

Chapter 3 : Linear Regression

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.

As we see above in table 3.4, the p values for all the predictors except newspaper is <0.05 making the values obtained by them, (Coefficient, Std. Error, t-stat) statistically significant. Whereas for newspaper, p value > 0.05, meaning there is insufficient evidence to conclude that a relationship between the predictor 'newspaper' and the response variable sales exists. A p value dictates that how much of the result one obtained from an experiment was by chance or random error, this difference arises as the sample and population datasets are different, so there can be more than 1 sample for a population and all the different samples can predict different outcomes. Hence having a lesser p value (<0.05) means that our null hypothesis of the predictor and response variable having no relationship among them ($\beta = 0$) {coefficients in this case} is false and there is in fact a relationship among them which is not by chance.

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

The major difference lies between what kind of **problem is being solved or Our Goal** , if the response variable is qualitative then we use a classifier , if it is quantitative then we use regression. | Both the methods are used for prediction and both are Non Parametric as well

A lower K value corresponds to a Model that has low bias but high variance due to it being highly dependent on just one dataset whereas

A high K value gives a smoother fit with high bias but low variance as now the predicted value is dependent on more than one data point , so even if one changes the other data points stabilizes it. Hence a Proper K value is chosen via Bias Variance tradeOff comparing Test MSE values for regression and error/accuracy test for Classification of different models with different K values

	KNN Classifier	KNN Regression
Output Produced	Qualitative/ Categorical	Quantitative / numerical
How is K used?	The nearest K data points are chosen and a majority rule is applied by default , if the differing variables are the same, then a tie breaker rule is applied.	The average of the closest K datapoints are taken and the predicted value is that average of the nearest K values.

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 $\beta_0=50$, $\beta_1=20$, $\beta_2=0.07$, $\beta_3=35$, $\beta_4=0.01$, $\beta_5= - 10$.

(a)Which answer is correct, and why?

Model formed

Starting sal = $50 + 20 \text{ gpa} + 0.07 \text{ iq} + 35 \text{ level} + 0.01 \text{ gpa iq} - 10 \text{ gpa level}$

For HS , level =0

Starting Sal = $50 + 20 \text{ gpa} + 0.07 \text{ iq} + 0 + 0.01 \text{ gpa iq} - 0$

For College Level = 1

Starting sal = $50 + 20 \text{ gpa} + 0.07 \text{ iq} + 35 + 0.01 \text{ gpa iq} - 10 \text{ gpa}$

CS > HS

CS - HS >0

$50 + 20 \text{ gpa} + 0.07 \text{ iq} + 35 + 0.01 \text{ gpa iq} - 10 \text{ gpa} - 50 - 20 \text{ gpa} - 0.07 \text{ iq} - 0.01 \text{ gpa iq} > 0$

$35 - 10 \text{ gpa} > 0$

$3.5 > \text{gpa}$

Hence for CS starting sal > HS starting sal their gpa must be < 3.5 , if it is > 3.5 then HS students will get more starting salary

HS>CS

$$50 + 20 \text{ gpa} + 0.07 \text{ iq} + 0 + 0.01 \text{ gpa iq} - 0 - 50 - 20 \text{ gpa} - 0.07 \text{ iq} - 35 - 0.01 \text{ gpa iq} + 10 \text{ gpa} > 0$$

$$-35 + 10 \text{ gpa} > 0$$

$$\text{gpa} > 3.5$$

Same result as above

i. For a fixed value of IQ and GPA, high school graduates earn more, on average, than college graduates.

True if GPA > 3.5

ii. For a fixed value of IQ and GPA, college graduates earn more, on average, than high school graduates.

True if GPA < 3.5

iii. 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.

True

iv. 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.

False

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

For College Level = 1

$$\text{Starting sal} = 50 + 20 \text{ gpa} + 0.07 \text{ iq} + 35 + 0.01 \text{ gpa iq} - 10 \text{ gpa}$$

Substituting for given values

$$\text{Starting sal} = 50 + 20 * 4 + 0.07 * 110 + 35 + 0.01 * 4 * 110 - 10 * 4$$

$$= 50 + 80 + 7.7 + 35 + 4.4 - 40$$

$$= 137.1 \Rightarrow \$ 137100$$

(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 assumption is false. Even if the coefficient of the interaction term is very small, its magnitude depends on the scale and since it is positive, we know that they have positive synergistic effect. But whether it needs to be taken into account or not is determined by Hypothesis testing and comparing p values determined by it, if the p value < 0.05 - that means there is an interaction effect

and we need to take this effect into account, whereas if $p \text{ value} > 0.05$ that means that the null hypothesis has enough evidence to be true and we will not take this effect into account. Another note to take into account is the **hierarchical principle** which dictates that whether the individual terms like GPA and IQ themselves have high p-value and interaction term has low p value, the interaction term will still be taken into account meaning, IF the interaction term has low p value and is statistically significant, The individual predictors having high p value don't mean anything and we must include them in our model as they are related to the interaction effect which in turn is related to the response variable, therefore leaving them out can alter the meaning of the model.

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, A separate cubic regression, i.e.

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon; \quad Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \epsilon$$

(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, the same, or is there not enough information to tell? Justify your answer.

For Training data: If the true relationship is Linear then for a non linear cubic regression model The RSS value would be lower than a linear model because of overfitting, The highly flexible model will learn not only the patterns but also the noise associated with the training data making it highly efficient but only on training data, it will be a case of low bias high variance, and it will fail to generalize the dataset as a whole, the highest reaching that of its linear counterpart as the linear equation is just the cubic regression with some coefficients $= 0$. But for a linear curve on a training dataset, the RSS value would be higher

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

For Test Dataset the RSS value of the cubic regression model would be higher than the generalized linear model as the highly flexible model now has learnt all the noise in the training data, and different datasets have different noise, so therefore, it will fail to reproduce appropriate results for newer data points whereas a more generalized model where true relationship is also linear will get a value of lower RSS on Test data.

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

For a training dataset , The RSS value of a non linear curve will always be lower than its linear counterpart simply because it can overfit the noise , and learnt the dataset very thoroughly regardless whether the true relationship is linear or not.

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

For a Test dataset the RSS value depends on which generalised the training data well enough to work on test dataset , here IF the true relationship is non linear then a non linear curve will have a lower RSS than linear simply because the linear model will fail to capture the complexity and will have a high Bias

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

$$\hat{\beta} = \left(\sum_{i=1}^n x_i y_i \right) / \left(\sum_{i'=1}^n x_{i'}^2 \right). \quad (3.38)$$

Show that we can write

$$\hat{y}_i = \sum_{i'=1}^n a_{i'} y_{i'}.$$

What is $a_{i'}$?

Note: We interpret this result by saying that the fitted values from linear regression are linear combinations of the response values.

given $\rightarrow \hat{y} = x_i \beta$

$$\hat{\beta} = \frac{\left(\sum_{i=1}^n x_i y_i \right)}{\left(\sum_{i=1}^n (x_i)^2 \right)}$$

Derive $\hat{y} = \sum_{i'=1}^n a_{i'} y_{i'}$

(1) substitute $\hat{\beta}$ in the equation set

$$\hat{y} = x_i \frac{\left(\sum_{i=1}^n x_i y_i \right)}{\left(\sum_{i=1}^n (x_i)^2 \right)}$$

(2) change i to i' we can rename the dummy variable as

↓

$$\hat{y} = x_i \frac{\left(\sum_{i'=1}^n x_{i'} y_{i'} \right)}{\left(\sum_{i'=1}^n (x_{i'})^2 \right)}$$

(3) now we know that both summation terms are constant numbers. we \rightarrow

$$\hat{y} = \frac{x_i}{\left(\sum_{i'=1}^n (x_{i'})^2 \right)} \sum_{i'=1}^n (x_{i'} y_{i'})$$

for a given value $y_{i'}$ the underlined term is a constant

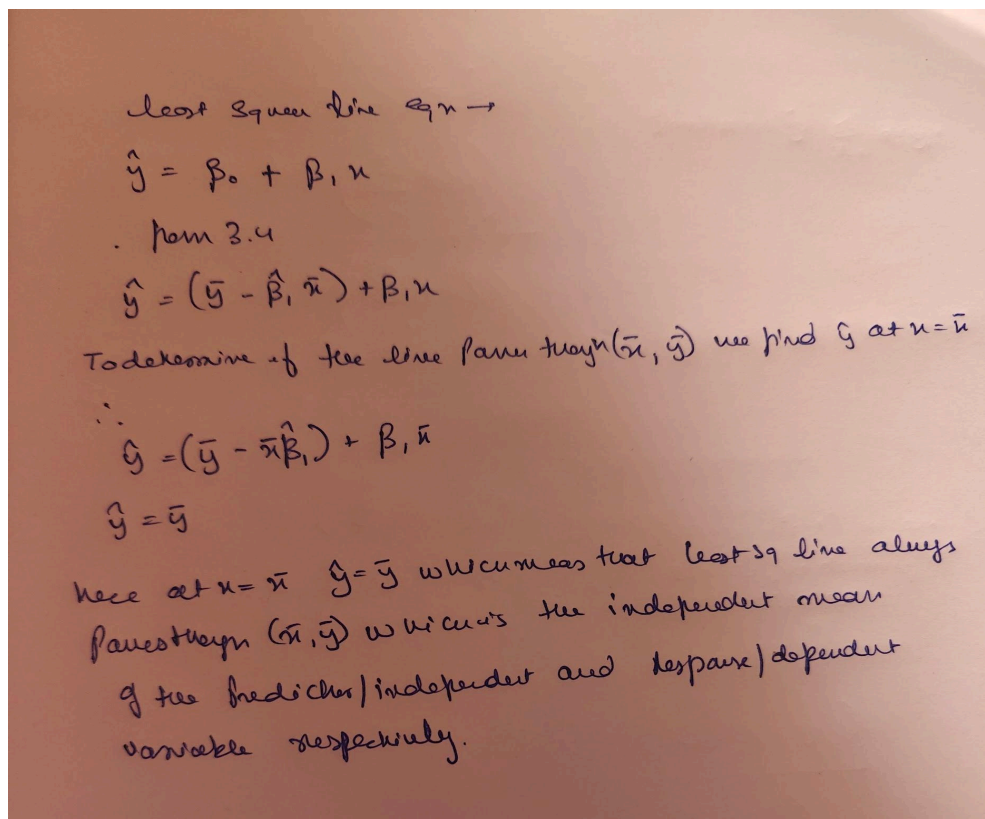
$$\therefore \hat{y} = \sum_{i'=1}^n \left[\frac{x_i x_{i'}}{\sum_{i'=1}^n (x_{i'})^2} \right] \cdot y_{i'} \Rightarrow \hat{y} = \sum_{i'=1}^n a_{i'} y_{i'}$$

$\rightarrow a_{i'}$

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})

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}, \quad (3.4)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x},$$



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 $\bar{x} = \bar{y} = 0$.

To calculate R^2 , we use the formula

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}} \quad (3.17)$$

Recall that *correlation*, defined as

$$\text{Cor}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}, \quad (3.18)$$

assumption given $\Rightarrow \bar{x} = \bar{y} = 0$ — (1)

Step 1 -

$$R^2 = 1 - \frac{RSS}{TSS}$$

Linear regression eqn

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_i \quad \text{--- (2)}$$

$$\text{Here } \beta_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

rem (1)

$$\beta_0 = 0$$

rem (2)

$$\hat{y} = \hat{\beta}_1 x_i \quad \text{--- (4)}$$

$$\text{Total Sum of Squares (TSS)} = \sum_{i=1}^n (y_i - \bar{y})^2$$

rem (1)

$$TSS \Rightarrow \sum_{i=1}^n (y_i)^2$$

rem (4)

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

$$RSS = \sum_{i=1}^n (y_i - \hat{\beta}_1 x_i)^2$$

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

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

$$\bar{x} = \bar{y} = 0$$

$$\Rightarrow \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \quad \text{--- (5)}$$

Step 2 -

$$RSS = \sum_{i=1}^n (y_i - \hat{\beta}_1 x_i)^2$$

$$\Rightarrow \sum_{i=1}^n y_i^2 + \hat{\beta}_1^2 x_i^2 - 2 y_i \hat{\beta}_1 x_i$$

Taking out the constants \rightarrow

$$\Rightarrow \sum_{i=1}^n y_i^2 - 2 \hat{\beta}_1 \sum_{i=1}^n x_i y_i + \hat{\beta}_1^2 \sum_{i=1}^n x_i^2$$

Now substituting $\hat{\beta}_1$ from (5)

$$RSS = \sum_{i=1}^n y_i^2 - 2 \left(\frac{\sum x_i y_i}{\sum x_i^2} \right) \sum x_i y_i + \left(\frac{\sum x_i y_i}{\sum x_i^2} \right)^2 \sum x_i^2$$

$$RSS \Rightarrow \sum_{i=1}^n y_i^2 - 2 \frac{(\sum x_i y_i)^2}{\sum x_i^2} + \frac{(\sum x_i y_i)^2}{(\sum x_i^2)^2} \sum x_i^2$$

$$\Rightarrow \sum_{i=1}^n y_i^2 - 2 \frac{(\sum x_i y_i)^2}{\sum x_i^2} + \frac{(\sum x_i y_i)^2}{\sum x_i^2}$$

$$\Rightarrow \sum_{i=1}^n y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2} \quad (6)$$

now we know

$$R^2 = 1 - \frac{RSS}{TSS} \Rightarrow 1 - \left(\frac{\sum_{i=1}^n y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2}}{\sum_{i=1}^n y_i^2} \right)$$

$$\Rightarrow \frac{\sum_{i=1}^n y_i^2 - \frac{(\sum x_i y_i)^2}{\sum x_i^2}}{\sum_{i=1}^n y_i^2} \Rightarrow \boxed{\frac{\sum_{i=1}^n (x_i y_i)^2}{(\sum x_i^2) (\sum y_i^2)}}$$

