_i &= 2y_i - \bar{y} - \bar{y} - \hat{\beta}_1 (x_i - \bar{x}) \\ &= 2(y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x}) \end{aligned} \quad (33)$$

After that, we can substitute the equation (32) and (33) to the expression of Z , that is:

$$\begin{aligned} Z &= \sum_{i=1}^n \hat{\beta}_1 (x_i - \bar{x}) \left[2(y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x}) \right] \\ &= \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x}) \left[2(y_i - \bar{y}) - \hat{\beta}_1 (x_i - \bar{x}) \right] \\ &= \hat{\beta}_1 \left[2 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})^2 \right] \end{aligned} \quad (34)$$

From the equation (30), we know that the equation can be simplified as:

$$\begin{aligned} Z &= \hat{\beta}_1 \left[2 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})^2 \right] \\ &= \hat{\beta}_1 \left[2 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right] \\ &= \hat{\beta}_1 \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{[\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})]^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned} \quad (35)$$

At this point, we can substitute our final reduced expression for Z into R^2 , that is:

$$R^2 = \frac{[\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})]^2}{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2} \quad (36)$$

Compare with the equation (26), we can finally get:

$$R^2 = Cor^2(X, Y) \quad (37)$$

So, in the case of simple linear regression of Y onto X , the R^2 statistic is equal to the square of the correlation between X and Y .