

1. Exercises

1.1 ISLR 2e (Gareth James, et al.): Section 3.7 (Exercises), page 120: Exercises 1, 3, 4-a.

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.

	Coefficient	Std. error	t -statistic	p -value
Intercept	2.939	0.3119	9.42	< 0.0001
TV	0.046	0.0014	32.81	< 0.0001
radio	0.189	0.0086	21.89	< 0.0001
newspaper	-0.001	0.0059	-0.18	0.8599

TABLE 3.4. For the Advertising data, least squares coefficient estimates of the multiple linear regression of number of units sold on TV, radio, and newspaper advertising budgets.

H_0^1 : The amount of budget dedicated to TV advertisements has no effect on sales.

H_0^2 : The amount of money in the radio advertisements budget does affect sales.

H_0^3 : There is no relationship between sales and newspaper advertisement.

Based on these p -values, I can conclude that the p -values for TV and radio are highly significant, while the p -value for newspaper advertising is not significant in comparison. This means I can reject H_0^1 and H_0^2 , but not H_0^3 . Therefore, I can say that the budget for newspaper advertising does not affect sales.

3. Suppose we have a data set with five predictors, $X_1 = \text{GPA}$, $X_2 = \text{IQ}$, $X_3 = \text{Level}$ (1 for College and 0 for High School), $X_4 = \text{Interaction between GPA and IQ}$, and $X_5 = \text{Interaction between GPA and Level}$. 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?

- For a fixed value of IQ and GPA, high school graduates earn more, on average, than college graduates.
- For a fixed value of IQ and GPA, college graduates earn more, on average, than high school graduates.
- For a fixed value of IQ and GPA, high school graduates earn more, on average, than college graduates provided that the GPA is high enough.
- For a fixed value of IQ and GPA, college graduates earn more, on average, than high school graduates provided that the GPA is high enough.

The least squares line is given by:

$$\hat{y} = 50 + 20 \cdot \text{GPA} + 0.07 \cdot \text{IQ} + 35 \cdot \text{Level} + 0.01 \cdot (\text{GPA} \times \text{IQ}) - 10 \cdot (\text{GPA} \times \text{Level})$$

For college students, (plug in Level=1) it becomes:

$$\hat{y} = 85 + 10 \cdot \text{GPA} + 0.07 \cdot \text{IQ} + 0.01 \cdot (\text{GPA} \times \text{IQ})$$

and for high school students, (plug in Level=0) it becomes:

$$\hat{y} = 50 + 20 \cdot \text{GPA} + 0.07 \cdot \text{IQ} + 0.01 \cdot (\text{GPA} \times \text{IQ})$$

With high school students and college students, their starting salary differ with the following equation: $50 + 20 \cdot \text{GPA} = 85 + 10 \cdot \text{GPA}$.

$$10 \cdot \text{GPA} = 35, \text{GPA} = 3.5$$

So the starting salary for high school students is higher than for college students on average if $50 + 20 \cdot \text{GPA} \geq 85 + 10 \cdot \text{GPA}$ which is equivalent to $\text{GPA} \geq 3.5$.

Therefore iii. is the right answer.

b. Predict the salary of a college graduate with IQ of 110 and a GPA of 4.0.

$$\begin{aligned}\hat{y} &= 85 + 10 \cdot \text{GPA} + 0.07 \cdot \text{IQ} + 0.01 \cdot (\text{GPA} \times \text{IQ}) \\ &= 85 + 10 \cdot 4.0 + 0.07 \cdot 110 + 0.01 \cdot (4.0 \cdot 110) = 137.1\end{aligned}$$

A college graduate with IQ of 110 and a GPA of 4.0 is predicted to get a salary of \$137,100.

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.

This statement is false because this isn't the proper way to test the impact of the interaction effect. In order to verify the impact of GPA and IQ on the quality of the model, I would have to test the null hypothesis of setting $\beta_4 = 0$. Then I will look at the p-value associated with the t of the F statistic to draw a conclusion.

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. $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$.

a. Suppose that the true relationship between X and Y is linear, i.e. $Y = \beta_0 + \beta_1 X + \epsilon$.

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.

Without knowing more details about the training data, it is difficult to know which training RSS is lower between linear or cubic. However, as the true relationship between X and Y is linear, we may expect the least-squares line to be close to the true regression line, and consequently, the RSS for the linear regression may be lower than for the cubic regression.