

CONCEPTUAL QUESTIONS

CHAPTER - 3

1. Describe the null hypotheses to which the p-values given in Table 3.4 correspond. Explain what conclusions you can draw based on these p-values. Your explanation should be phrased in terms of sales, TV, radio, and newspaper, rather than in terms of the coefficients of the linear model.

Roughly speaking, p value is the probability that the given result is true, considering the null hypothesis. If the p value is very less, we say that the p value is significant and we reject the null hypothesis. In a given linear model, null hypothesis is that there is no relationship between a predictor and response, and there low or significant p means that we reject null hypothesis and that there is a relationship between predictor and response.

In the given table. Table 3.4, the p value of TV, radio and newspaper are <0.001 , <0.001 , and 0.8599 respectively. P-values for TV and radio are significant, and hence there is a relationship between TV and sales, and , radio and sales, However there is a very large p value for newspaper, and hence we can conclude that there is no relationship between newspaper and radio.

2. Carefully explain the differences between the KNN classifier and KNN regression methods.

KNN classifier is used for classification and for a given input X_0 , it selects the k points which are closest to X_0 , and outputs the label contained by the maximum of the k points. In KNN regression, which is used to predict quantitative values, it selects k points closest to the input X_0 , and returns the average of all the values.

3. Suppose we have a data set with five predictors, X_1 = GPA, X_2 = IQ, X_3 = Gender (1 for Female and 0 for Male), X_4 = Interaction between GPA and IQ, and X_5 = Interaction between GPA and Gender. The response is starting salary after graduation (in thousands of dollars). Suppose we use least squares to fit the model, and get $\hat{\beta}_0 = 50$, $\hat{\beta}_1 = 20$, $\hat{\beta}_2 = 0.07$, $\hat{\beta}_3 = 35$, $\hat{\beta}_4 = 0.01$, $\hat{\beta}_5 = -10$.

(a) Which answer is correct, and why?

(iii)

$$\begin{aligned} \text{Sales} &= \hat{\beta}_0 + \hat{\beta}_1 \text{GPA} + \hat{\beta}_2 \text{IQ} + \hat{\beta}_3 \text{Gender} + \hat{\beta}_4 \text{GPA} * \text{IQ} \\ &\quad + \hat{\beta}_5 \text{GPA} * \text{GENDER} \\ \text{Sales}(\text{man}) &= \hat{\beta}_0 + \hat{\beta}_1 \text{GPA} + \hat{\beta}_2 \text{IQ} + \hat{\beta}_4 \text{GPA} * \text{IQ} \\ \text{Sales}(\text{woman}) &= \hat{\beta}_0 + \hat{\beta}_1 \text{GPA} + \hat{\beta}_2 \text{IQ} + \hat{\beta}_3 + \hat{\beta}_4 \text{GPA} * \text{IQ} + \hat{\beta}_5 \text{GPA} \\ \text{Sales}(\text{woman}) - \text{Sales}(\text{man}) &= \hat{\beta}_3 + \hat{\beta}_5 \text{GPA} \\ &= 10 - 35 \text{GPA} \\ &= 35 - 10 * \text{GPA} \end{aligned}$$

Therefore, for fixed value of IQ and GPA, if the value of $(10 - 35 * \text{GPA}) > 0$, then sales of male will be greater than female, ie, for $\text{GPA} > 3.5$, sales of male is larger than sale of female.

(b) Predict the salary of a female with IQ of 110 and a GPA of 4.0.

$$\text{Sales} = 50 + 20 * \text{GPA} + 0.007 * \text{IQ} + 35 * \text{GENDER} + 0.01(\text{GPA} * \text{IQ}) - 10(\text{GPA} * \text{GENDER})$$

Python code

```
x = [4,110,1]
```

$$\text{Sales} = 50 + x[0] * 20 + x[1] * 0.07 + x[2] * 35 + (x[0] * x[1]) * 0.01 + (x[0] * x[2]) * (-10)$$

```
Print(Sales)
```

137.1

(c) True or false: Since the coefficient for the GPA/IQ interaction term is very small, there is very little evidence of an interaction effect. Justify your answer.

False. For calculating the significance the value of coefficient is not the best choice. Hence we can't conclude that it is there is little evidence of an interaction effect.

4. I collect a set of data ($n = 100$ observations) containing a single predictor and a quantitative response. I then fit a linear regression model to the data, as well as a separate cubic regression, i.e. $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \dots$

(a) Suppose that the true relationship between X and Y is linear, i.e. $Y = \beta_0 + \beta_1 X + \dots$. Consider the training residual sum of squares (RSS) for the linear regression, and also the training RSS for the cubic regression.

Would we expect one to be lower than the other, would we expect them to be the same, or is there not enough information to tell? Justify your answer.

For training data, RSS decreases as we increase model complexity. The model will overfit to find the patterns and will overfit the data. So, RSS for cubic regression will be lower than RSS of linear regression

(b) Answer (a) using test rather than training RSS.

In the case of test data, for cubic regression which is overfitted, will produce results far from the true values, whereas linear regression, which is close to the true function will perform better on the test data. And hence RSS will be lower on test data for linear regression

(c) Suppose that the true relationship between X and Y is not linear, but we don't know how far it is from linear. Consider the training RSS for the linear regression, and also the training RSS for the cubic regression. Would we expect one to be lower than the other, would we expect them to be the same, or is there not enough information to tell? Justify your answer.

Since, training RSS decreases with increasing complexity, we will observe lower RSS for cubic regression.

(d) Answer (c) using test rather than training RSS

We can't predict whether RSS for test data will be lower for cubic regression or linear regression, since we have no clue what the true function is like.

We can find out the value of RSS for test data for cubic regression and linear regression, If test RSS of cubic is less than test RSS of linear regression, then we can say that true relationship is more non-linear.

5. Consider the fitted values that result from performing linear regression without an intercept. In this setting, the i th fitted value takes the form $\hat{y}_i = x_i \hat{\beta}$, where $\hat{\beta} = \sum_{i=1}^n x_i y_i / \sum_{i=1}^n x_i^2$. (3.38) Show that we can write $\hat{y}_i = \sum_{j=1}^n a_{ij} y_j$. What is a_i ?

$$\underline{3.5} \quad y_i = x_i \hat{\beta}$$

$$\hat{\beta} = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} = \frac{\sum_{j=1}^n x_j y_j}{\sum_{j=1}^n x_j^2} \quad - (1)$$

$$y_i = x_i \hat{\beta}$$

$$= x_i \frac{\sum_{j=1}^n x_j y_j}{\sum_{j=1}^n x_j^2} = \frac{\sum_{j=1}^n x_i x_j y_j}{\sum_{j=1}^n x_j^2}$$

$$= \sum_{j=1}^n \frac{x_i x_j}{\sum_{j=1}^n x_j^2} y_j$$

$$= \sum_{j=1}^n a_j y_j$$

$$\text{where } a_j = \frac{x_i x_j}{\sum_{j=1}^n x_j^2}$$

6. Using (3.4), argue that in the case of simple linear regression, the least squares line always passes through the point (\bar{x}, \bar{y}) .

3.6 The RSS Line is \rightarrow

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X$$

where $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad \text{--- (1)}$

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

$$Y = \hat{\beta}_0 + \hat{\beta}_1 X$$

$$Y = \hat{\beta}_0 \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 X \quad (\text{from eq (1)})$$

for $X = \bar{x}$,

$$Y = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 \bar{x}$$

$$Y = \bar{y}$$

hence (\bar{x}, \bar{y}) lies on RSS line.

7. It is claimed in the text that in the case of simple linear regression of Y onto X , the R^2 statistic (3.17) is equal to the square of the correlation between X and Y (3.18). Prove that this is the case. For simplicity, you may assume that $\bar{x} = \bar{y} = 0$.

3.7 To prove: $R^2 = \text{cor}(x, y)^2$ when $\bar{x} = \bar{y} = 0$.

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum (y_i - \hat{y})^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\sum (y_i - \beta_0 - \beta_1 x_i)^2}{\sum y_i^2}$$

$$\text{where, } \beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum x_i y_i}{\sum x_i^2}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x} = 0,$$

hence,

$$R^2 = 1 - \frac{\sum (y_i - \beta_1 x_i)^2}{\sum y_i^2} = 1 - \frac{\sum (y_i^2 - 2\beta_1 x_i y_i + \beta_1^2 x_i^2)}{\sum y_i^2}$$

$$R^2 = \frac{\sum y_i^2 - \sum y_i^2 + 2\beta_1 \sum x_i y_i - \beta_1^2 \sum x_i^2}{\sum y_i^2} = \frac{2\cancel{\sum x_i y_i} \sum x_i y_i}{\sum y_i^2}$$

$$R^2 = \frac{2 \cdot \frac{\sum x_i y_i}{\sum x_i^2} \cdot \sum x_i y_i - \left(\frac{\sum x_i y_i}{\sum x_i^2} \right)^2 (\sum x_i^2)}{\sum y_i^2}$$

$$R^2 = \frac{2 \frac{(\sum x_i y_i)^2}{\sum x_i^2} - \frac{(\sum x_i y_i)^2}{\sum x_i^2}}{\sum y_i^2}$$

$$R^2 = \frac{\sum (x_i y_i)^2}{\sum x_i^2 \sum y_i^2} = \frac{\sum (x_i y_i)^2}{\sum x_i^2 \times \sum y_i^2} \quad \text{--- (1)}$$

$$\text{COR}(X, Y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}}$$

$$= \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \cdot \sum y_i^2}}$$

$$\text{COR}(X, Y)^2 = \frac{\sum (x_i y_i)^2}{\sum x_i^2 \cdot \sum y_i^2} \quad \text{--- (2)}$$

from (1) & (2)

$$R^2 = \text{COR}(X, Y)^2$$