ANUSHA MUPPALLA

1628631|

ISL HW-2 Assignment.

null hypothesis means, that there is no combination of individual features (or) parameter like TV, radio, newspaper or internet advertisement will influence sales.

$$sales \ c = TV (o) + radio (o) + Newspaper (o) + internet \times o + incpt (\beta_1) y$$

the p value fo each term tests null hypothesis that the coefficient is equal to zero i.e no effect on them.

→ A p value ($<0.05$) indicates that you can reject (null hypothesis)

⊙ A prediction that has lower p value likely to be a meaningful addition to your model because changes in features valu are related to response value

→ from table p values are significant for (TV) $P<0.0001$ and internet ($P<0.0001$) where as it is not significant for radio and newspaper because of large p values so, we can reject TV, internet as they are null hypothesis. and no effect what matters

→ But we cannot reject the radio and newspaper as they do affect the advertising as it is a inference problem, when internet and TV are not going up.

# 5) a.

## kNN classification

1) it is a kNN method using k-Nearest neighbor but the final result i.e o/p of $y$ (response) is **qualitative** i.e

   majority class of k-Nearest neighbor

ii) knowing neighbours as $x_0$ then estimating conditional probability $P(y=j/x=x_0)$ for j class

## kNN regression

1) it is also using k-NN but the final response is **quantitative** one for $f(x)$ i.e as average value of corresponding y's in k-nearest neighbour of X.

ii) knowing neighbours then estimating $f(x_0)$ as the average of all training responses.

parametric approach:-

$$\hat{y} = \beta_0 + \beta_A X_A + \beta_w X_w + \beta_q X_q + \beta_{Aw} \cdot X_{Aw}$$

$$\hat{y} = 50 + 0.5 X_A + (1.0) X_w + (-40) X_q + (0.01) X_{Aw} + (1.0) X_{Aq}$$

where $X_A = 30$, $X_w = 150$, $X_q = male = 0$

$$\hat{y} = 50 + (0.5)(30) + (1.0)(150) - 40(0) + (0.01)(30 \times 150) + (1.0)(30(0))$$

$$\hat{y} = 260$$

b) $\hat{y} = \beta_0 + \beta_A X_A + \beta_w X_w + \beta_q X_q + \beta_{Aw} X_{Aw}$:

same as @ question where $X_q = 1$ (female)

$$X_A = 30,$$
$$X_w = 150$$

$$\hat{y} = 50 + 0.5(30) + (1.0)(150) - 40(1) + (0.01)(30 \times 150) + (1.0)(30)$$

$$\hat{y} = 250$$

c) $X_A = 60$, $X_w = 150$, $X_q = 0$ (male)

$$\hat{y} = 50 + 0.5(60) + (1.0)(150) - 40(0) + 0.01(60 \times 150) + (1.0)(60 \times 0)$$

$$= 50 + 30 + 150 - 0 + 90 + 0$$

$$\hat{y} = 320$$

d) $X_A = 60$, $X_w = 150$, $X_q = 1$ (female)

$$\hat{y} = 50 + 0.5(60) + (1.0)(150) - 40(1) + 0.01(60 \times 150) + (1.0)(60 \times 1)$$

$$\hat{y} = 340$$

e) $X_G = 1, X_G = 0 \implies X_A = ? \quad X_\omega = 150$

$\hat{y}) = 50 + (0.5)(X_A) + (1.0)(150) - 40(1) + 0.01(X_A * 150) + (1.0)(X_A)$

$= 50 + (0.5)(X_A) + (1.0)(150) - 40(0) + 0.01[X_A * 150] + (1.0)(X_A)$

$X_A - 40 = 0$

$$\boxed{X_A = 40}$$

Age of 40 bcm both gender have same blood levels

$y = \hat{\beta_0} + \hat{\beta_1} x + \varepsilon$

without knowing more about training data, it is difficult to know which training RSS is lower b/n linear or cubic but if relationship is consider as linear then least squares line to be close to true regression line.

b) Answer to the previous question smaller than training Rss

To this case, test Rss depends upon test data so as of now we have not enough information to decide. but we may assume that polynomial regression will have higher test Rss as the overfit and more error than linear regression than linear regression

c) if relationship is not linear, I would know cubic regression to better fit the training data due to increased flexibility.
so I can't judge that cubic regression may have some non linearity

d) at is difficult to judge whether non linear or cubic regression will lead to a lower Rss in the data (test) as it depends on the true relationship of X and Y. if it is far from linear, cubic regression will lead to lower Rss in test data and better representation of data generating process

Given table have

$n = 4$ (no of samples)

$b = 5$ (no of groups)

mean of total sum of squares (ss)

$$[T] = \frac{(7+6+8+9+4+3+2+4+11+10+10+10+5+4+5+4+6+5+7+6)}{20}$$

$$= \frac{(30+13+41+18+24)^2}{20} = \frac{(126)^2}{20}$$

$$= 793.8$$

mean b/n groups

$$[B] = \frac{(\Sigma x_1)^2 + (\Sigma x_2)^2 + (\Sigma x_3)^2 + (\Sigma x_4)^2 + (\Sigma x_5)^2}{b-1}$$

$$= \frac{(7+6+8+9)^2 + (4+3+2+4)^2 + (11+10+10+10)^2 + (5+4+5+4)^2 + (6+5+7+6)^2}{4}$$

$$= \frac{3650}{4}$$

$$= 912.5$$

mean within groups.

$$[W] = \frac{7^2+6^2+8^2+9^2+4^2+3^2+2^2+4^2+11^2+10^2+10^2+10^2+5^2+4^2+5^2+4^2+6^2+5^2+7^2+6^2}{15 \; (4-1)}$$

$$= 924$$

| Source | Expression | degree of freedom | SS | mean square Error | F | pvalue |
|---|---|---|---|---|---|---|
| B/n groups | [B]=912·5 | 5-1=4 | [B]-[T] 118·7 | $\frac{SS}{df}$ = 29·675 | 38·74 | ·F-table |
| within groups | [W]=924 | 20-5=15 | [W]-[B] =11.5 | 0·766 | | |
| Total | [T]=793·8 | 20-1=19 | 136·2 (W)-(F) | | | |

pvalue = 0. 00000009712

pvalue < 0.0001

It is significant at P<0.1