

## ASSIGNMENT II

Date \_\_\_\_\_  
Page \_\_\_\_\_

Ques 1) Modern email servers and anti-spam filter attempt to identify spam email and direct them to a junk folder. There are various ways to detect spam and research still continues. In this regard, an information security officer tries to confirm that the chance for an email to be spam depends on whether it contains image or not. The following data were collected on  $n = 1000$  random email messages.

spam status	Image Containing status			Total
	With image	No image		
With spam	160	240		400
No spam	140	460		600
Total	300	700		1000

Assess whether being spam and containing images are independent factors at 1% level of significance.

→ Soln, (Using Chi-Square Test for independence)

- Problems to test:

$H_0$ : Being spam and containing images are independent factor.

$H_1$ : Being spam and containing images are dependent factors.

- Expected Frequency:

~~Spam~~ formula:  $Ef = (\text{row total} * \text{column total}) / \text{total emails}$

status	observed frequency(OE)	Expected Frequency
Spam with image	160	$400 \times 300 / 1000 = 120$
Spam without image	240	$400 \times 700 / 1000 = 280$
No spam with image	140	$600 \times 300 / 1000 = 180$
No spam without image	460	$600 \times 700 / 1000 = 420$

- Calculating Chi-square ( $\chi^2$ ):

$$\text{Chi-Square } (\chi^2) = \sum \frac{(O_F - E_F)^2}{E_F}$$

$$\begin{aligned} \chi^2 &= (160 - 120)^2 / 120 + (240 - 280)^2 / 280 + (140 - 180)^2 / 180 + \\ &\quad \cancel{(460 - 420)^2 / 420} \\ &= 13.33 + 5.7142 + 8.88889 + 3.8095 \\ &= 31.746 \end{aligned}$$

- Degree of Freedom (df):

$$df = (\text{no of rows} - 1) * (\text{no of columns} - 1)$$

$$df = (2 - 1) * (2 - 1) = 1$$

- P-value:

$$\alpha = 0.01 (1\%)$$

Here

$$\chi_{\alpha}^2 < \chi_{\text{calc}}^2 \quad \text{in the degree of freedom 1}$$

i.e.,  $6.835 < 31.746$

- Decision:

The calculated chi-square is greater than the table value for  $df = 1$  and  $\alpha = 0.01$

hence, we reject  $H_0$  at 1% level of significance level.

- Conclusion:

Hence the presence of image is not independent of whether an email is spam. There is a strong relation ship between the two factors.

No 2) The following data related to the number of children classified according to the type of feed and the nature of teeth.

Type of feed	Nature of teeth.	
	Normal	Defective
Breast	18	12
Bottle	2	13

Do the information provide sufficient evidence to conclude that type of feeding and nature of teeth are dependent?  
Use Chi-Square test at 5% level of significance.

→ Soln,

Using Chi-square Test for independence:

- Problem to test:

$H_0$ : There is no association

$H_1$ : There is association.

- Data

Type of feed	Normal teeth	Defective teeth	Total
Breast	18	12	30
Bottle	2	13	15
Total	20	25	45

Level of Significance:  $\alpha = 0.05$  (5%)

- Expected Frequency (EF):

Formula:  $EF = (\text{row total} * \text{Column total}) / \text{Total}$

Type	Observed Frequency (OF)	Expected Frequency (EF)
Brest, Normal	18	$(20 \times 30)/45 = 13.33$
Breast, Defective	12	$(30 \times 25)/45 = 16.67$
Bottle, Normal	2	$(15 \times 20)/45 = 6.67$
Bottle, Defected	13	$(15 \times 25)/45 = 8.33$

- Calculating chi-Square ( $\chi^2$ ):

$$\begin{aligned}
 \chi^2 &= \sum \frac{(OF - EF)^2}{EF} \\
 &= \frac{(18 - 13.33)^2}{13.33} + \frac{(12 - 16.67)^2}{16.67} + \frac{(2 - 6.67)^2}{6.67} + \frac{(13 - 8.33)^2}{8.33} \\
 &= 1.636 + 1.30827 + 3.2697 + 2.618115 \\
 &= \underline{\underline{8.832085}}
 \end{aligned}$$

- Degree of freedom (df)

$$\begin{aligned}
 df &= (\text{no of rows} - 1) * (\text{no of column} - 1) \\
 &= 2-1 \times 2-1 = 1
 \end{aligned}$$

- P-value:

$$\alpha = 0.05 (5\%)$$

Here

$$\boxed{\chi^2_{\alpha} < \chi^2_{\text{calc}}}
 \quad \text{in the } (df = 1)$$

i.e.  $3.841 < 8.832$

- Decision:

Since  $\chi^2(\text{calculated})$  is greater than the table value,  
We reject  $H_0$  at 5% level of significance.

### ~~Ques.~~ Conclusion:

Hence there is an association between type of feeding and the nature of teeth.

Qno3) Social media users use a variety of devices to access social networking, mobile phones are increasingly popular. However, is there a difference in the various age groups in the proportion of social media users who use their mobile phone to access social ~~media~~ networking? A study showed the following results for the different age groups.

	Age	Age	Age
Use mobile phones to access social networking	18-34	35-64	65+
yes	60	37	14
No	40	63	86

At the 0.05 level of significance, is there evidence of a different age group with respect to use of mobile phone for accessing social networking?

→ So in,

Using Chi-Square Test.

### • Problem to Test

$H_0$ : Evidence of different age group

$H_0$ : ~~enough~~ No Evidence of different age group.

- **Datas:**

Age group	Yes	No	Total
18-34	60	40	100
35-64	37	63	100
65+	14	86	100
Total	111	189	300

Level of significance  $\alpha = 0.05$

- **Expected frequency (EF):**

Formula:  $EF = (\text{Total row} \times \text{Column Total}) / \text{all total}$

Age	Observed frequency (OF)	Expected frequency (EF)
18-34 (yes)	60	$(100 \times 111) / 300 = 37$
18-34 (No)	40	$(100 \times 189) / 300 = 63$
35-64 (yes)	37	$(100 \times 111) / 300 = 37$
35-64 (No)	63	$(100 \times 189) / 300 = 63$
65+ (yes)	14	$(100 \times 111) / 300 = 37$
65+ (No)	86	$(100 \times 189) / 300 = 63$

- **Chi-square ( $\chi^2$ ):**

$$\chi^2 = \sum \frac{(OF - EF)^2}{EF}$$

$$= 14.29729 + 8.396825 + 0 + 0 + 14.29729 + 8.396825 \\ = 45.3882309$$

- Degree of freedom (Df)

$$df = (\text{row} - 1) * (\text{Column} - 1)$$

$$= (3 - 1) \times (2 - 1)$$

$$= 2$$

- P-value:

$$\alpha = 0.05$$

Here

$$1 X^2_{\alpha} < X^2_{\text{calc}}$$
 in  $df = 2$

$$\text{i.e., } 5.991 < 15.3882$$

- Decision:

Since  $X^2$  is greater than the table value at  $\alpha = 0.05$  level significance, ~~Accept H<sub>0</sub>~~. Reject H<sub>0</sub>.

- Conclusion:

Hence, there is sufficient evidence to suggest that there is a different age group in the proportion of social media users who use their mobile phone to access social networking at 0.05 level of significance.

Qno 4) A random sample of 200 married men, all retired, were classified according to education and number of children.

Education	No. of children			Total
	0-1	2-3	Over 3	
Elementary	14	37	32	83
Secondary	19	42	17	78
College	12	17	10	39
Total	45	96	59	200

Test the hypothesis, at the 1% level of significance, that the number of children is independent of the level of education attained by the father.

→ Soln,

Using Chi-Square Test

- Problem to test:

$H_0$ : The no of children is independent of level of education attained by the father.

$H_1$ : The no of children is dependent of level of education attained by the father.

- Expected frequency ( $E_f$ )

Formula:

$$E_f = (\text{Row total} \times \text{Column total}) / \text{all total}$$

Age Education	Observed frequency (OF)	Expected frequency (EF)
Elementary (0-1)	14	$(83 \times 45) / 200 = 18.875$
Elementary (2-3)	37	$(83 \times 96) / 200 = 39.84$
Elementary (3+)	832	$(83 \times 59) / 200 = 24.485$
Secondary (0-1)	19	$(78 \times 45) / 200 = 17.55$
Secondary (2-3)	42	$(78 \times 96) / 200 = 37.44$
Secondary (3+)	17	$(78 \times 59) / 200 = 27.01$
College (0-1)	12	$(39 \times 45) / 200 = 8.775$
College (2-3)	17	$(39 \times 96) / 200 = 18.72$
College (3+)	10	$(39 \times 59) / 200 = 11.505$

- Chi square ( $\chi^2$ ):

$$\chi^2 = \sum \left[ \frac{(OF - EF)^2}{EF} \right]$$

$$= 1.170314 + 0.202448 + 2.306523 + 0.1198 + 0.55538 \\ + 3.70974 + 1.185256 + 0.158034 + 0.196873 \\ = 9.60436$$

- Degree of freedom (df):

$$df = (\text{row} - 1) \times (\text{column} - 1) \\ = 2 \times 2 = 4$$

- P-value

$$\alpha = 0.01 \quad (1\%)$$

Here

$$\boxed{\chi^2 \alpha < \chi^2_{\text{cal}}} \quad \text{in } df = 4 \\ \text{i.e., } 13.277 > 9.60436$$

• Decision:

Hence, we accept  $H_0$  at 1% level of significance.

• Conclusion:

Hence, the no. of children is independent of level of education attained by the father.

Qn05) A psychologist wishes to verify that a certain drug increases the reaction time to given stimulus. The following reaction times (in tenth of seconds) were recorded before and after injection of the drug for each of four subjects.

Subject	1	2	3	4
Reaction time Before	7	2	12	12
After	13	3	18	13

→ ~~Soln~~, Test at 5% level of significance to determine whether the drug significantly increases reaction time.  
Use non parametric test.

→ Soln

performing non-parametric (Wilcoxon Signed Rank Test)

- Problems to test

$H_0$ : Drug has no effect on reaction time.

$H_1$ : The drug increases reaction time. (after reaction time greater)

- Data

Subject	Before RT	After RT	Difference (After - before)
1	7	13	6
2	2	3	1
3	12	18	6
4	12	13	1

Significance level  $\alpha = 0.05$  (5%)

Sample size( $n$ )= 4

- Ranking :

Subject	Difference	Rank
1	6	4
2	1	1
3	6	4
4	1	1

→ Here we only consider +ve ranks ( $T+$ ) as only increasing or not is asked.

- Calculating  $T+$ :

$$\begin{aligned} T+ &= \text{Rank(Sub1)} + \text{Rank(Sub3)} \\ &= 4 + 4 \\ &= 8 \end{aligned}$$

- Wilcoxon Signed-Rank statistic (W):

$$\begin{aligned} W &= T+ (\text{for increasing}) \\ &= 8 \end{aligned}$$

- Critical value:

$$\alpha = 0.05$$

for two-tail and 4 sample size,  $W_{\alpha} = 0$

- Decision

$$W_{0.05} < W_{\text{calc}}$$

$$\text{i.e. } 0 < 8$$

hence, we reject  $H_0$  at 0.05 level of significance.

- Conclusion:

Hence, we have strong evidence to suggest that the drug significantly increases reaction time.

- Conclusion :

Hence, we have strong evidence to suggest that the drug significantly increases reaction time.

(Qno 6) What do you mean by ~~para~~ non-parametric test? write down advantages of non-parametric tests over parametric tests.

- Non-parametric test, also known as distribution-free tests, are statistical method used in hypothesis testing that do not make assumptions about the frequency distribution of the variables being evaluated.
- It's advantages are:
  - a) Can be applied to qualitative (rank, ordinal, categorical data) as well as quantitative data.
  - b) It is the only possible test for respectively small sample.
  - c) It is simple to understand, quicker and easier to apply
  - d) It is less time consuming.
  - e) It needs no assumption about the population from which sample size is selected.
  - f) It has greater rank of applicability because of milder assumption.
  - g) It does not require complicated Sampling theory.

Qno7) Bank of Nepal recorded the sex of first 30 customers who appeared last Monday with notation MMFMMF MFFMMMF MFFMFF MFFMFFM FMMMF FF. At the 0.05 level of significance, test the randomness of this sequence.

→ Soln,

$$\text{no of } M(n_1) = 15$$

$$\text{no of } F(n_2) = 15$$

$$\text{no of runs (r)} = 16$$

$$\text{level of significance } (\alpha) = 0.05$$

- Problem to test:

$H_0$ : The sequence is random order

$H_1$ : The sequence isn't in random order.

- Test statistic

$$\rightarrow r = 16$$

$$Mr = \frac{2n_1 n_2}{n_1 + n_2} = \frac{2 \times 15 \times 15}{30} = 15$$

$$\sigma_r^2 = \frac{2n_1 n_2 (2n_1 n_2 - n_1 - n_2)}{(n_1 + n_2)^2 (n_1 + n_2 - 1)} = \frac{2 \times 15 \times 15 (2 \times 15 \times 15 - 15 - 15)}{(15+15)^2 \times (15+15-1)}$$

$$\therefore \sigma_r = \sqrt{4.2413} = 2.0909$$

$$Z = \frac{r - Mr}{\sigma_r} = \frac{16 - 15}{2.0909} = 0.3716$$

• Critical value:

At  $\alpha = 0.05$ , level of significance critical value is,

$$Z_{\text{tabulated}} = 2\alpha$$

• Level of significance

$$\text{At } \alpha = 0.05, (\underline{r}, \bar{r}) = (10, 22), r = 16$$

• Decision

Accept  $H_0$  at  $\alpha = 0.05$ . as  $r \in (\underline{r}, \bar{r})$  i.e  $16 \in (10, 22)$

• Conclusion:

Hence, the sample are in random order.

Qno 8) Define level of significance. Describe run test with some relevant examples.

→ ~~the~~ Level of significance, denoted by  $\alpha$  (alpha), is a critical value used in the hypothesis testing to determine the cutoff point for rejecting the null hypothesis. It represents the maximum probability of making type I error, which occurs when the null hypothesis is incorrectly rejected when it is actually true.

→ Run test is a ~~parametric~~ non-parametric test used to determine the randomness of the selected sample.

Run is a set of identical or related symbols contained between two different symbols or none at all.

Let us consider a random sample of size  $n$  is selected.

Let  $x_1, x_2, x_3, x_4, \dots, x_n$  be samples.

Step 1) Problem to test

$H_0$ : In random order

$H_1$ : ~~is~~ Not in random order.

Step 2) Finding median of sample and assign a symbols, number of observations of symbol A denoted by  $n_1$ , and number of observations of symbol B denoted by  $n_2$ .

Step 3) If,

~~For small sample ( $n_1, n_2 \leq 20$ )~~

Test statistic

Number of runs (r)

~~Step 2~~) level of significance.

$\alpha = 0.05$  is taken unless we are given.

~~Step 3~~) Critical value:

$r$  and  $\bar{r}$  are obtained from table according to level of significance, degree of freedom  $n_1$  and  $n_2$ , and alternative hypothesis.

Decision

Accept  $H_0$  at  $\alpha$  level of significance if  $r \in (l, r)$ , reject else.

Step 3) Else,

For large sample size.

In this case ~~large sample~~  $r$  is approx normally distributed with mean  $M_r = \frac{2n_1 n_2}{n_1 + n_2} + 1$

$$\text{And variance } \sigma_r^2 = \frac{2n_1 n_2 (2n_1 n_2 - n_1 - n_2)}{(n_1 + n_2)^2 (n_1 + n_2 - 1)}$$

Test statistic

$$z = \frac{r - M_r}{\sigma_r} \sim N(0, 1)$$

level of significance

$\alpha = 0.05$  is taken unless we are given

Critical value

$Z$  is obtained from table for level of significance, and alternative hypothesis.

Decision.

Reject  $H_0$  if  $|z| > z_{\text{tabulated}}$ . accept otherwise

Note For  $\alpha$  other than 5% even if sample size is small use Z test in which  $Z = |r - M_r| / 0.5$   
Or.

(Qno 9) What do you mean by run? Marks secured by a sample of 15 students in final exam of statistics II are found to be 27, 34, 46, 21, 7, 56, 44, 32, 25, 42, 33, 48, 28, 41, 5, 49. Are marks in random order? use 5% level of significance.

→ Run is a set of identical or related symbols contained between two different symbols or none at all.

→ So 15,

$$n = 15$$

Arranging data in ascending order

5, 7, 21, 25, 27, 28, 32, 33, 34, 41, 42, 44, 48, 49, 56

$$\text{median} = \left( \frac{n+1}{2} \right)^{\text{th}} \text{ item}$$

$$= \frac{15+1}{2} = 8^{\text{th}} \text{ item}$$

$$\therefore \text{median} = 33$$

Now number greater than 33 be B & less than be A.

A B B A A B B A A B B A B A B

no of A ( $n_1$ ) = 7

no of B ( $n_2$ ) = 8

no of runs ( $r$ ) = 9

Problem to test

$H_0$ : The sequence is random.

$H_1$ : The sequence isn't random.

level of significance.

$\alpha = 0.05$ ,  $(\underline{r}, \bar{r}) = (4, 13)$ ,  $r = 9$

- Decision:

Accept  $H_0$  at  $\alpha = 0.05$  as  $r \in (\underline{r}, \bar{r})$  i.e.  $9 \in (4, 13)$

- Conclusion:

Hence, the sample are in random order.

Qno 10) What is median test? Following data represents marks secured by students of section A and section B of a college in mid-term exam of statistics II.

Section A	30	27	19	22	28	25	9	13	20
Section B	24	28	16	22	19	29	7	11	

Is there any significant difference in marks of section A and section B? Use median test at 5% level of significance.

→ Median Test is the non-parametric test used to test the significance difference between two independent distribution.

→ Soln  
arranging combined data in ascending order

7 9 11 13 16 19 19 20 22 22 24 25 27 28 28 29 30

$$n = 18$$

$$\text{median } M_d = \left( \frac{n+1}{2} \right)^{\text{th item}}$$

$$= \frac{18}{2}^{\text{th item}}$$

$$= \cancel{9}^{\text{th item}} \text{ } 9^{\text{th item}}$$

$$\therefore \text{median} = 22$$

no of obs in 2<sup>nd</sup> sample less or equal to median ( $\alpha \leq 5\%$ )  
 1<sup>st</sup> sample size ( $n_1$ ) = 9  
 2<sup>nd</sup> sample size ( $n_2$ ) = 8

$$K = \frac{n_1 + n_2}{2} = \frac{9+8}{2} = 8.5 \approx 8$$

let  $Md_1$  and  $Md_2$  be median of 1 population and 11 population respectively.

- Problem to test

$$\begin{aligned} H_0: Md_1 = Md_2 & \quad \left. \begin{array}{l} \text{Two tailed} \\ \text{ } \end{array} \right. \\ H_a: Md_1 \neq Md_2 & \end{aligned}$$

- Test statistics

$$P(A=a) = \frac{c(n_1, a) c(n_2, K-a)}{c(n_1 + n_2, K)} = \frac{c(9, a) c(8, 8-a)}{c(17, 8)}, a=0, 1, 2, 3, 4, 5, 6, 7$$

- Critical value:

$$P = P(A \geq a) = P(A \geq 5)$$

$$\begin{aligned} &= \frac{\sum c(9, a) c(8, 8-a)}{c(17, 8)} \\ &= \frac{c(9, 5) c(8, 3)}{c(17, 8)} + \frac{c(9, 6) c(8, 2)}{c(17, 8)} + \frac{c(9, 7) c(8, 1)}{c(17, 8)} \\ &\approx 0.398849 \end{aligned}$$

- Decision:

$$2P = 2 \times 0.398849 = 0.7984 > \alpha = 5\% = 0.05 \text{. Accept } H_0$$

- Conclusion:

There is no significant difference in marks of section A and section B

Qno 11) Two computer makers, A and B, compete for a certain market. Their users rank the quality of computers on a 4-point scale as "Not satisfied", "Satisfied", "Good quality", and "Excellent quality, will recommend to others. The following counts were observed.

Computer maker	Not satisfied	Satisfied	Good quality	Excellent quality
A	20	40	70	20
B	10	30	40	20

Is there a significant difference in customer satisfaction of the computers produced by A and by B using Mann-Whitney U test at 5% level of significance.

→ Soln

A	Ranks	B	Ranks
20	3.5	10	1
40	6.5	30	5
70	8	40	6.5
20	3.5	20	3
5	<del>8.5</del> R <sub>1</sub> = 20.5		R <sub>2</sub> = 15.5

Sample size of A = 4 (n<sub>1</sub>)

Sample size of B = 4 (n<sub>2</sub>)

Sum of Ranks of A (R<sub>1</sub>) = 20.5

Sum of Ranks of B (R<sub>2</sub>) = 15.5

$$U_1 = n_1 n_2 + \frac{n_1(n_1+1)}{2} - R_1 = 4 \times 4 + \frac{(4)(4+1)}{2} - 20.5 \\ = \cancel{16} 5.5$$

$$U_2 = n_1 n_2 - U_1 = 4 \times 4 - 5.5 \\ = 10.5$$

$$U_0 = \min \{U_1, U_2\} = 5.5$$

Let  $Md_1$  and  $Md_2$  be median ~~mean~~ of A and B resp.

- Problem to test:

$$H_0: Md_1 = Md_2$$

$$H_1: Md_1 \neq Md_2$$

- Test statistics,

$$U_0 = 5.5$$

- Critical value

$$\text{at } \alpha = 0.05, P = 0 \text{ (critical value at } (4,4))$$

- Decision.

$P < U_0$  i.e.,  $U_{0.05}(4,4) = 0 < U_0 = 5.5$ . Accept  $H_0$  at 5% level of significance.

- Conclusion:

There is no significant difference in customer satisfaction of computer produced by A & B.

At  $\alpha = 5\% = 0.05$ , critical value is  $U_{\alpha}(n_A, n_B) = U_{0.05}(4, 4) = 0$

\* Decision

$U_0 = 5.5 > U_{0.05}(4, 4) = 0 \Rightarrow$  Accept  $H_0$  at 0.05 level of significance.

\* Conclusion

There is no significant difference in customer satisfaction of computers produced by A and B.

Q12 A chemist uses 3 catalyst for distilling alcohol and layout were tabulated below

Catalyst	Alcohol (in cc)				
C <sub>1</sub>	380	430	410		
C <sub>2</sub>	290	350	270	250	270
C <sub>3</sub>	400	380	450		

Are there any significant differences between catalyst? Test at 5% level of significance. Use Kruskal Wallis H test.

SOLN:

We are given

Catalyst	Alcohol (in cc)					R
C <sub>1</sub>	380	430	410			25
Rank C <sub>1</sub>	6.5	10	9			25
C <sub>2</sub>	290	350	270	250	270	
Rank C <sub>2</sub>	4	5	2.5	1	2.5	15
C <sub>3</sub>	400	380	450			
Rank C <sub>3</sub>	8	6.5	11			25

here,

Sample size of C<sub>1</sub> = n<sub>1</sub> = 3

" " C<sub>2</sub> = n<sub>2</sub> = 5

" " C<sub>3</sub> = n<sub>3</sub> = 3

$$n = n_1 + n_2 + n_3 = 3 + 5 + 3 = 11$$

And,

$t_i$  :- No. of observations having same rank  $[t_i : 1, 2, 3, \dots]$

i.e.  $t_1 = 2$ ,  $t_2 = 2$  [for 270 and 380].

\* Problem to test

$$H_0: Md_1 = Md_2 = Md_3$$

$$H_1: \text{At least one } Md_i \text{ is different}$$

$[i = 1, 2, 3]$

$$\begin{cases} Md_1: \text{Median of alcohol for C}_1 \text{ (in cc)} \\ Md_2: " " C_2 " \\ Md_3: " " C_3 " \end{cases}$$

\* Test statistic

$$H = \frac{12}{n(n+1)} \left[ \sum_i \frac{R_i^2}{n_i} \right] - 3(n+1)$$

$$1 - \frac{\sum (t_i^3 - t_i)}{n^3 - n}$$

$$= \frac{12}{11(11+1)} \left[ \frac{25.5^2}{3} + \frac{15^2}{5} + \frac{25.5^2}{3} \right] - 3(11+1)$$

$$1 - \frac{(2^3 - 2) + (2^3 - 2)}{11^3 - 11}$$

$$= 7.5688$$

\* Critical value

At  $n_1 = 3$ ,  $n_2 = 5$  and  $n_3 = 3$  and  $H = 7.5688$ , critical value is

$$P_{\phi} = 0.009$$

\* Decision

$$P = 0.009 \Rightarrow <\alpha = 5\% = 0.05$$

$\Rightarrow$  Reject  $H_0$  at 5% level of significance

\* Conclusion

There is significant difference between catalyst.

Q13 There are 3 brands of computers Dell, Lenovo and HP. The following are lifetime of 15 computers in years

Serial No.	Computer brand	(lifetime in years)
1	Dell	15
2	Lenovo	10
3	HP	9
4	Dell	12
5	Lenovo	6
6	HP	7
7	Dell	4
8	Lenovo	8
9	HP	13
10	Dell	11
11	HP	5
12	Lenovo	7
13	Dell	3
14	HP	5
15	Lenovo	4

Apply appropriate statistical test to identify whether the average life of time (in years) is significantly different across 3 brands of computers at 5% level of significance. You can again tabulate data initially in required format for statistical analysis.

SOLN:

We are given

$$\text{No. Sample size of Dell brand } (n_1) = 5$$

$$\text{ " " Lenovo " } (n_2) = 5$$

$$\text{ " " HP " } (n_3) = 5$$

$$\therefore n = n_1 + n_2 + n_3 = 15$$

Computer brands

	Lifetime (in yrs)					Rank (R <sub>i</sub> )
Dell	15	12	9	11	3	
Rank for Dell (R <sub>1</sub> )	15	13	2.5	12	1	43.5
Lenovo	10	6	8	7	4	
Rank for Lenovo (R <sub>2</sub> )	11	6	9	7.5	2.5	36
HP	9	7	13	5	5	
Rank for HP (R <sub>3</sub> )	10	7.5	14	4.5	4.5	40.5

Here,

$$R_1 = 43.5, R_2 = 36, R_3 = 40.5$$

$$t_1 = 2(4, 4), t_2 = 2(5, 5), t_3 = 2(7, 7)$$

\* Problem to test

$$H_0: M_{d1} = M_{d2} = M_{d3}$$

$$H_1: \text{At least one } M_{di} \text{ is different. } [i: 1, 2, 3]$$

\* Test statistic

$$H = \frac{\frac{12}{n(n+1)} \sum_i \frac{R_i^2}{n_i} - 3(n+1)}{1 - \frac{\sum (t_i^3 - t_i)}{n^3 - n}}$$

$$= \frac{\frac{12}{15(15+1)} \left[ \frac{R_1^2}{n_1} + \frac{R_2^2}{n_2} + \frac{R_3^2}{n_3} \right] - 3(15+1)}{1 - \frac{(t_1^3 - t_1) + (t_2^3 - t_2) + (t_3^3 - t_3)}{15^3 - 15}}$$

$$= \frac{\frac{12}{15 \times 16} \left[ \frac{43.5^2}{5} + \frac{36^2}{5} + \frac{40.5^2}{5} \right] - 3 \times 16}{1 - \frac{(2^3 - 2) + (2^3 - 2) + (2^3 - 2)}{15^3 - 15}}$$

$$= 0.2865$$

\* Critical value

At  $\alpha = 5\% = 0.05$ , critical value is  $\chi^2_{\alpha(k-1)} = \chi^2_{0.05(3-1)} = \chi^2_{0.05} = 5.991$

\* Decision

$$H = 0.2865 < \chi^2_{\alpha(k-1)} = 5.991$$

$\Rightarrow$  Accept  $H_0$  at 5% level of significance

\* Conclusion

Average lifetime is not significantly different across 3 brands of computers at 5% level of significance.

Q14 Marks secured by 8 students in 3 chapter tests in a subject are as follows:

Student	A	B	C	D	E	F	G	H
test I	13	11	16	19	6	14	18	5
test II	14	10	18	11	12	9	18	7
test III	15	19	13	10	11	5	17	4

Is there any significant difference in marks in 3 chapter tests? Use Friedman's two way ANOVA test at 10% level of significance.

SOLN:

We are given

Student	A	B	C	D	E	F	G	H	R <sub>i</sub>
test I	13	11	16	19	6	14	18	5	
Rank I	1	2	2	3	1	3	2.5	2	16.5
test II	14	10	18	11	12	9	18	7	
Rank II	2	1	3	2	3	2	2.5	3	18.5
test III	15	19	13	10	11	5	17	4	
Rank III	3	3	1	1	2	1	1	1	19.3

Here,

$$\text{No. of samples } (k) = 3$$

$$\text{Total no. of samples } \Rightarrow 8+8+8 = 24 = n \quad \begin{cases} \text{Sample size of} \\ \text{each test } (n) = 8 \end{cases}$$

$$\text{Tied cases} = 2 \quad [18 \text{ and } 18 \text{ in G}] \Rightarrow t_1 = 2$$

\* Problem to test

$$H_0: M_{dI} = M_{dII} = M_{dIII}$$

$$H_1: \text{At least one } M_{di} \text{ is different } [i = I, II, III]$$

\* Test statistic

$$F_r = \frac{\frac{1}{12} \left[ \sum_i R_i^2 \right] - 3n(k+1)}{nk(k+1)}$$

$$1 - \frac{\sum t_i^2 - t_1^2}{n(k^2 - k)}$$

i.e.

$$Fr = \frac{12}{8 \times 3(3+1)} \left[ 16.5^2 + 18.5^2 + 13^2 \right] - 3 \times 8 \times (3+1)$$

$$= \frac{(2^3 - 2)}{8 \times (3^3 - 3)}$$

$$= 2$$

## \* Critical Value

From Friedman table, at  $k=3$ ,  $n=8$  and  $Fr=2$ , critical value is  $p \approx$

## \* Decision

$p > \alpha = 0.05$ . Accept  $H_0$  at 5% level of significance.

## \* Conclusion

There is no significant difference in marks in 3 chapter tests.

Q15 It was reported somewhere that children whenever plays game in computer, they used the computer very roughly which may reduce the lifetime of computer. The random access memory (RAM) of a computer also plays a crucial role on the lifetime of a computer. A researcher wanted to examine how the lifetime of a personal computer which is used by children is affected by the time (in hours) spends by children per day to play games and available random access memory (RAM) measured in mega bytes (MB) of a used computer. Data is provided in the following table:

Lifetime (yrs)	5	1	7	2	3	4	6
Play time (hrs/day)	2	8	1	5	6	3	2
RAM (MB)	8	2	6	3	2	4	7

Date \_\_\_\_\_  
Page \_\_\_\_\_

Identify which one is dependent variable? Solve this problem using multiple linear regression model and provide problem specific interpretation based on the regression model developed.

Soln:

→ Lifetime (in years) is dependent variable =  $y$  (say)

Play time ( $x_1$ ) [say] and RAM ( $x_2$ ) [say] are independent variables.

And,

Required linear regression model is

$$y = b_0 + b_1 x_1 + b_2 x_2 \quad \textcircled{*}$$

Here,

$y$	$x_1$	$x_2$	$x_1 y$	$x_2 y$	$x_1^2$	$x_2^2$	$x_1 x_2$
5	2	8	10	40	4	64	16
1	8	2	8	2	64	4	16
7	1	6	7	42	1	36	6
2	5	3	10	6	25	9	15
3	6	2	18	6	36	4	12
4	3	4	12	16	9	16	12
6	2	7	12	42	4	49	14

where,

$$\sum y = 28$$

$$\sum x_2 y = 154 \quad n = 7$$

$$\sum x_1 = 27$$

$$\sum x_1^2 = 143$$

$$\sum x_2 = 32$$

$$\sum x_2^2 = 182$$

$$\sum x_1 y = 77$$

$$\sum x_1 x_2 = 91$$

Now,

To fit  $\textcircled{*}$ , we solve the following

$$\rightarrow \sum y = n b_0 + b_1 \sum x_1 + b_2 \sum x_2$$

$$\text{or, } 28 = 7 b_0 + 27 b_1 + 32 b_2 \quad \textcircled{1}$$

$$\rightarrow \Sigma x_1 y = b_0 \Sigma x_1 + b_1 \Sigma x_1^2 + b_2 \Sigma x_1 x_2$$

$$\text{or, } 77 = 27b_0 + 143b_1 + 91b_2 - \textcircled{II}$$

$$\rightarrow \Sigma x_2 y = b_0 \Sigma x_2 + b_1 \Sigma x_1 x_2 + b_2 \Sigma x_2^2$$

$$\text{or, } 154 = 32b_0 + 91b_1 + 182b_2 - \textcircled{III}$$

Then,

Using Cramer's rule

coeff. of $b_0$	coeff. of $b_1$	coeff. of $b_2$	constant
-----------------	-----------------	-----------------	----------

7	27	32	28
---	----	----	----

27	143	91	77
----	-----	----	----

32	91	182	154
----	----	-----	-----

i.e.

$D$ :	7	27	32
	27	143	91
	32	91	182

$D_0$ :	28	27	32
	77	143	91
	154	91	182

$D_{x_1}$ :	7	28	32
	27	77	91
	32	154	182

$D_{x_2}$ :	7	27	28
	27	143	77
	32	91	154

And,

$$b_0 = \frac{D_0}{D} = 6.9613$$

$$b_1 = \frac{D_1}{D} = -0.7854$$

$$b_2 = \frac{D_2}{D} = 0.0149$$

50.

Q) becomes

$$y = 6.9613 + (-0.7854) b_0 x_1 + 0.0149 b_2 x_2$$

$$\text{i.e. } y = 6.9613 - 0.7854 x_1 + 0.0149 x_2$$

Now,

$b_0 = 6.9613$  means that the lifetime of computer is 6.9613 years when play time is 0 hours/day and 0 MB RAM is used.

$b_1 = -0.7854$  means that lifetime of computer decreases by 0.7854 years when play time increases by 1 hours/day and constant RAM is used.

$b_2 = 0.0149$  means that the lifetime of computer increases by 0.0149 years when RAM usage increases by 1 MB at constant playtime.

Q16

What are required conditions for error variable in multiple regression analysis? The Internal Revenue Service is trying to estimate the monthly amount of unpaid taxes discovered by its auditing division. The Internal Revenue Service estimated this figure on the basis of field auditing labour hours and no. of hours its computers are used. The table given below presents these data for the last 10 months:

Month	Field audit labour hours in 100 ( $x_1$ )	Computer hours in 100 ( $x_2$ )	Annual unpaid taxes discovered million of dollars ( $y$ )
Jan	45	16	29
Feb	42	14	24
Mar	44	15	27
Apr	45	13	25
May	43	12	26

Jun	46	14	28
Jul	44	16	30
Aug	45	16	28
Sep	44	15	28
Oct	43	15	27

Given,

$$\sum y = 12005, \sum y x_1 = 4013, \sum x_1 x_2 = 6485,$$

$$\sum y^2 = 7428, \sum x_1^2 = 19461, \sum x_2^2 = 2173$$

- Develop the estimating equation best describing these data.
- Interpret the value of regression coefficients.
- Estimate the annual unpaid tax for field audit labour 4200 hours and computer 1600 hours.

SOLN:

The required linear regression equation is

$$y = b_0 + b_1 x_1 + b_2 x_2 - \textcircled{1}$$

where,

$y$  = amount of unpaid taxes discovered (in million £)

$x_1$  = field audit labour hours (in 100)

$x_2$  = computer hours (in 100)

Now,

To fit  $\textcircled{1}$  we use the following equations

$$\sum y = n b_0 + b_1 \sum x_1 + b_2 \sum x_2$$

where,

$$\sum y = 272, \sum x_1 = 441, \sum x_2 = 147, n = 10$$

i.e.

$$272 = 10 b_0 + 441 b_1 + 147 b_2 - \textcircled{1}$$

And,

$$\sum y x_1 = b_0 \sum x_1 + b_1 \sum x_1^2 + b_2 \sum x_1 x_2$$

where,

$$\sum y x_1 = 12005, \sum x_1^2 = 19461, \sum x_1 x_2 = 6485$$

i.e.

$$12005 = 441 b_0 + 19461 b_1 + 6485 b_2 - (i)$$

and,

$$\sum y x_2 = b_0 \sum x_2 + b_1 \sum x_1 x_2 + b_2 \sum x_2^2$$

where,

$$\sum y x_2 = 4013, \sum x_2 = 147, \sum x_1 x_2 = 6485, \sum x_2^2 = 2173$$

i.e.

$$4013 = 147 b_0 + 6485 b_1 + 2173 b_2 - (ii)$$

Now,

Using Cramers rule to solve (i), (ii) and (iii)

coeff. of  $b_0$

coeff. of  $b_1$

coeff. of  $b_2$

constant

10

441

147

272

441

19461

6485

12005

147

6485

2173

4013

So,

$$D = \begin{vmatrix} 10 & 441 & 147 \\ 441 & 19461 & 6485 \\ 147 & 6485 & 2173 \end{vmatrix}$$

$$D_D = \begin{vmatrix} 272 & 441 & 147 \\ 12005 & 19461 & 6485 \\ 4013 & 6485 & 2173 \end{vmatrix}$$

$$D_{b_1} = \begin{vmatrix} 10 & 272 & 147 \\ 441 & 12005 & 6485 \\ 147 & 4013 & 2173 \end{vmatrix}$$

$$D_2 = \begin{vmatrix} 10 & 441 & 272 \\ 441 & 19461 & 12005 \\ 147 & 6485 & 4013 \end{vmatrix}$$

Then,

$$b_0 = \frac{D_0}{D} = -13.8196$$

$$b_1 = \frac{D_1}{D} = 0.5637$$

$$b_2 = \frac{D_2}{D} = 1.0995$$

Thus,

(\*) can be written as

$$\textcircled{i} \quad y = -13.8196 + 0.5637x_1 + 1.0995x_2$$

\textcircled{ii}

$b_0 = -13.8196$  means annual unpaid taxes discovered is ~13.8196 million dollars when field audit labour hour is 0 and computer hours is 0.

$b_1 = 0.5637$  means annual unpaid taxes discovered increase by 0.5637 million dollars when field audit labour hour increases by 1 and computer hour is constant.

$b_2 = 1.0995$  means annual unpaid taxes discovered increase by 1.0995 million dollars when computer hour increases by 1 and field audit labour hour is kept constant.

\textcircled{iii}

$y = ?$  when  $x_1 = 42$  (in 100 hrs) and  $x_2 = 16$  (in 100 hrs)  
i.e.

$$\begin{aligned} y &= 42 - 13.8196 + 0.5637 * 42 + 1.0995 * 16 \\ &= 27.4478 \end{aligned}$$

$\Rightarrow$  Actual unpaid tax for field labour hour is 4200 and computer hour 1600 is 27.4478 million dollars

9.17 A computer manager needs to know how efficiency of her new computer program depends on the size of incoming data and how many tables are used to arrange each data set. Efficiency will be measured on basis of number of processed requests per hour. Applying program to data set of different sizes and number of tables used, she gets the following results:

Processed requests (y)	16	26	17	41	50	55	40
Data size (GB) ( $x_1$ )	15	10	10	8	7	7	6
No. of tables ( $x_2$ )	1	2	10	10	20	20	4

The regression equation obtained is  $y = 52.7 - 2.87x_1 + 0.85x_2$

Total sum of square = 1452

Sum of square due to regression = 1143.3

- Interpret the values of regression coefficients  $b_1$  and  $b_2$ .
- Test significance of regression model at 0.05 level of significance.
- Is there significant relationship between processed requests and no. of tables at 0.05 level of significance? Given standard error of  $b_2 = 0.55$ .
- What percentage of variation of processed requests is explained by data size and number of tables?
- Compute standard error of estimate.
- Estimate the number of processed requests if data size is 9 GB and number of tables used are 8.

A computer manager needs to know how efficiency of her new computer program depends on the size of incoming data and how many tables are used to arrange each data set. Efficiency will be measured on basis of number of processed requests per hour. Applying program to data set of different sizes and number of tables used, she gets the following results:

Processed requests (y)	16	26	17	41	50	55	40
Data size (GB) ( $x_1$ )	15	10	10	8	7	7	6
No. of tables ( $x_2$ )	1	2	10	10	20	20	4

The regression equation obtained is  $y = 52.7 - 2.87 x_1 + 0.85 x_2$

Total sum of square = 1452

Sum of square due to regression = 1143.3

- Interpret the values of regression coefficients  $b_1$  and  $b_2$ .
- Test significance of regression model at 0.05 level of significance.
- Is there significant relationship between processed requests and no. of tables at 0.05 level of significance? Given standard error of  $b_2 = 0.55$ .
- What percentage of variation of processed requests is explained by data size and number of tables?
- Compute standard error of estimate.
- Estimate the number of processed requests if data size is 9 GB and number of tables used are 8.

Soln:

The given regression equation is

$$y = 52.7 - 2.87 x_1 + 0.85 x_2$$

where,

$$b_0 = 52.7, \quad b_1 = -2.87, \quad b_2 = 0.85$$

for ①

$b_1 = -2.87$  means that the processed request decreases by  $\approx 2.87$  GB when data size increases by 1 GB, while number of tables are kept constant.

$b_2 = 0.85$  means that the processed request increases by 0.85 when no. of table increases by 1 while data size is kept constant.

For ④

Here,

$$\text{y} \quad x_1 \quad x_2 \quad \cdot \quad \text{Total sum of square (TSS)} = 1452$$

$$16 \quad 15 \quad 1 \quad \text{Sum of square due to regression (SSR)} = 1143.3$$

$$26 \quad 10 \quad 2 \quad \text{Total no. of observations (n)} = 7$$

$$17 \quad 10 \quad 10 \quad \text{No. of independent variables (k)} = 2$$

41 8 10 We know,

$$50 7 20 \quad \text{SSR} = TSS - SSE$$

$$55 7 20 \quad \text{or} \quad SSE = TSS - SSR$$

$$40 6 4 \quad \Rightarrow SSE = 308.7 = \text{sum of square due to error}$$

For ⑥

\* Problem to test

$H_0: \beta_j = 0$  (There is no linear relationship between dependent variable y and independent variables)

$H_1: \text{At least one } \beta_j \text{ is different from zero } [j = 1, 2]$

\* Test statistic

$$F = MSR$$

$$MSE$$

Where,

$$MSR = SSR/k = 1143.3/2 = 571.65$$

$$MSE = SSE/(n-k-1) = 308.7/(7-2-1) = 77.175$$

So,

$$F = \frac{571.65}{77.175} = 7.4072 = F_{\text{calculated}}$$

\* Critical value

At  $\alpha = 0.05$ ,  $k=2$ ,  $n-k-1 = 4$ , critical value is  $F_{\alpha}(k, n-k-1)$   
 $= F_{0.05}(2, 4) = 6.944$

\* Decision?

$$F = 7.4072 > F_{0.05}(2, 4) = 6.944$$

Reject  $H_0$  at 0.05 level of significance.

\* Conclusion?

Regression model is significant at 0.05 level of significance.

For ②

Check for significant relationship between  $y$  and  $x_2$  [ $\alpha = 0.05$  and  $s_{b_2} = 0.55$  = standard error of  $b_2$ ].

Here,

\* Problem to test

$$H_0: B_2 = 0$$

$$H_1: B_2 \neq 0$$

\* Test statistic

$$t = \frac{b_2}{s_{b_2}} = \frac{0.85}{0.55} = 1.5455$$

\* Critical value

At  $\alpha = 0.05$ ,  $(n-k-1) = 4$ , critical value (2-tail) is  $t_{\alpha}(n-k-1)$

$$= t_{0.05}(4) = 2.776$$

\* Decision?

$$t = 1.5455 < t_{0.05}(4) = 2.776$$

∴ Accept  $H_0$  at 0.05 level of significance

#### \* Conclusion

There is no significant relationship between processed requests and no. of tables at 0.05 level of significance.

For ①

$$R^2 = \frac{SSR}{TSS} = \frac{1143.3}{1452} = 0.7874$$

⇒ 78.74%. Variation in processed request is explained by data size and no. of tables.

For ②

Standard error of estimate is given by

$$S_e = \sqrt{\frac{SSE}{n-k-1}} = \sqrt{MSG} = \sqrt{77.175} = 8.7849$$

For ③

Data size ( $x_1$ ) = 9 GB

No. of tables used ( $x_2$ ) = 8

No. of processed request ( $y$ ) = ?

where,

$$y = 52.7 - 2.87 x_1 + 0.85 x_2$$

$$= 52.7 - 2.87 \times 9 + 0.85 \times 8$$

$$= 33.67$$

i.e.

No. of processed request is 33.67.

- Q18 A computer manager is keenly interested to know how efficiency of her new computer program depends on size of incoming data.

nd data structure. Efficiency will be measured by no. of processed requests per hour. Data structure may be measured on how many tables were used to arrange each data set. All the information was put together as follows:

Data size (GB)	6	7	7	8	10	10	15
No. of tables	4	20	20	10	10	2	1
Processed requests	40	55	50	41	17	26	16

Identify which one is dependent variable? Fit appropriate multiple regression model and provide problem specific interpretations of the fitted regression coefficients.

Soln:

→ Processed request is dependent variable.

Here,

Required linear regression equation is

$$y = b_0 + b_1 x_1 + b_2 x_2 \quad \textcircled{4}$$

where,

$y$  = processed requests

$x_1$  = Data size (GB)

$x_2$  = No. of tables processed.

Now,

To fit  $\textcircled{4}$  we require certain values for certain equations

which are calculated as:

$y$	$x_1$	$x_2$	$x_1 y$	$x_2 y$	$x_1^2$	$x_2^2$	$x_1 x_2$
40	6	4	240	160	36	16	24
55	7	20	385	1100	49	400	140
50	7	20	250	1000	49	400	140
41	8	10	328	410	64	100	80
17	10	10	170	170	100	100	100
26	10	2	260	52	100	4	20
16	15	1	240	16	225	1	15

Then,

$$\sum y = 245$$

$$\sum x_2 y = 2908$$

$$\sum x_1 = 63$$

$$\sum x_1^2 = 623$$

$$n = 7$$

$$\sum x_2 = 67$$

$$\sum x_2^2 = 1021$$

$$\sum x_1 y = 1973$$

$$\sum x_1 x_2 = 519$$

Then,

$$\sum y = n b_0 + b_1 \sum x_1 + b_2 \sum x_2$$

$$\text{or, } 245 = 7 b_0 + 63 b_1 + 67 b_2 \quad \text{--- (1)}$$

And,

$$\sum y x_1 = b_0 \sum x_1 + b_1 \sum x_1^2 + b_2 \sum x_1 x_2$$

$$\text{or, } 1973 = 63 b_0 + 623 b_1 + 519 b_2 \quad \text{--- (11)}$$

And,

$$\sum y x_2 = b_0 \sum x_2 + b_1 \sum x_1 x_2 + b_2 \sum x_2^2$$

$$\text{or, } 2908 = 67 b_0 + 519 b_1 + 1021 b_2 \quad \text{--- (111)}$$

Now,

Solving (1) (11) and (111) using Cramer's rule

Coeff. of $b_0$	Coeff. of $b_1$	Coeff. of $b_2$	Constant
-----------------	-----------------	-----------------	----------

7	63	67	245
---	----	----	-----

63	623	519	1973
----	-----	-----	------

67	519	1021	2908
----	-----	------	------

$D =$	7	63	67	
	63	623	519	
	67	519	1021	
	245	63	67	
$D_1 =$	1973	623	519	
	2908	519	1021	
$D_1 =$	7	245	67	
	63	1973	519	
	67	2908	1021	
$D_2 =$	7	63	245	
	63	623	1973	
	67	519	2908	

Then,

$$b_0 = D_0 = 51.1685$$

D

$$b_1 = D_1 = -2.7309$$

D

$$b_2 = D_2 = 0.8786$$

D

i.e.,

(\*) can be written as

$$y = 51.1685 - 2.7309x_1 + 0.8786x_2$$

Here,

$b_0 = 51.1685$  means processed request is 51.1685 when data size is 0 GB and no. of tables processed is 0.

$b_1 = -2.7309$  means no. of processed requests decreases by 2.7309 as data size increases by 1 GB keeping no. of tables processed constant.

$b_2 = 0.8786$  means no. of processed request increases by 0.8786 as no. of tables increase by 1 keeping data size constant.

Q19 What is Multiple Linear Regression (MLR)? From following information of variables  $x_1$ ,  $x_2$  and  $y$ , fit a regression equation of  $y$  on  $x_1$  and  $x_2$ .

$$\sum x_1 = 272, \sum x_2 = 441, \sum y = 147, \sum x_1^2 = 7428$$

$$\sum x_2^2 = 19461, \sum y^2 = 2137, \sum x_1 x_2 = 4013, \sum x_1 x_2 = 12005$$

$$\sum x_1 y = 6485, n=10.$$

Soln:

Required linear regression equation is

$$y = b_0 + b_1 x_1 + b_2 x_2 \quad \textcircled{1}$$

Here,

To fit  $\textcircled{1}$ , we use the following equations

$$\sum y = n b_0 + b_1 \sum x_1 + b_2 \sum x_2$$

$$147 = 10 b_0 + 272 b_1 + 441 b_2 \quad \textcircled{1}$$

And,

$$\sum y x_1 = b_0 \sum x_1 + b_1 \sum x_1^2 + b_2 \sum x_1 x_2$$

$$\text{or, } 4013 = 272 b_0 + 19461 b_1 + 12005 b_2 \quad \textcircled{11}$$

And,

$$\sum y x_2 = b_0 \sum x_2 + b_1 \sum x_1 x_2 + b_2 \sum x_2^2$$

$$\text{or, } 6485 = 441 b_0 + 12005 b_1 + 19461 b_2 \quad \textcircled{111}$$

Then,

Solving  $\textcircled{1}$   $\textcircled{11}$  and  $\textcircled{111}$  using Cramer's rule

Coeff. of  $b_0$       Coeff. of  $b_1$       Coeff. of  $b_2$       Constant

10	272	441	147
272	19461	12005	4013
441	12005	19461	6485

$D =$	10	272	441
	272	<del>7428</del> 19461	12005
	441	12005	19461

$D_0 =$	147	272	441
	4013	<del>7428</del> 19461	12005
	6485	12005	19461

$D_1 =$	10	147	441
	272	<del>7428</del> 19461	12005
	441	6485	19461

$D_2 =$	10	272	147
	272	<del>7428</del> 19461	4013
	441	12005	6485

Now,

$$b_0 = \frac{D_0}{D} = 10.4934$$

$$b_1 = \frac{D_1}{D} = 0.5801$$

$$b_2 = \frac{D_2}{D} = -0.2624$$

Thus,

(\*) can be written as

$$y = 10.4934 + 0.5801x_1 - 0.2624x_2$$

### Multiple linear regression (MLR)

MLR can be defined as linear functional relationship of a variable with 2 or more variables.

Q20

Suppose we are given the following information with  $n=7$ . Multiple regression model is

$$\hat{y} = 8.15 + 0.56x_1 + 0.54x_2$$

Here, Total sum of square = 1493

Sum of square due to error = 91

Find

(i)  $R^2$  and interpret it.

(ii) Test overall significance of model.

SOLN:

We are given

Total sum of square (TSS) = 1493

Sum of square due to error (SSE) = 91

$n = 7$ ,

$x_1$  &  $x_2$  are 2 independent variable  $\Rightarrow k = 2$

Then,

Sum of square due to regression is given by

$$SSR = TSS - SSE \\ = 1493 - 91$$

$$\Rightarrow SSR = 1402$$

Now,

$$R^2 = \frac{SSR}{TSS} = \frac{1402}{1493} = 0.9390 = 93.9\%$$

It means 93.9% variation in  $y$  is explained by  $x_1$  &  $x_2$

Here,

$$\text{Mean Square Regression (MSR)} = \frac{SSR}{k}$$

$$\Rightarrow MSR = 1402 / 2 = 701$$

$$\text{Mean Square Error (MSE)} = \frac{SSE}{n-k-1}$$

$$\Rightarrow MSE = 91 = 22.75$$

For 11)

\* Problem to test

$H_0: \beta_j = 0$  [Dependent and independent variables are insignificant]

$H_1: \text{At least one } \beta_j \text{ is different from zero } [j=1, 2]$  i.e.

There is linear relationship between dependent and at least one

independent variable  $\Rightarrow$  Dependent and independent variables are significant

\* Test statistic

$$F = \frac{MSR}{MSE} = \frac{701}{22.75} = 30.8132$$

\* Critical value

Let  $\alpha = 5\% = 0.05$ . At  $\alpha = 0.05$ ,  $k = 2$  and  $(n - k - 1) = 4$ , critical value is  $F_{\alpha}(k, n - k - 1) = F_{0.05}(2, 4) = 6.944$

\* Decision

$$F = 30.8132 > F_{0.05}(2, 4) = 6.944$$

Reject  $H_0$  at 0.05 level of significance

\* Conclusion

Regression model is significant.

Q22 The following ANOVA summary table was obtained from a multiple regression model with 2 independent variable

SV	SS	df	MS	F-ratio
Regression	12.62	2	?	?
Error	0.78	12	?	
Total	13.4	14		

- Determine Mean sum of square due to regression, the mean sum of square due to error and F-value.
- Test significance of overall regression model at 5% level of significance

- iii. Compute coefficient of determination and interpret its value.  
 iv. Find standard error of estimate.

SOLN:

From the given ANOVA table

$$\text{Sum of square due to regression (SSR)} = 12.62$$

$$\text{Sum of square due to error (SSE)} = 0.78$$

$$df(\text{regression}) = 2 = k$$

$$df(\text{error}) = 12 = n - k - 1$$

$$\therefore n = 15$$

①

Then,

Mean sum of square due to Regression is given by

$$MSR = \frac{SSR}{k} = \frac{12.62}{2} = 6.31$$

And,

Mean sum of square due to error is given by

$$MSE = \frac{SSE}{n-k-1} = \frac{0.78}{12} = 0.065$$

Now,

F-value is given by

$$F = \frac{MSR}{MSE} = \frac{6.31}{0.065} = 97.0769$$

②

\* Problem to test

$$H_0: \beta_i = 0$$

$H_1: \text{At least one } \beta_i \text{ is different from zero } \quad \{i = 1, 2\}$

\* Test Statistic

$$F = 97.0769$$

\* Critical value

$H_0: \alpha = 5\% = 0.05, k = 2$  and  $n - k - 1 = 12$ , critical value is  
 $F_{(k, n-k-1)} = F_{0.05}(2, 12) = 3.885$

\*  $F = 97.0769 > F_{0.05}(2, 12) = 3.885$

Reject  $H_0$  at  $5\%$  level of significance.

\* Conclusion  
The overall regression model is significant at  $5\%$  level of significance.

(ii) We know,

$$TSS = SSR + SSE = 13.4$$

Then,

$$\text{Coeff. of determination} = R^2 = \frac{SSR}{TSS} = \frac{12.62}{13.4} = 0.9418 = 94.18\%$$

It means  $94.18\%$  variation in  $y$  is explained by  $x_1$  and  $x_2$ .

(iii) Standard error of estimate ( $s_e$ ) =  $\sqrt{MSE} = \sqrt{0.065} = 0.2550$

Hence,

①  $MSR = 6.31$

$MSE = 0.065$

$F = 97.0769$

(iv) The overall regression model is significant at  $5\%$  level of significance.

(v)  $R^2 = 0.9418 = 94.18\%$ .

(vi) Standard error of estimate =  $0.2550$ .