

Conceptual Questions

① In 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-Sales and Radio-Sales . However, there is a very large p value for newspaper hence no relationship between Newspaper-Sales .

If the p-value is very less , we say p is significant and we reject null-hypothesis . And null hypothesis is that there is no relationship between a predictor and response . Simultaneously we accept Alternative hypothesis .

KNN classifier

Used for classification . For a given x_0 , it selects k points which are closest to x_0 and outputs the label contained by maximum of the k points .

KNN Regression

Used for prediction of quantitative value , it selects k points closest to x_0 and returns the average of all the values .

$$\textcircled{3} \text{ (a) Sales} = \hat{\beta}_0 + \hat{\beta}_1 \text{ GPA} + \hat{\beta}_2 \text{ IQ} + \hat{\beta}_3 \text{ Gender} + \hat{\beta}_4 \text{ GPA} * \text{IQ}$$

$$+ \hat{\beta}_5 \text{ GPA} * \text{Gender}$$

$$\text{Sales (Male)} = \hat{\beta}_0 + \hat{\beta}_1 \text{ GPA} + \hat{\beta}_2 \text{ IQ} + \hat{\beta}_4 \text{ GPA} * \text{IQ}$$

$$\text{Sales (Female)} = \hat{\beta}_0 + \hat{\beta}_1 \text{ GPA} + \hat{\beta}_2 \text{ IQ} + \hat{\beta}_3 + \hat{\beta}_4 \text{ GPA} * \text{IQ} +$$

$$\hat{\beta}_5 \text{ GPA}$$

$$= \hat{\beta}_0$$

$$\text{Sales (Female)} - \text{Sales (Male)}$$

$$= \hat{\beta}_3 + \hat{\beta}_5 \text{ GPA}$$

$$= 35 - 10 \text{ GPA} \rightarrow \textcircled{A}$$

If the eq \textcircled{A} is +ve then female earn more on average than male.

Looking at all conditions -

(iii) when GPA is high enough eq \textcircled{A} is -ve
hence males earn more on average than females provided that GPA is high enough.

$$\text{(b) Sales} = \hat{\beta}_0 + \hat{\beta}_1 \text{ GPA} + \hat{\beta}_2 \text{ IQ} + \hat{\beta}_3 \text{ Gender} + \hat{\beta}_4 \text{ GPA} * \text{IQ}$$

$$+ \hat{\beta}_5 \text{ GPA} * \text{Gender}$$

$$\text{Sales} = 50 + 20 \text{ GPA} + 0.07 \text{ IQ} + 35 \text{ GENDER} +$$

$$0.01 \text{ GPA} * \text{IQ} - 10 \text{ GPA} * \text{GENDER}$$

$$\text{Sales (females)} = 50 + 20(4) + 0.07(110) + 35 + 0.01(4)(110)$$

$$= 10(4)$$

$$= \underline{\underline{137.1}} + 50 = \underline{\underline{187.1}}$$

(c) false, for significance the value of coefficients is not the best choice so, we can't conclude.

(4) (a) complexity \downarrow \rightarrow training error

for training data, RSS decrease with increase in complexity of model. For cubic regression the model will overfit hence for cubic regression RSS will be lower than RSS for linear regression.

(b) for test data, the cubic model will not be able to predict as better as linear regression and RSS for test data will be for linear regression will be lower.

(c) The RSS for cubic regression will be lower. Reason is stated in part (a).

(d) As the model is unknown so it is difficult to define RSS of which model will be lower. As the nature of model becomes more non-linear test RSS for cubic will be lower.

$$⑤ \hat{y}_i = x_i \hat{\beta} \quad (\text{Given}) \quad \text{where,}$$

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

To, Prove

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

Proof -

$$\begin{aligned} \hat{y}_i &= x_i \frac{\sum_{j=1}^n x_j y_j}{\sum_{j=1}^n x_j^2} \\ &= \frac{\sum_{j=1}^n x_i x_j y_j}{\sum_{j=1}^n x_j^2} = c_i' \sum_{i=1}^n x_i y_i \\ &= \sum_{i=1}^n c_i' x_i y_i = \sum_{i=1}^n a_{i'}' y_{i'} \end{aligned}$$

Hence Proved

$$\text{where } a_{i'}' = \sum_{i=1}^n c_i' x_i \quad \text{where } c_i = \frac{x_i}{\sum_{i=1}^n x_i^2}$$

(6) The least square line in case of simple regression is

$$y = \hat{\beta}_0 + \hat{\beta}_1 x \rightarrow ①$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \rightarrow ② \text{ (Given)}$$

Putting ② in ①

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

Now for $x = \bar{x}$

we get.

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

$$y = \bar{y}$$

Hence $y = \hat{\beta}_0 + \hat{\beta}_1 x$ passes through (\bar{x}, \bar{y})

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

$$\begin{aligned} R^2 &= 1 - \frac{\text{RSS}}{\text{TSS}} = 1 - \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \\ &= 1 - \frac{\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2}{\sum_{i=1}^n y_i^2} \quad \text{As } (\bar{y} = 0) \end{aligned}$$

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum x_i y_i}{\sum x_i^2}$$

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

So putting β_1, β_0 in R^2

$$= 1 - \frac{\sum (y_i - \beta_1 x_i)^2}{\sum y_i^2} \quad (\text{Solving ahead})$$

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

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

$$(\text{cor}(x, y))^2 = \frac{\left(\sum x_i y_i \right)^2}{\sum x_i^2 \sum y_i^2} \quad (\text{when } \bar{x} = \bar{y} = 0)$$

$$\text{Hence } R^2 = \underline{\underline{(\text{cor}(x, y))^2}}$$