

2)



	Coefficient	Std. error	t-statistic	p-value
(Constant)	6.188	0.819	7.556	<.0001
TV	0.008	0.004	2.059	<.0500
Rad	0.149	0.006	25.248	<.0001
House	0.001	0.007	0.138	>.0500

TABLE 8.4. For the interesting data, least squares coefficient estimates of the average house (dependent) number of units sold on TV, radio, and newspaper advertising budgets.

Null hypotheses:

There is no  $R^2$

Simple Regression:  $H_0: \beta_1 = 0$

$H_a: \beta_1 \neq 0$

Multiple Regression:  $H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$

$H_a$ : atleast one  $\beta_j$  is not zero

A. Null Hypo:

- TV has no  $R^2$  w.r.t. Sales
- Rad. has no  $R^2$  w.r.t. sales
- House has no  $R^2$  w.r.t. sales

A. conclusion:

- As the p value for TV is  $> 0.05$  is very less we can say that there is  $\neq$  between those parameter and Sales.
- P (House) is high  $\therefore$  do not reject  $H_0$

2). ANOVA classification:

- get the data
- Tab the  $X_i$
- find manual  $t$  point to case
- tabbed point

for a check and point is:

$$P(V_j | X_j) = \frac{1}{K} \sum_{i=1}^K I(y_i = j)$$



$$\therefore P(y_{\text{true}}(x) = \frac{1}{S} \sum_{i=1}^S I(y_i = j)$$

$$P(y_{\text{pred}}(x) = \frac{1}{S} \sum_{i=1}^S I(y_i = j)$$

ANOVA classification:-

Sum of:

$$\sum (y_i - \bar{y})^2 = \sum \frac{1}{K} \sum_{i=1}^K y_i^2 - \bar{y}^2$$

find diff. bet. points

mean: i.e. loss

median: i.e. loss

3).  $X_1 = 60A$   $X_2 = 70$   $X_3 = 1$   $X_4 = 1$   $X_5 = 1$   $X_6 = 1$   $X_7 = 1$   $X_8 = 1$   $X_9 = 1$   $X_{10} = 1$

$\mu_1 = 20$   $\mu_2 = 0.01$   $\mu_3 = 70$   $\mu_4 = 0.01$   $\mu_5 = 10$   $\mu_6 = 10$   $\mu_7 = 10$   $\mu_8 = 10$   $\mu_9 = 10$   $\mu_{10} = 10$

$Y$  (study salary) in 1000 \$

A:

$$\text{Salary}_{\text{avg}} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots \rightarrow (1)$$

$$= 50 + 20 X_1 + 0.01 X_2 + 35 X_3 + 0.01 X_4 + (-10 X_5)$$

$$\text{Salary}_{\text{avg}} = 50 + 20 X_1 + 0.01 X_2 + 0 + 0.01 X_4 + 0$$

$$\text{Salary}_{\text{avg}} = \text{Salary}_{\text{avg}}$$

$$50 + 20 X_1 + 0.01 X_2 + 35 X_3 + 0.01 X_4 - 10 X_5 = 50 + 20 X_1 + 0.01 X_2 + 0.01 X_4$$

$$35 X_3 - 10 X_5 = 0$$

$$35 X_3 > 10 X_5$$

$$35 X_3 > 10 X_5$$

$$\frac{35}{10} < X_1 \quad \text{or} \quad 3.5 < X_1$$

Only (iii) is correct.

$$(i) \quad 20:100 \quad \text{or} \quad 0.2:1$$

$$\text{Pay} = 1000 \cdot 0$$

$$\# 1000 \cdot 0$$

(ii) The coeff tells that how much effect it has on unit increase

Every unit  $\uparrow$  of  $X_1, X_2$  &  $X_3$

(iii) Do  $S.E$  to check how good our estimate is

$S.E \ll \ll \ll$  reliable parameter

$Pr \ll \ll$  of  $S.E \ll \ll$  effects  $\hat{\beta}$

$Pr \ll \ll$  of  $S.E \gg \gg$  not affecting  $\hat{\beta}$

(iv)

$n = 100$

Linear:  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \epsilon$

Cubic:  $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \epsilon$

(a)  $X$  to  $Y$  is linear

$RSS_{\text{train}}(n)$   $RSS_{\text{train}}(n)$  what would be  $R^2$ ?

From my learning:

$RSS_{\text{train}}$  will be less than linear because

$$RSS_{\text{train}} = e_1^2 + e_2^2 + e_3^2$$

Sum: error of train data

Var: error of test data

While testing linear would have High  $R^2$  & High loss

During cubic would have Low loss & High loss

$$RSS_{\text{train}} > RSS_{\text{train}}$$

(b) For testing:

Cubic - overfit & Linear - underfit

$$RSS_{\text{linear}} < RSS_{\text{cubic}}$$

$\therefore$  explains more variance

(c)

$RSS_{\text{train}} \downarrow$  as the model is of high variance

$RSS_{\text{train}} \ll \ll$  New model

(d)

For test we cannot say because if the new model is like

more pay the variance  $\uparrow$  (i)