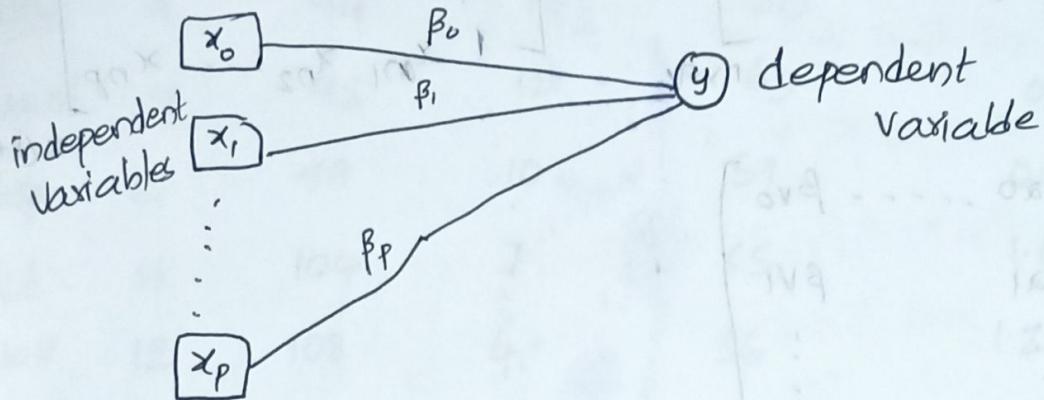


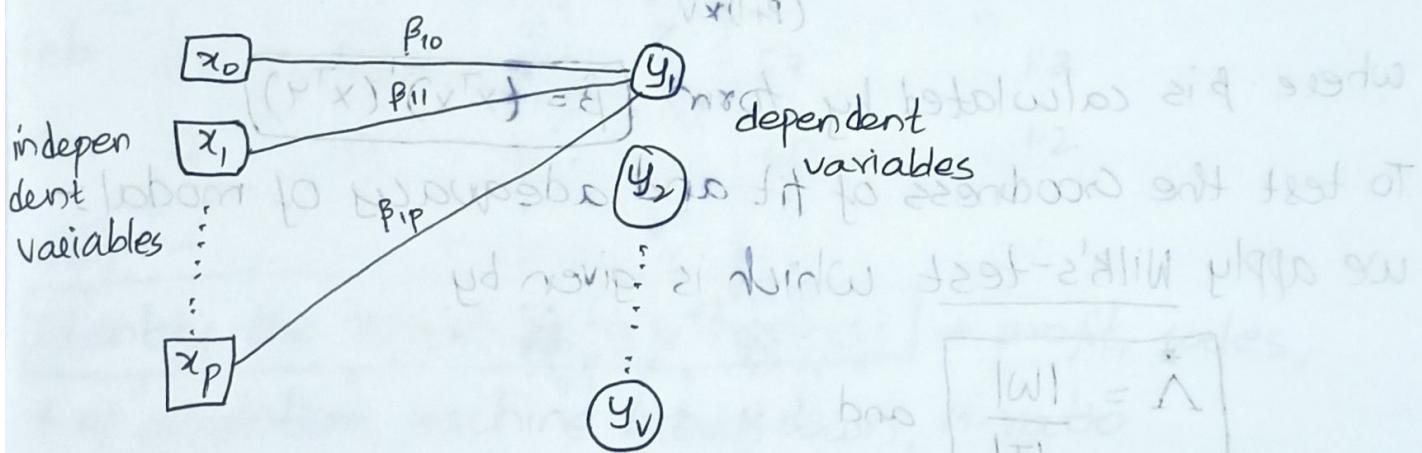
UNIT - 3

Multi-variate Regression:-

Multiple linear regression:- (MLR)



Multi-variate linear regression:- (MVLR)



- Multivariate linear regression model is given by

$$y_1 = \beta_{10} + \beta_{11}x_1 + \beta_{12}x_2 + \dots + \beta_{1p}x_p + \epsilon_1$$

$$y_2 = \beta_{20} + \beta_{21}x_1 + \beta_{22}x_2 + \dots + \beta_{2p}x_p + \epsilon_2$$

$$y_v = \beta_{v0} + \beta_{v1}x_1 + \beta_{v2}x_2 + \dots + \beta_{vp}x_p + \epsilon_v$$

Let the multivariate linear regression is of the form

$$Y = X\beta + \epsilon$$

$$\Rightarrow Y_{n \times v} = X_{n \times (p+1)} \cdot \beta_{(p+1) \times v} + \epsilon_{n \times v}$$

where

$$y = \begin{bmatrix} y_{11} & y_{12} & \dots & y_{1q} \\ y_{21} & y_{22} & \dots & y_{2q} \\ \vdots & \vdots & & \vdots \\ y_{n1} & y_{n2} & \dots & y_{nq} \end{bmatrix}_{n \times q}$$

$$x = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}_{n \times (p+1)}$$

$$\beta = \begin{bmatrix} \beta_{10} & \beta_{11} & \dots & \beta_{1q} \\ \beta_{20} & \beta_{21} & \dots & \beta_{2q} \\ \vdots & \vdots & & \vdots \\ \beta_{p0} & \beta_{p1} & \dots & \beta_{pq} \end{bmatrix}_{(p+1) \times q}$$

where β is calculated by form $\hat{\beta} = (x^T x)^{-1} (x^T y)$

To test the Goodness of fit and adequacy of model :-
we apply Wilks's-test which is given by

$$\lambda^* = \frac{|W|}{|I|} \quad \text{and}$$

$$F = \frac{n-k-1}{k-1} * \frac{1-\sqrt{\lambda^*}}{\sqrt{\lambda^*}} \sim F_{P(k-1), P(n-k-1)}$$

which is obtained by the multivariate analysis of variance (MANOVA)

Eg:- Fit a multivariate linear regression model with profit & sales depending on percentage of absenteeism, machine breakdown and M ratio.

<u>Month</u>	<u>Profit</u>	<u>Sales</u>	<u>% of Absentis</u>	<u>Machine Break</u>	<u>M-Ratio</u>
April	10	100	9	62	1.6
May	12	110	8	58	1.3
June	11	105	7	64	1.2
July	9	94	14	60	0.8
August	9	95	12	63	0.8
Sep	10	99	10	57	0.9
Oct	11	104	7	55	1.0
Nov	12	108	4	56	1.2
Dec	11	105	8	59	1.1
Jan	10	98	5	61	1.0
Feb	11	103	7	57	1.2
Mar	12	110	6	60	1.2

Step 1:-

Identify the variables of interest → profit, sales, % of absenteeism, machine breakdown, M-ratio

Step 2:- Identify the response variables →

$y_1 \rightarrow$ profit
 $y_2 \rightarrow$ sales } dependent variables

Step 3:- Identify the explanatory variables →

$x_1 \rightarrow$ % of absenteeism

$x_2 \rightarrow$ machine breakdowns } independent variables

$x_3 \rightarrow$ M-ratio } variables

Step 4:- Find the dependent relationships →

$$y = f(x)$$

Step 5:- Solve the above system of relations by applying the formula.

$$\hat{\beta} = (X^T X)^{-1} (X^T Y)$$

Let the multivariate linear regression modal in this problem the of the form

$$Y_1 = \beta_{10} + \beta_{11} X_1 + \beta_{12} X_2 + \beta_{13} X_3$$

$$Y_2 = \beta_{20} + \beta_{21} X_1 + \beta_{22} X_2 + \beta_{23} X_3$$

where

$$Y = \begin{bmatrix} \text{Profit} & \text{Sales} \\ 10 & 100 \\ 12 & 110 \\ 11 & 105 \\ 9 & 94 \\ 9 & 95 \\ 10 & 99 \\ 11 & 104 \\ 12 & 108 \\ 11 & 105 \\ 10 & 98 \\ 11 & 103 \\ 12 & 110 \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & 9 & 62 & 1.0 \\ 1 & 8 & 58 & 1.3 \\ 1 & 7 & 64 & 1.2 \\ 1 & 14 & 60 & 0.8 \\ 1 & 12 & 63 & 0.8 \\ 1 & 10 & 57 & 0.9 \\ 1 & 7 & 55 & 1.0 \\ 1 & 6 & 56 & 1.2 \\ 1 & 5 & 61 & 1.0 \\ 1 & 7 & 57 & 1.2 \\ 1 & 6 & 60 & 1.2 \end{bmatrix}$$

where $\hat{\beta} = (X^T X)^{-1} (X^T Y)$

$$\Rightarrow \hat{\beta} = \begin{bmatrix} 10.897 & 91.097 \\ -0.045 & -0.064 \\ -0.087 & -0.294 \\ 5.035 & 27.835 \end{bmatrix}$$

$$(X)I = E$$

$$Y_1(\text{Profit}) = 10.2917 - 0.045x_1 - 0.087x_2 + 5.035x_3$$

$$Y_2(\text{Sales}) = 91.097 - 0.064x_1 - 0.294x_2 + 27.835x_3$$

Eg: Construct one way MANOVA table to the following data.

Treatment 1: $\begin{bmatrix} 9 \\ 3 \end{bmatrix} \begin{bmatrix} 6 \\ 2 \end{bmatrix} \begin{bmatrix} 9 \\ 7 \end{bmatrix}$ $\bar{y}_1 = \frac{8}{3}$ $\bar{y} = \frac{10}{3}$

Treatment 2: $\begin{bmatrix} 0 \\ 4 \end{bmatrix} \begin{bmatrix} 2 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 1 \\ 2 \end{bmatrix}$ $\begin{bmatrix} 4 \\ 5 \end{bmatrix}$

Treatment 3: $\begin{bmatrix} 3 \\ 8 \end{bmatrix} \begin{bmatrix} 1 \\ 9 \end{bmatrix} \begin{bmatrix} 2 \\ 7 \end{bmatrix}$ $\begin{bmatrix} 2 \\ 8 \end{bmatrix}$

H₀: $(\mu_{11} = \mu_{12} = \mu_{13}) = (\mu_{21} = \mu_{22}) = (\mu_{31} = \mu_{32} = \mu_{33})$

H₁: $(\mu_{11} \neq \mu_{12} \neq \mu_{13}) \neq (\mu_{21} \neq \mu_{22}) \neq (\mu_{31} = \mu_{32} = \mu_{33})$

Level of Significance = 5% (chosen)

To Test the above hypothesis, the procedure is as follows:-

For y_{ij} :

① Sum of squares due to error

$$\begin{aligned} SSE &= \sum (y_{ij} - \bar{y}_i)^T (y_{ij} - \bar{y}_i) = \sum (y_{ij} - \bar{y}_i)^2 \\ &= (9-8)^2 + (6-8)^2 + (9-8)^2 + (0-1)^2 + (2-1)^2 + (3-2)^2 + \\ &\quad (1-2)^2 + (2-2)^2 = 10 \end{aligned}$$

② Sum of squares due to total

$$\begin{aligned} SST &= \sum (y_{ij} - \bar{y})^T (y_{ij} - \bar{y}) = \sum (y_{ij} - \bar{y})^2 \\ &= (9-4)^2 + (6-4)^2 + (9-4)^2 + (0-4)^2 + (2-4)^2 + (3-4)^2 \\ &\quad + (1-4)^2 + (2-4)^2 = 88 \end{aligned}$$

③ Sum of squares due to Regression

$$SSR = \sum n_i (\bar{y}_i - \bar{y})^T (\bar{y}_i - \bar{y}) = SST - SSE = 88 - 10 = 78$$

FOR y_2 :

$$\textcircled{1} \quad SSe = \sum (y_{ij} - \bar{y}_i)^2 = \sum (y_{ij} - \bar{y}_i)^2 \\ = (3-4)^2 + (2-4)^2 + (7-4)^2 + (4-2)^2 + (0-2)^2 + (8-8)^2 + \\ (9-8)^2 + (7-8)^2 = 24$$

$$\textcircled{2} \quad SST = \sum (y_{ij} - \bar{y})^2 = \sum (y_{ij} - \bar{y})^2 \\ = (3-5)^2 + (2-5)^2 + (7-5)^2 + (4-5)^2 + (0-5)^2 + (8-5)^2 + \\ (9-5)^2 = 72$$

$$\textcircled{3} \quad SS_B = \sum n_i (\bar{y}_i - \bar{y}) (\bar{y}_i - \bar{y}) = SST - SSe = 72 - 24 = 48$$

Cross product values of y_1 & y_2 :

$$\textcircled{1} \quad SSe = (9 \times 3 - 8 \times 4) + (6 \times 2 - 8 \times 4) + (9 \times 7 - 8 \times 4) \\ + (0 \times 4 - 1 \times 2) + (2 \times 0 - 1 \times 2) + (3 \times 8 - 2 \times 8) + (1 \times 9 - 2 \times 8) \\ + (2 \times 7 - 2 \times 8) = 1$$

$$\textcircled{2} \quad SST = (9 \times 3 - 4 \times 5) + (6 \times 2 - 4 \times 5) + (9 \times 7 - 4 \times 5) + (0 \times 4 - 4 \times 5) \\ + (2 \times 0 - 4 \times 5) + (3 \times 8 - 4 \times 5) + (1 \times 9 - 4 \times 5) + (2 \times 7 - 4 \times 5) = 7 + (-8) + 43 - 20 - 20 + 4 - 11 - 6 = -11$$

$$\textcircled{3} \quad SS_B = SST - SSe = -11 - 1 = -12$$

MANOVA One Way Classification Table:

Source of Variation	sum of squares	degrees of freedom	Wilks' Value	F-Test
Regression	$\sum n_i (\bar{y}_i - \bar{y})^2$ $(\bar{y}_i - \bar{y}) = B$	$K-1$	$\lambda^* = \frac{ W }{ T }$	$F = \frac{D-K-1}{K-1} *$
error	$\sum (y_{ij} - \bar{y}_i)^2$	$D-K$		$\frac{1-\sqrt{\lambda^*}}{\sqrt{\lambda^*}} \sim F_{P(K-1), P(n-K-1)}$
Total	$\sum (y_{ij} - \bar{y})^2$	$D-1$		

Inference:-

If $F_{cal} > F_{table}$ value, we reject H_0

Hence we conclude that otherwise we accept H_0

MANOVA One way classification Table :-

source of variation	sum of squares	degrees of freedom	wil's value	F-test
Regression	$\begin{pmatrix} 78 & -12 \\ -12 & 48 \end{pmatrix} = B$	$3-1=2$	$\lambda = \frac{B^T B}{ B } = \frac{239}{6215} = 0.0384$	$F = \frac{(n-k-1) * 1 - \sqrt{\lambda}}{k-1} * \frac{1}{\sqrt{\lambda}}$ $= \frac{8-3-1}{3-1} * \frac{1 - \sqrt{0.0384}}{\sqrt{0.0384}}$
Error	$\begin{pmatrix} 10 & 1 \\ 1 & 24 \end{pmatrix} = W$	$8-3=5$		
Total	$\begin{pmatrix} 88 & -11 \\ -11 & 72 \end{pmatrix} = T$	$8-1=7$		$= 8.2$ $\sim F_{2(3-1), 7(8-3-1)}$ $\sim F_{4, 8}$

Inference:-

The table value at 5% L.O.S at 4, 8 degrees of freedom = 3.84

$\therefore F_{cal} > F_{table}$, we reject H_0

Hence we conclude that there is no homogeneity among the treatments.

② Construct one-way MANOVA table to the following data.

$$\text{Treatment A: } \begin{bmatrix} 2 \\ 3 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} 5 \\ 4 \end{bmatrix} \begin{bmatrix} 2 \\ 5 \end{bmatrix}$$

$$\text{Treatment B: } \begin{bmatrix} 4 \\ 8 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix} \begin{bmatrix} 6 \\ 7 \end{bmatrix}$$

$$\text{Treatment} \leftarrow [7][8][10][9][7] \\ [6][7][8][5][6]$$

$$\text{Sd: } H_0: (U_{11}=U_{12}=U_{13}=U_{14}) = (U_{21}=U_{22}=U_{23}) = \\ (U_{31}=U_{32}=U_{33}=U_{34}=U_{35})$$

$$H_1: (U_{11}=U_{12}=U_{13}=U_{14}) \neq (U_{21}=U_{22}=U_{23}) \neq \\ (U_{31}=U_{32}=U_{33}=U_{34}=U_{35})$$

Level of significance = s.l. (chosen)

To test the above hypothesis, the procedure is as follows:-

	\bar{y}_i	$\bar{\bar{y}}$
A	$[2][3][5][2]$	$[3]$
B	$[4][5][6]$	$[5]$
C	$[7][8][10][9][7]$	$[8.8]$

For \bar{y}_i :

$$\begin{aligned} \textcircled{1} \text{ SSE} &= \sum (y_{ij} - \bar{y}_i)^T (y_{ij} - \bar{y}_i) = \sum (y_{ij} - \bar{y}_i)^2 \\ &= (2-3)^2 + (3-3)^2 + (5-3)^2 + (2-3)^2 + (4-5)^2 + (5-5)^2 + \\ &\quad (6-5)^2 + (7-8)^2 + (8-8)^2 + (10-8)^2 + (9-8)^2 + \\ &\quad (7-8)^2 = 14.8 \end{aligned}$$

$$\begin{aligned} \textcircled{2} \text{ SST} &= \sum (y_{ij} - \bar{\bar{y}})^T (y_{ij} - \bar{\bar{y}}) = \sum (y_{ij} - \bar{\bar{y}})^2 \\ &= (2-5.67)^2 + (3-5.67)^2 + (5-5.67)^2 + (2-5.67)^2 + (4-5.67)^2 \\ &\quad + (5-5.67)^2 + (6-5.67)^2 + \dots + (7-5.67)^2 \\ &= 76.67 \end{aligned}$$

$$\textcircled{3} \cdot SSe = \sum n_i (\bar{y}_i - \bar{\bar{y}})^T (\bar{y}_i - \bar{\bar{y}}) = SST - SSE = 61.87$$

For y_2 :

$$\begin{aligned} \textcircled{1} SSE &= \sum (y_{ij} - \bar{y}_i)^T (y_{ij} - \bar{y}_i) = \sum (y_{ij} - \bar{y}_i)^2 \\ &= (3-4)^2 + (4-4)^2 + (4-4)^2 + (5-4)^2 + (8-7)^2 + (6-7)^2 + \\ &\quad (7-7)^2 + (6-6.4)^2 + (7-6.4)^2 + (8-6.4)^2 + (5-6.4)^2 + (6-6.4)^2 \\ &= 9.2 \end{aligned}$$

$$\begin{aligned} \textcircled{2} SST &= \sum (y_{ij} - \bar{\bar{y}})^T (y_{ij} - \bar{\bar{y}}) = \sum (y_{ij} - \bar{\bar{y}})^2 \\ &= (3-5.75)^2 + (4-5.75)^2 + (5-5.75)^2 + \dots + (6-5.75)^2 \\ &= 28.25 \end{aligned}$$

$$\textcircled{3} SSt = \sum n_i (\bar{y}_i - \bar{\bar{y}}) (\bar{y}_i - \bar{\bar{y}}) = SST - SSE = 19.05$$

For both y_1 & y_2 :

$$\begin{aligned} \textcircled{1} SSE &= (2 \times 3 - 3 \times 4)^2 + (3 \times 4 - 3 \times 4)^2 + (5 \times 4 - 3 \times 4)^2 + (2 \times 5 - 3 \times 4)^2 \\ &\quad + (4 \times 8 - 5 \times 7)^2 + (5 \times 6 - 5 \times 7)^2 + (6 \times 7 - 5 \times 7)^2 + \\ &\quad (7 \times 6 - 8 \times 6.4)^2 + (8 \times 7 - 8 \times 6.4)^2 + \dots = 1.6 \\ \textcircled{2} SST &= (2 \times 3 - 5.67 \times 5.75)^2 + (3 \times 4 - 5.67 \times 5.75)^2 + \dots \\ &\quad + (7 \times 6 - 5.67 \times 5.75)^2 = 26 \end{aligned}$$

$$\textcircled{3} SSt = SST - SSE = 26 - 1.6 = 24.4$$

MANOVA ONE WAY CLASSIFICATION TABLE

Source of Variation	Sum of Squares	Degrees of freedom	Wtll's value	F-test
Regression	$(61.87 \ 24.4)$ $(24.4 \ 19.05)$ = B	$3-1=2$	$* = \frac{ W }{ T }$	$F = \frac{12-3-1}{3-1} \times$
Error	$(14.8 \ 1.6)$ $(1.6 \ 9.2)$ = W	$12-3=9$	$= \frac{133.6}{148.9} \cdot 9$	$1 - \frac{\sqrt{0.0896}}{\sqrt{0.0896}}$
Total	$(76.67 \ 26)$ $(26 \ 28.25)$ = T	$12-1=11$	$= 0.0896$	$F = 9.36$ $\sim F_{4,16}$

Inference:-

The table value of F at 5% L.O.S for 4,16 degrees of freedom is $F_{tab} = 3.01$

$F_{cal} > F_{tab}$, we reject H_0

hence we conclude that there is no homogeneity between the constraints.

- ③ Carry out MANOVA one way classification at 1% level of significance and comment.

$$\bar{y}_i \quad \bar{y}_{\bar{i}}$$

Treatment 1: $\begin{bmatrix} 4 \\ 6 \end{bmatrix} \begin{bmatrix} 3 \\ 7 \end{bmatrix} \begin{bmatrix} 7 \\ 5 \end{bmatrix} \begin{bmatrix} 4 \\ 8 \end{bmatrix} \begin{bmatrix} 5 \\ 8 \end{bmatrix} \begin{bmatrix} 4.6 \\ 6.8 \end{bmatrix}$

Treatment 2: $\begin{bmatrix} 4 \\ 5 \end{bmatrix} \begin{bmatrix} 2 \\ 7 \end{bmatrix} \begin{bmatrix} 2 \\ 5 \end{bmatrix} \begin{bmatrix} 2.667 \\ 5.667 \end{bmatrix} \begin{bmatrix} 3.75 \\ 5.33 \end{bmatrix}$

Treatment 3: $\begin{bmatrix} 3 \\ 6 \end{bmatrix} \begin{bmatrix} 6 \\ 2 \end{bmatrix} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \begin{bmatrix} 3.5 \\ 3.85 \end{bmatrix}$

Sol:- $H_0: (\mu_{11} = \mu_{12} = \mu_{13} = \mu_{14} = \mu_{15}) = (\mu_{21} = \mu_{22} = \mu_{23}) = (\mu_{31} = \mu_{32} = \mu_{33} = \mu_{34})$

$H_1: (\mu_{11} = \mu_{12} = \mu_{13} = \mu_{14} = \mu_{15}) \neq (\mu_{21} = \mu_{22} = \mu_{23}) \neq (\mu_{31} = \dots = \mu_{34})$

Level of significance = 1% (Given)

To test the above hypothesis, the procedure is as follows,

table

For y_{ij} :

$$\begin{aligned} ① SSE &= \sum (y_{ij} - \bar{y}_i)^2 = \sum (y_{ij} - \bar{y}_i)^2 \\ &= (4-4.6)^2 + (3-4.6)^2 + (7-4.6)^2 + (4-4.6)^2 + (5-4.6)^2 + \\ &\quad (4-2.667)^2 + (2-2.667)^2 + (2-2.667)^2 + (3-3.5)^2 + \\ &\quad (6-3.5)^2 + (2-3.5)^2 + (3-3.5)^2 = 0.36 + 2.56 + 5.76 + 0.36 + \end{aligned}$$

$$0.16 + 1.776 + 0.4356 + 0.4356 + 0.25 + 6.25 + 2.25 + 0.25 = \\ = 20.8472 \quad (20.8472)$$

$$\textcircled{2} \quad SST = \sum (y_{ij} - \bar{y})^T (y_{ij} - \bar{y}) = \sum (y_{ij} - \bar{y})^2 \\ = (4-3.75)^2 + (3-3.75)^2 + (7-3.75)^2 + (4-3.75)^2 + (5-3.75)^2 + \\ (4-3.75)^2 + (2-3.75)^2 + (2-3.75)^2 + (3-3.75)^2 + (6-3.75)^2 + \\ (2-3.75)^2 + (3-3.75)^2 = 0.0625 + 0.5625 + 10.5625 + 0.0625 + \\ 1.5625 + 0.0625 + 3.0625 + 3.0625 + 0.25 + 6.25 + 3.0625 + \\ 5.0625$$

$$0.25 = 28.8125 \quad (28.8125)$$

$$\textcircled{3} \quad SSE = \sum n_i (\bar{y}_i - \bar{y})^T (\bar{y}_i - \bar{y}) = SST - SSE = 28.8125 - 20.8472 \\ = 7.9653 \quad (7.9653)$$

For y_1 :

$$\textcircled{1} \quad SSE = \sum (y_{ij} - \bar{y}_i)^T (y_{ij} - \bar{y}_i) = \sum (y_{ij} - \bar{y}_i)^2 \\ = (6-6.8)^2 + (7-6.8)^2 + (5-6.8)^2 + (8-6.8)^2 + (5-5.667)^2 + \\ (7-5.667)^2 + (5-5.667)^2 + (6-3.25)^2 + (2-3.25)^2 + (1-3.25)^2 + \\ (4-3.25)^2 = 0.64 + 0.04 + 3.24 + 1.44 + 1.44 + 0.444 + 1.776 + \\ 0.444 + 7.5625 + 1.5625 + 5.0625 + 0.5625 = 24.214$$

$$\textcircled{2} \quad SST = 52.66$$

$$\textcircled{3} \quad SSt = 52.66 - 24.214 = 28.446.$$

For both y_1 & y_2 :

$$\textcircled{1} \quad SSE = (4 \times 6 - 4.6 \times 6.8) + (3 \times 7 - 4.6 \times 6.8) + (7 \times 5 - 4.6 \times 6.8) + \\ \dots + (3 \times 4 - 3.5 \times 3.25) = -7.32$$

$$\textcircled{2} \quad SST = (4 \times 6 - 3.75 \times 5.33) + (3 \times 7 - 3.75 \times 5.33) + \dots + \\ \dots + (3 \times 4 - 3.75 \times 5.33) = 0.15$$

$$\textcircled{3} \quad SSt = SST - SSE = 7.47$$

MANOVA ONE WAY CLASSIFICATION TABLE :-

Source of Variation	Sum of squares	Degrees of freedom	Wil's value	F-statistic
Regression	$\begin{bmatrix} 7.38 & 7.47 \\ 7.47 & 28.41 \end{bmatrix} = B$	$3-1=2$	$\lambda = \frac{ W }{ T }$	$F = \frac{12-3-1}{3-1} *$
Error	$\begin{pmatrix} 20.87 & -7.32 \\ -7.32 & 24.22 \end{pmatrix} = W$	$12-3=9$	$= \frac{451.889}{1487.905}$	$\frac{1-\sqrt{0.3037}}{\sqrt{0.3037}}$
Total	$\begin{pmatrix} 28.25 & 0.15 \\ 0.15 & 52.66 \end{pmatrix} = T$	$12-1=11$	≈ 0.3037	$= 3.26$ $\sim F_{2(3-1), 2(12-3)}$ $\sim F_{4,6}$

Inference:-

The table value of F at 1% L.O.S for 4, 16 degrees of freedom is $F_{table} = 4.77$

$\therefore F_{cal} < F_{table}$, we accept H_0

Hence we conclude that there is homogeneity in the groups and also among the groups.

i.e., the multivariate linear regression modal (MVLR) is applicable for future predictions.

Theory Questions:-

- ① Write the procedure to fit the Multivariate linear regression modal.
- ② Discuss the statistical Analysis of multivariate analysis of variance for ONE way classification (OR) What are the statistics commonly used in testing

the hypothesis in MANOVA.

③ what is the difference between ANOVA and MANOVA?

Sol: ANOVA:-

- ANOVA stands for analysis of variance.
- In statistics, when two or more than two means are compared simultaneously, the statistical method used to make comparison is called ANOVA.
- It is a method which gives values and the results which can be tested in order to determine if a relationship of any significance exists between different variables
- The name ANOVA has been given to the comparison of means because in order to determine any relationship between the means, the variances are actually being compared under currently

Procedure for ANOVA:-

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

$$H_1: \mu_1 \neq \mu_2 \neq \dots \neq \mu_k$$

Level of significance is $\alpha\%$ (Given/choosen)

To test the above hypothesis, the ANOVA table is given by

Source of Variation	Sum of squares	degrees of freedom	mean sum of squares	variance ratio
Treatments	S_{tr}^2	$k-1$	$MS_{tr}^2 = \frac{S_{tr}^2}{k-1}$	$F = \frac{MS_{tr}^2}{MS_e^2} \sim F_{k-1, n-k}$
Error	S_e^2	$n-k$	$MS_e^2 = \frac{S_e^2}{n-k}$	
Total	S_T^2	$n-1$		

Inference:-

$F_{cal} \geq F_{\alpha, k-1, n-k}$, we reject H_0 .
otherwise we accept H_0 .

MANOVA :-

- MANOVA stands for multivariate Analysis of Variance.
- In statistics, the MANOVA methods contain multiple dependent variables. They help in determining the differences between either two or more dependent variables.
- MANOVA assists in determining these differences simultaneously.
- In MANOVA method determines the dependent variable gets significant effected by the change in independent variables, it also determines the interactions taking place amongst dependent variables.

Procedure:-

$$H_0: (\mu_{11} = \mu_{12} = \dots) = (\mu_{21} = \mu_{22} = \dots) = \dots = (\mu_{k1} = \mu_{k2} = \dots)$$

$$H_1: (\mu_{11} = \mu_{12} = \dots) \neq (\mu_{21} = \mu_{22} = \dots) \neq \dots \neq (\mu_{k1} = \mu_{k2} = \dots)$$

Level of significance α (Given/choosen)
To test the above hypothesis the procedure is
as follows:-

Source of variation	sum of squares	D.O.F	W.H.'s value	F-statistic
Treatments	$\sum n_i (\bar{y}_i - \bar{\bar{y}})^T (\bar{y}_i - \bar{\bar{y}}) = B$	K-1	$\lambda = \frac{ W }{ T }$	$F = \frac{n-K-1}{D-K} *$
Error	$\sum (y_{ij} - \bar{y}_i)^T (y_{ij} - \bar{y}_i) = W$	N-K		$\frac{1 - \sqrt{\lambda}}{\sqrt{\lambda}} \sim$
Total	$\sum (y_{ij} - \bar{\bar{y}})^T (y_{ij} - \bar{\bar{y}}) = T$	N-1		$F_{P(K-1), P(n-K-1)}$

where

k indicates the number of treatments.

p indicates the number of values in each sub group.

n indicates total number of groups in the experiment.

Inference:-

Summary of MULR:

- ① There will be more than one independent variable
- ② Simultaneous occurrence of dependent variable
- ③ Independent variables are same
- ④ Estimation of the parameters is similar to multiple linear Regression Model (MLR)
- ⑤ Model Adequacy Test is countless.
- ⑥ The Assumptions in MLR are equally important to MULR.
- ⑦ Does not give much advantage over MLR but protects against the event handling.

⑧ The tests applied for goodness of fit are

MLR - ANOVA, R²:

MVLR - MANOVA:

Discriminant Analysis :- (Data classification)

Discriminant Analysis is a statistical technique used to classify observations into non overlapping groups based on one or more quantitative predictor variables.

There are several different ways to conduct the discriminant analysis. One such approach is based on linear regression and the other is linear discriminant analysis.

Discriminant Analysis Based on Linear Regression:-

When there are only 2 classification groups, the linear discriminant analysis based on linear regression is applied when

- The dependent variable is a categorical variable
(A categorical variable takes only 2 values like Yes/No, True/False, 0/1, Male/Female etc)
- The dependent variable is expressed as a dummy variable (values having zeroes or ones)
- Observations are assigned to groups based on whether the predicted score is closer to 0 or 1.
- The regression equation is called the discriminant function.

The biggest difference between discriminant analysis and standard multiple regression analysis is the use of categorical variable as a dependent variable.

e.g. The SAT is an aptitude test taken by high school juniors and seniors. The college administration use the SAT along with high school grade point avg. to predict the academic success in the college.

The following table shows the SAT and GPA scores for 10 students in a college and shows whether each student ultimately graduated from the college.

<u>SAT</u>	<u>GPA</u>	<u>Graduate</u>
1300	2.7	Yes
1260	3.7	Yes
1220	2.9	Yes
1180	2.5	Yes
1060	3.9	Yes
1140	2.1	No
1100	3.5	No
1020	3.3	No
980	2.3	No
940	3.1	No

Define a discriminant function that classifies an incoming student with SAT 1000 and GPA = 2.9 as graduate or non-graduate.

Sol: The dependent variable (Graduate) is a categorical variable (having Yes/No).

To use that variable in regression Analysis, we have to convert it into a quantitative variable.

We can make a categorical variable into a quantitative variable through a dummy variable recoding, by replacing Yes with 1 and No with 0.

Now the data table looks like

x_1 SAT	x_2 GPA	y Graduate	$Y =$	$X =$
1300	2.7		[1]	[1 1300 2.7]
1260	3.7		[1]	[1 1260 3.7]
1220	2.9		[1]	[1 1220 2.9]
1180	2.5		[1]	[1 1180 2.5]
1060	3.9		[1]	[1 1060 3.9]
1140	2.1	0	[0]	[1 1140 2.1]
1100	3.5	0	[0]	[1 1100 3.5]
1020	3.3	0	[0]	[1 1020 3.3]
980	2.3	0	[0]	[1 980 2.3]
940	3.1	0	[0]	[1 940 3.1]

let $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ where,

$$\hat{\beta} = (\mathbf{x}^T \mathbf{x})^{-1} (\mathbf{x}^T Y) = \begin{bmatrix} -3.8392 \\ 0.0032 \\ 0.2395 \end{bmatrix}$$

$$Y = -3.8392 + 0.0032 \text{ SAT} + 0.2395 \text{ GPA}$$

Then $Y = 0.05335$ and the value is near to zero
 \therefore The candidate will not be graduated.

Note:-

If the recoding of the categorical variable changes then the parametric values in the discriminant function changes. To have the commonality in the discriminant function, we apply the Fisher's Linear Discriminant Function.

17/08/2022

Eg:-

and their qualities are measured in terms of curvature and diameter. Results of quality control by experts are given in the following table.

<u>Curvature</u>	<u>Diameter</u>	<u>Quality control result</u>
2.95	6.63	Passed
2.53	7.71	Passed
3.57	5.65	Passed
3.16	5.47	Passed
2.58	7.46	Not Passed
2.16	6.22	Not Passed
3.27	3.32	Not Passed

The new chip ring has curvature 2.81 and diameter 5.46. Can you solve this problem by employing Fisher's Linear Discriminant Analysis.

Sol:-

Step-1:

$$X = \begin{bmatrix} 2.95 & 6.63 \\ 2.53 & 7.79 \\ 3.57 & 5.65 \\ 3.16 & 5.47 \\ 2.58 & 4.46 \\ 2.16 & 6.22 \\ 3.27 & 3.52 \end{bmatrix}$$

$$\mu = [2.89 \ 5.68]$$

$$Y = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 2 \\ 2 \\ 2 \end{bmatrix}$$

$$X_1 = \begin{bmatrix} 2.95 & 6.63 \\ 2.53 & 7.79 \\ 3.57 & 5.65 \\ 3.16 & 5.47 \end{bmatrix}$$

$$\mu_1 = [3.0525 \ 6.385]$$

$$X_2 = \begin{bmatrix} 2.58 & 4.46 \\ 2.16 & 6.22 \\ 3.27 & 3.52 \end{bmatrix}$$

$$\mu_2 = [2.67 \ 4.733]$$

Step-2: Polled co-variance Matrix.

$$C = \frac{(x - \mu)^T(x - \mu)}{n}$$

$$x - \mu = \begin{bmatrix} 0.06 & 0.95 \\ -0.36 & 2.11 \\ 0.68 & -0.03 \\ 0.27 & -0.21 \\ -0.31 & -1.22 \\ -0.73 & 0.54 \\ 0.38 & -2.16 \end{bmatrix}$$

$$C = \begin{bmatrix} 0.2059 & -0.2309 \\ -0.2309 & 1.6921 \end{bmatrix}$$

$$\bar{C}^{-1} = \begin{bmatrix} 5.7342 & 0.7824 \\ 0.7824 & 0.6977 \end{bmatrix}$$

Fisher's Linear Discriminant Function is given by

$$f_i = \mu_i^T \bar{C}^{-1} x_k^T - \frac{1}{2} \mu_i^T \bar{C}^{-1} \mu_i + \ln(p_i)$$

$$f_1 = \mu_1^T \bar{C}^{-1} x_k^T - \frac{1}{2} \mu_1^T \bar{C}^{-1} \mu_1 + \ln(p_1)$$

$$f_1 = [3.0525 \ 6.385] \begin{bmatrix} 5.7342 & 0.7824 \\ 0.7824 & 0.6977 \end{bmatrix} \begin{bmatrix} 2.81 \\ 5.46 \end{bmatrix}$$

$$-\frac{1}{2} [3.0525 \quad 6.385] \begin{bmatrix} 5.7342 & 0.7824 \\ 0.7824 & 0.6977 \end{bmatrix} \begin{bmatrix} 3.0525 \\ 6.385 \end{bmatrix} + \ln\left(\frac{4}{7}\right)$$

$$f_1 = 100.5862 - 56.186 + \left(\frac{1}{2} \cdot 0.5596 \right) = 43.8406$$

$$f_2 = u_2 \bar{c}' x_K^T - \frac{1}{2} u_2 \bar{c}' u_2^T + \ln\left(\frac{3}{7}\right)$$

$$f_2 = [2.67 \quad 4.733] \begin{bmatrix} 5.7342 & 0.7824 \\ 0.7824 & 0.6977 \end{bmatrix} \begin{bmatrix} 2.67 \\ 4.733 \end{bmatrix}$$

$$-\frac{1}{2} [2.67 \quad 4.733] \begin{bmatrix} 5.7342 & 0.7824 \\ 0.7824 & 0.6977 \end{bmatrix} \begin{bmatrix} 2.67 \\ 4.733 \end{bmatrix} + \ln\left(\frac{4}{7}\right)$$

$$f_2 = 43.86$$

$\therefore f_2 > f_1$, so we classify the new chip ring into the second group.

i.e., the new chip ring is not passed.

Eg: classify the following

samples of class $w_1 = x_1 = (x_1, x_2)$

$$= \{(4,2)(9,4)(2,3)(3,6)(4,4)\}$$

$w_2 = x_2 = (x_1, x_2)$

$$= \{(9,10)(6,8)(9,5)(8,7)(10,8)\}$$

$x = 4$

$$(4)u_1 + (4)u_2 = (4,4)$$

$$(9)u_1 + (6)u_2 = (9,6)$$

$$X = \begin{bmatrix} 4 & 2 \\ 2 & 4 \\ 2 & 3 \\ 3 & 6 \\ 4 & 4 \\ 9 & 10 \\ 6 & 8 \\ 9 & 5 \\ 8 & 7 \\ 10 & 8 \end{bmatrix} \quad x_1 = \begin{bmatrix} 4 & 2 \\ 2 & 4 \\ 2 & 3 \\ 3 & 6 \\ 4 & 4 \end{bmatrix} \quad x_2 = \begin{bmatrix} 9 & 10 \\ 6 & 8 \\ 9 & 5 \\ 8 & 7 \\ 10 & 8 \end{bmatrix}$$

$$U_1 = [3 \ 3.8] \quad U_2 = [8.4 \ 7.6]$$

$$U = [5.7 \ 5.7]$$

The pooled covariance matrix,

$$\boxed{C = \frac{(x - U)^T(x - U)}{n}}$$

$$x - U = \begin{bmatrix} -1.7 & -3.7 \\ -3.7 & -1.7 \\ -3.7 & -2.7 \\ -2.7 & 0.3 \\ -1.7 & -1.7 \\ 3.3 & 4.3 \\ 0.3 & 2.3 \\ 3.3 & -0.7 \\ 2.3 & 1.3 \\ 4.3 & 2.3 \end{bmatrix}$$

$$C = \begin{bmatrix} 8.61 & 5.01 \\ 5.01 & 5.81 \end{bmatrix}$$

$$\bar{C}^{-1} = \begin{bmatrix} 0.2033 & -0.1494 \\ -0.1494 & 0.2567 \\ -0.2010 & 0.3454 \end{bmatrix}$$

Fisher's Linear Discriminant Function is given by

$$\boxed{f_i = U_i \bar{C}^{-1} x_k^T - \frac{1}{2} U_i \bar{C}^{-1} U_i^T + \ln(P_i)}$$

$$f_1 = U_1 \bar{C}^{-1} x_k^T - \frac{1}{2} U_1 \bar{C}^{-1} U_1^T + \ln(P_1)$$

$$f_1 = [3 \ 3.8] \begin{bmatrix} 0.2331 & -0.2010 \\ -0.2010 & 0.3454 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix} - \frac{1}{2} [3 \ 3.8]$$

$$\begin{bmatrix} 0.2331 & -0.2010 \\ -0.2010 & 0.3454 \end{bmatrix} \begin{bmatrix} 3 \\ 3.8 \end{bmatrix} + \ln\left(\frac{5}{10}\right)$$

$$= 3.9346 - 1.2513 + (-0.6931) = 1.9902$$

$$f_2 = U_2 C^{-1} x_k^T - \frac{1}{2} U_2 C^{-1} U_2^T + \ln\left(\frac{5}{10}\right)$$

$$\equiv [8.4 \ 7.6] \begin{bmatrix} 0.2331 & -0.2010 \\ -0.2010 & 0.3454 \end{bmatrix} \begin{bmatrix} 5 \\ 6 \end{bmatrix}$$

$$- \frac{1}{2} [8.4 \ 7.6] \begin{bmatrix} 0.2331 & -0.2010 \\ -0.2010 & 0.3454 \end{bmatrix} \begin{bmatrix} 8.4 \\ 7.6 \end{bmatrix} + \ln\left(\frac{5}{10}\right)$$

$$= 7.7720 - 5.3670 - 0.6931 = 1.7119$$

$\therefore f_1 > f_2$, the observation (3,6) is classify into first group.

③ compute the linear discriminant function and classify the observation 5.1 & 3.2.

Samples of class $w_1 = \{(1,2) (2,3) (3,3) (4,5) (5,5)\}$

$w_2 = \{(4,2) (5,0) (5,2) (3,2) (5,3) (6,3)\}$

$$x = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 3 \\ 4 & 5 \\ 5 & 5 \\ 4 & 2 \\ 5 & 0 \\ 5 & 2 \\ 3 & 2 \\ 3 & 3 \end{bmatrix}$$

$$U = [3.909 \ 2.7272]$$

$$x_1 = \begin{bmatrix} 1 & 2 \\ 2 & 3 \\ 3 & 3 \\ 4 & 5 \\ 5 & 5 \end{bmatrix}$$

$$U_1 = [3 \ 3.6]$$

$$x_2 = \begin{bmatrix} 4 & 2 \\ 5 & 0 \\ 5 & 2 \\ 3 & 2 \\ 5 & 3 \\ 6 & 3 \end{bmatrix}$$

$$U_2 = [4.666 \ 2]$$

The pooled covariance matrix,

$$C = \frac{(x - u)^T (x - u)}{n}$$

$$x - u = \begin{bmatrix} -2.909 & -0.7272 \\ -1.909 & 0.2728 \\ -0.909 & 0.2728 \\ 0.1091 & 0.2728 \\ 1.1091 & 0.2728 \\ -0.091 & -0.7272 \\ 1.091 & -0.7272 \\ 1.091 & -0.7272 \\ -0.909 & -0.7272 \\ 1.091 & 0.2728 \\ 2.091 & 0.2728 \end{bmatrix}$$

$$C = \begin{bmatrix} 0.0826 & 0.1570 \\ 0.1570 & 1.8347 \end{bmatrix}$$

$$\bar{C} = \begin{bmatrix} 0.4832 & -0.0413 \\ -0.0413 & 0.5485 \end{bmatrix}$$

$$u = [0.0744]$$

Fisher's linear Discriminant Function is given by

$$f_i = u_i^T \bar{C}^{-1} x_k^T - \frac{1}{2} u_i^T \bar{C} u_i + \ln(P_i)$$

$$f_1 = u_1^T \bar{C}^{-1} x_k^T - \frac{1}{2} u_1^T \bar{C} u_1 + \ln(P_1)$$

$$f_1 = [3 3.6] \begin{bmatrix} 0.4832 & -0.0413 \\ -0.0413 & 0.5485 \end{bmatrix} \begin{bmatrix} 5.1 \\ 3.2 \end{bmatrix}$$

$$- \frac{1}{2} [3 3.6] \begin{bmatrix} 0.4832 & -0.0413 \\ -0.0413 & 0.5485 \end{bmatrix} \begin{bmatrix} 3 \\ 3.6 \end{bmatrix} + \ln\left(\frac{5}{11}\right)$$

$$= 6.486$$

$$f_2 = [8 11.1] = 11.1$$

$$[0.8 8] = 11.1$$

$$\therefore f_2 > f_1$$

Therefore, the observation (5.1, 3.2) is classified into second group.

④ Consider the 2 datasets

$$x_1 = \begin{bmatrix} 3 & 7 \\ 2 & 4 \\ 4 & 7 \end{bmatrix}, x_2 = \begin{bmatrix} 6 & 9 \\ 5 & 7 \\ 4 & 8 \end{bmatrix} \text{ for which } \bar{x}_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix} \text{ & } \bar{x}_2 = \begin{bmatrix} 5 \\ 8 \end{bmatrix}$$

and pooled = $\begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$. classify the observation

$x_0' = [2 \ 7]$ as a population π_1 (or) π_2 .

Sol: Given,

$$x_1 = \begin{bmatrix} 3 & 7 \\ 2 & 4 \\ 4 & 7 \end{bmatrix}, x_2 = \begin{bmatrix} 6 & 9 \\ 5 & 7 \\ 4 & 8 \end{bmatrix}, c = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$$

$$u_1 = [3 \ 6], u_2 = [5 \ 8], \bar{c}' = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix}$$

Fisher's linear discriminant function:-

$$f_i = u_i \bar{c}' x_k^T - \frac{1}{2} u_i \bar{c}' u_i^T - \ln(p_i)$$

$$f_1 = u_1 \bar{c}' x_k^T - \frac{1}{2} u_1 \bar{c}' u_1^T - \ln(p_1)$$

$$= [3 \ 6] \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 7 \end{bmatrix} - \frac{1}{2} [3 \ 6] \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 3 \\ 6 \end{bmatrix} + \ln\left(\frac{3}{6}\right)$$

$$= 21 - \frac{1}{2}(18) - 0.693 = 11.36$$

$$f_2 = u_2 \bar{c}' x_k^T - \frac{1}{2} u_2 \bar{c}' u_2^T - \ln(p_2)$$

$$= [5 \ 8] \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 7 \end{bmatrix} - \frac{1}{2} [5 \ 8] \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 5 \\ 8 \end{bmatrix} - \ln\left(\frac{3}{6}\right)$$

$$= 7.30$$

$f_1 > f_2$

the observation $(2, 7)$ into first group (or) population.
∴ It is into π_1 population.

- ⑤ Consider the observations on $p=2$ variables from $g=3$ populations assuming that the population have a common covariance matrix (Σ). Given that $P_1 = P_2 = 0.25$ and $P_3 = 0.50$ classify $[2 \ 7]$ using discriminant scores.

$$x_1 = \begin{bmatrix} -2 & 5 \\ 0 & 3 \\ -1 & 1 \end{bmatrix} \quad x_2 = \begin{bmatrix} 0 & 6 \\ 2 & 4 \\ 1 & 2 \end{bmatrix} \quad x_3 = \begin{bmatrix} 1 & -2 \\ 0 & 0 \\ -1 & -4 \end{bmatrix}$$

$$\mu_1 = [-1 \ 3] \quad \mu_2 = [1 \ 4] \quad \mu_3 = [0 \ -2]$$

$$C = \frac{(x - \mu)^T(x - \mu)}{n}$$

$$C = \begin{bmatrix} 1.33 & 0.11 \\ 0.11 & 9.56 \end{bmatrix}$$

$$\bar{C}^{-1} = \begin{bmatrix} 0.75 & -0.008 \\ -0.008 & 0.10 \end{bmatrix}$$

$$x = \begin{bmatrix} -2 & 5 \\ 0 & 3 \\ 1 & 1 \\ 0 & 6 \\ 2 & 4 \\ 1 & 2 \\ 1 & -2 \\ 0 & 0 \\ -1 & -4 \end{bmatrix} \quad x - \mu = \begin{bmatrix} -2 & 3.33 \\ 0 & 1.33 \\ 1 & 0.66 \\ 0 & 4.33 \\ 2 & 2.33 \\ 1 & 0.33 \\ 1 & -3.66 \\ 0 & -1.66 \\ -1 & -5.66 \end{bmatrix}$$

$$\mu = [0 \ 1.66]$$

Fisher's linear Discriminant Function is given by

$$f_i = \mu_i^T \bar{C}^{-1} x_K^T - \frac{1}{2} \mu_i^T \bar{C}^{-1} \mu_i + \ln(P_i)$$

$$f_1 = \mu_1^T \bar{C}^{-1} x_K^T - \frac{1}{2} \mu_1^T \bar{C}^{-1} \mu_1 + \ln(P_1)$$

$$= [-1 \ 3] \begin{bmatrix} 0.75 & -0.008 \\ -0.008 & 0.10 \end{bmatrix} \begin{bmatrix} -2 \\ -1 \end{bmatrix} - \frac{1}{2} [-1 \ 3] \begin{bmatrix} 0.75 & -0.008 \\ -0.008 & 0.10 \end{bmatrix} \begin{bmatrix} -1 \\ 3 \end{bmatrix} + \ln(0.25)$$

$$= 1.24 - \frac{1}{2}(1.698) - 1.386$$

$$= -1.028$$

$$f_2 = \mu_2 \bar{C}^T x_K^T - \frac{1}{2} \mu_2 \bar{C}^T \mu_2^T + \ln(P_2)$$

$$= [1 \ 4] \begin{bmatrix} 0.75 & -0.008 \\ -0.008 & 0.10 \end{bmatrix} \begin{bmatrix} -2 \\ -1 \end{bmatrix} - \frac{1}{2} [1 \ 4] \begin{bmatrix} 0.75 & -0.008 \\ -0.008 & 0.10 \end{bmatrix} \begin{bmatrix} 1 \\ 4 \end{bmatrix} + \ln(0.25)$$

$$= -1.828 - \frac{1}{2}(2.286) - 1.386 = -4.357$$

$$f_3 = \mu_3 \bar{C}^T x_K^T - \frac{1}{2} \mu_3 \bar{C}^T \mu_3^T + \ln(P_3)$$

$$= [0 \ -2] \begin{bmatrix} 0.75 & -0.008 \\ -0.008 & 0.10 \end{bmatrix} \begin{bmatrix} -2 \\ -1 \end{bmatrix} - \frac{1}{2} [0 \ -2] \begin{bmatrix} 0.75 & -0.008 \\ -0.008 & 0.10 \end{bmatrix} \begin{bmatrix} 0 \\ -2 \end{bmatrix} + \ln(0.5)$$

$$= -0.728$$

$$\therefore f_3 > f_1 > f_2,$$

The observation $(-2, -1)$ belongs to the 3rd group.

24/08/2022

- ⑥ Consider the observations on $P=2$ variables from $g=3$ populations. Assuming that the populations have common covariance matrix (Σ). Find the fisher's discriminants, the data path

$$x_1 = \begin{bmatrix} -2 & 5 \\ 0 & 3 \\ -1 & 1 \end{bmatrix} \quad x_2 = \begin{bmatrix} 0 & 6 \\ 2 & 4 \\ 1 & 2 \end{bmatrix} \quad x_3 = \begin{bmatrix} 1 & -2 \\ 0 & 0 \\ -1 & -4 \end{bmatrix}$$

$$\bar{x}_1 = \begin{bmatrix} -1 \\ 3 \end{bmatrix} \quad \bar{x}_2 = \begin{bmatrix} 1 \\ 4 \end{bmatrix} \quad \bar{x}_3 = \begin{bmatrix} 0 \\ -2 \end{bmatrix}$$

Find ii) solve $S \leq \min(g-1, P)$

i) \bar{x} , B and w

iii) and eigen values & fisher's discriminant functions

isdt: let $x = \begin{bmatrix} -2 & 5 \\ 0 & 3 \\ -1 & 1 \\ 0 & 6 \\ 2 & 4 \\ 1 & 2 \\ 1 & -2 \\ 0 & 0 \\ -1 & -4 \end{bmatrix}$ $B = \frac{1}{2}(\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})^T$

$$\bar{x}_i - \bar{x} = \begin{bmatrix} -1 & 1.33 \\ 1 & 2.33 \\ 0 & -3.66 \end{bmatrix}$$

$$B = \begin{bmatrix} 2 & 1 \\ 1 & 20.67 \end{bmatrix}$$

$$W = \frac{1}{2}(\bar{x}_{ij} - \bar{x}_{ij})(\bar{x}_{ij} - \bar{x}_{ij})^T$$

$$\bar{x} = \begin{bmatrix} 0 \\ 1.66 \end{bmatrix}$$

$$\bar{x}_{ij} - \bar{x}_i = \begin{bmatrix} -2+1 & 5-3 \\ 0+1 & 3-3 \\ -1+1 & 1-3 \\ 0-1 & 6-4 \\ 2-1 & 4-4 \\ 1-1 & 2-4 \\ 1-0 & -2+2 \\ 0-0 & 0+2 \\ -1-0 & -4+2 \end{bmatrix} = \begin{bmatrix} -1 & 2 \\ 1 & 0 \\ 0 & -2 \\ -1 & 2 \\ 1 & 0 \\ 0 & -2 \\ 1 & 0 \\ 0 & 2 \\ -1 & -2 \end{bmatrix} \quad W = \begin{bmatrix} 6 & -2 \\ -2 & 24 \end{bmatrix}$$

ii) \rightarrow

Note:- To solve $S \leq \min(g-1, P)$

$$S \leq \min(3-1, 2) = \min(2, 2) = 2$$

we have to find two non-zero eigen values for $W^T B$.

$$\bar{W}^T B = \frac{1}{140} \begin{bmatrix} 24 & 2 \\ 2 & 6 \end{bmatrix} * \begin{bmatrix} 2 & 1 \\ 1 & 20.67 \end{bmatrix}$$

$$\bar{W}^T B = \begin{bmatrix} 0.3571 & 0.4667 \\ 0.0714 & 0.8997 \end{bmatrix}$$

To find the eigen values of $\bar{W}^T B$

$$|\bar{W}^T B - \lambda I| = 0$$

$$\left| \begin{bmatrix} 0.3571 & 0.4667 \\ 0.0714 & 0.8997 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right| = 0$$

$$\left| \begin{array}{cc} 0.3571 - \lambda & 0.4667 \\ 0.0714 & 0.8997 - \lambda \end{array} \right| = 0$$

$$(0.3571 - \lambda)(0.8997 - \lambda) - (0.0714)(0.4667) = 0.$$

$$\lambda^2 - 1.2571\lambda + 0.2881 = 0$$

$$\frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$\boxed{\lambda_1 = 0.9556, \lambda_2 = 0.3015}$$

To find the eigen vectors,

$$\begin{pmatrix} 0.3571 - \lambda & 0.4667 \\ 0.0714 & 0.9 - \lambda \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

$$\begin{pmatrix} -0.5985 & 0.4667 \\ 0.0714 & -0.0555 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

$$x_1(-0.5985) + (0.4667)x_2 = 0 \rightarrow ① \quad \checkmark$$

$$0.0714x_1 - 0.0555x_2 = 0 \rightarrow ②$$

$$x_1 = 0, x_2 = 0$$

$$\frac{1 - 0.5985}{\sqrt{(-0.5985)^2 + (0.4667)^2}}, \frac{0.4667}{\sqrt{(-0.5985)^2 + (0.4667)^2}}$$

$\boxed{[0.7886, 0.6149]}$

is a set of eigen vector for $\lambda_1 = 0.9556$.

$$\lambda_2 = 0.3015$$

$$\begin{pmatrix} 0.0556 & 0.4667 \\ 0.0714 & 0.5985 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

$$0.0556 x_1 + 0.4667 x_2 = 0 \quad \checkmark$$

$$0.0714 x_1 + 0.5985 x_2 = 0$$

$$x_1 = 0, x_2 = 0$$

$$\frac{1 - 0.5985}{\sqrt{(0.0556)^2 + (0.4667)^2}}, \frac{0.4667}{\sqrt{(0.0556)^2 + (0.4667)^2}}$$

$$\boxed{[0.1184, 0.9930]}$$

is a set of eigen vector for $\lambda_2 = 0.3015$.

Now, the fishes discriminant function

$$f_1 = 0.7886 x_1 + 0.6149 x_2$$

$$f_2 = 0.1184 x_1 + 0.9930 x_2$$

$$\textcircled{1} \leftarrow 0 = x(F334.0) + (2882.0), x$$

$$\textcircled{4} \leftarrow 0 = x(2220.0) - x(1150.0)$$

25/8
 Q) Classify a new observation into one of the populations, let us assign an observation x_0 into one of $g=3$ populations π_1, π_2, π_3 . Given the following hypothetical probabilities, misclassification class and density values.

		True Population		
		π_1	π_2	π_3
Classify as: π_1	π_1	$c(1 1)=0$	$c(1 2)=500$	$c(1 3)=100$
	π_2	$c(2 1)=10$	$c(2 2)=0$	$c(2 3)=50$
	π_3	$c(3 1)=50$	$c(3 2)=200$	$c(3 3)=0$
Pair Probab		$P_1 = 0.05$	$P_2 = 0.60$	$P_3 = 0.35$
Densities at x_0		$f_1(x_0) = 0.01$	$f_2(x_0) = 0.85$	$f_3(x_0) = 2$

Use minimum expected cost of Misclassification (ECM) procedures to classify a new observation x_0 .

Procedure 1:-

ECM based on pair probabilities, densities and misclassification costs.

$$\sum_{\substack{i=1 \\ i \neq K}}^n P_i \cdot f_i(x_0) \cdot c(K|i)$$

Allot the new observation to the smallest value.

$$\underline{k=1}, \quad P_2 f_2(x_0) \cdot C(1/2) + P_3 f_3(x_0) \cdot C(1/3)$$

$$0.60 * 0.85 * 500 + 0.35 * 2 * 100 = 325$$

$$\underline{k=2}, \quad P_1 f_1(x_0) \cdot C(2/1) + P_3 f_3(x_0) \cdot C(2/3)$$

$$0.05 * 0.01 * 10 + 0.35 * 2 * \frac{100}{50} = \frac{70.005}{35.0005}$$

$$\underline{k=3}, \quad P_1 f_1(x_0) \cdot C(3/1) + P_2 f_2(x_0) \cdot C(3/2)$$

$$0.05 * 0.01 * 50 + 0.60 * 0.85 * \frac{100}{50} = 102.025.$$

If class of misclassification are equal then the formula based on is given by

$$P_1 f_1(x_0) = 0.0005$$

$$P_2 f_2(x_0) = 0.510$$

$$P_3 f_3(x_0) = 0.7$$

Allot the observation to the biggest value, so allot the observation to the third population.

If all classes of misclassification are equal then we use baye's theorem.

$$P(\pi_1/x_0) = \frac{P_1 f_1(x_0)}{\sum P_i f_i(x_0)}$$

$$= \frac{0.05 * 0.01}{(0.05 * 0.01) + (0.60 * 0.85) + (0.35 * 2)} = 0.0004$$

$$P(\pi_2/x_0) = \frac{P_2 f_2(x_0)}{\sum_{i=1}^3 P_i f_i(x_0)} = \frac{0.60 * 0.85}{...} = 0.4197$$

$$P(\pi_3/x_0) = 0.576$$

Allot the observation to the biggest value.
So allot to the third population.

Confusion Matrix:

The apparent error rate in the confusion matrix is calculated by

$$\frac{D_{1m} + D_{2m}}{D_1 + D_2}$$

		Predicted class		
		π_1	π_2	
Actual class	π_1	D_{1c} ↓ Popul. correct classification	D_{1m} ↓ Popul. misclassification	D_1
	π_2	D_{2m} ↓ Popul. misclassification	D_{2c} ↓ Pop. correct classification	D_2

Eg:- Calculate apparent error rate.

		Predicted		
		π_1	π_2	
Actual class	π_1	10	2	12
	π_2	2	10	12

$$\therefore \text{Apparent error rate} = \frac{2+2}{12+12} = 16.7\%$$

Q-4 Principal Component Analysis - (Data Reductⁿ)

- Regardless ones area of study, the first decision that the researcher bases concern to which variables are to be measured. Generally the researcher adopts a heat and miss strategy, collecting information on the typically large number of variables that might be relevant.
- How many variables are to be measured, some practical problems may arise for e.g. with as . as 10 variables there are $^{10C_2} = 45$ correlation coefficients, with 20 variables, there are $^{20C_2} = 190$ correlation coefficients to be considered.
- so, the number of correlation coefficients keeps on accelerating as the number of variables increases. Obviously, with large numbers of variables, the correlation coefficients are so large as to beyond the comprehension and some data reduction technique that can systematically summarize the large correlation matrices is needed which is possible with the principle component analysis.

Q:- Is there any technique which serves the purpose of data reduction?

Sol:- The data reduction is possible by PCA (Principle Component Analysis).

- ① The basic idea of principal component analysis is to describe the variation of a set of multivariate data in terms of a set of uncorrelated variables each of which is a particular linear combination of a set of correlated variables.
- ② The new variables are derived in decreasing order of importance so that the first transformed variable accounts for as much as possible of the variation in the original data.
- ③ The usual objective of this type of analysis is to see whether the first few components account for most of the variation in the original data.
- ④ If they do so, then they can be used to summarize the data with little loss of information, thus providing a reduction in the dimensionality of the data which might be useful in simplifying the later analysis.

Simple Algorithm for PCA:-

$$(u - \bar{u})(u - \bar{u})^T = \Sigma$$

Step 1 :- Standardize the data set

Step 2 :- calculate the covariance matrix of the Standardized dataset

Step 3 :- Calculate the eigen values and eigen vectors of the covariance matrix.

Step 4 :- sort the eigen values and their corresponding eigen vectors.

Step 5: Pick up k eigen values to form a eigen matrix of eigen vectors

Step 6: Transform the original matrix to find the new data set

Steps to PCA:

Step 1: Decide the variables to be studied

For eg:- let x_1, x_2, \dots, x_p be the p number of variables under step

Step 2: collect the observations under each variable

$$\begin{matrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{matrix}$$

Step 3: calculate the variance covariance matrix which is given by

$$C = \frac{(x - \mu)^T (x - \mu)}{n}$$

Step 4: obtain the eigen values and eigen vectors of the variance covariance matrix

Let $\lambda_1, \lambda_2, \dots, \lambda_p$ be the eigen values

Sort the eigen values in the descending order.

$$\lambda_{(1)} \geq \lambda_{(2)} \geq \dots \geq \lambda_{(p)}$$

Find the corresponding eigen vectors of the eigen values.

Step 5:- How many principal components to retain?

Principal component	Variance explained	
z_1	λ_1	$\frac{\lambda_1}{\sum \lambda_i} \times 100$
z_2	λ_2	$\frac{\lambda_1 + \lambda_2}{\sum \lambda_i} \times 100$
\vdots	\vdots	\vdots
z_p	λ_p	$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_p}{\sum \lambda_i} \times 100$

Note:-

Generally we take the thresholds as the 90%, 95% or 99% (depends upon the researcher or the statistician)

Step 6: Transform the variables into new set of reduced variables say z_1, z_2, \dots, z_a where ($a < p$)

Q5: Find the principal components of the following data.

$$\begin{array}{l} \text{let } x = \begin{bmatrix} 2 & 4 \\ 1 & 3 \\ 0 & 1 \\ -1 & 0.5 \end{bmatrix}, \quad x - \mu = \begin{bmatrix} 1.5 & 1.875 \\ 0.5 & 0.875 \\ -0.5 & -1.125 \\ -1.5 & -1.625 \end{bmatrix} \\ \mu = [0.5 \quad 2.125] \end{array}$$

• Variance covariance matrix is given by

$$C = \frac{(x - u)^T(x - u)}{n}$$

$$C = \begin{bmatrix} 1.25 & 1.5625 \\ 1.5625 & 2.0468 \end{bmatrix}$$

To find the eigen values of variance covariance matrix C is

$$(C - \lambda I) = 0$$

$$\begin{vmatrix} 1.25-\lambda & 1.5625 \\ 1.5625 & 2.0468-\lambda \end{vmatrix} = 0$$

$$(1.25-\lambda)(2.0468-\lambda) - (1.5625)^2 = 0$$

$$2.5585 - 1.25\lambda - 2.0468\lambda + \lambda^2 - 2.4414 = 0$$

$$\lambda^2 - 3.2968\lambda + 0.1171 = 0$$

$$\boxed{\lambda_1 = 3.2608 \quad \lambda_2 = 0.0359}$$

To find the eig
Principal component Analysis Table:- (How many components are to be retained):

Principal component	Variance explained	Cumulative proportion of total variance
z_1	$\lambda_1 = 3.261$	$\frac{3.261}{3.2969} \times 100 = 98.91\%$
z_2	$\lambda_2 = 0.0359$	$\frac{3.2969}{3.2969} \times 100 = 100\%$
	$\sum \lambda_i = 3.2969$	

Note:- The threshold is fixed as 99.1. (depends upon the researcher)

Here we have to retain 2 principal components.
To find the Eigen vectors for $\lambda_1 = 3.261$

$$\begin{pmatrix} 1.25 - 3.261 & 1.5625 \\ 1.5625 & 2.0468 - 3.261 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

$$-2.011x_1 + 1.5625x_2 = 0 \rightarrow ①$$

$$1.5625x_1 - 1.2142x_2 = 0 \rightarrow ②$$

$$x_1 = 0, x_2 = 0$$

$$+2.011x_1 = 1.5625x_2$$

$$\frac{x_1}{1.5625} = \frac{x_2}{2.011}$$

Normalization,

$$\frac{1.5625}{\sqrt{1.5625^2 + 2.011^2}}, \frac{2.011}{\sqrt{1.5625^2 + 2.011^2}}$$

$$0.6135, 0.7396$$

is a set of eigen vector.

Find a eigen vector for $\lambda_2 = 0.0359$

$$\begin{pmatrix} 1.25 - 0.0359 & 1.5625 \\ 1.5625 & 2.0468 - 0.0359 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

$$1.2141x_1 + 1.5625x_2 = 0 \rightarrow ①$$

$$1.5625x_1 + 0.0109x_2 = 0 \rightarrow ②$$

$$1.2141x_1 = -1.5625x_2$$

$$\frac{x_1}{-1.5625} = \frac{-x_2}{1.2141}$$

Normalization,

$$\frac{-1.5625}{\sqrt{1.5625^2 + 1.2141^2}}, \frac{1.2141}{\sqrt{1.5625^2 + 1.2141^2}}$$

$$[-0.7896, 0.6135]$$

is a set of eigen vector when $\lambda_2 = 0.0359$

$$\therefore z_1 = 0.6135 x_1 + 0.7896 x_2$$

$$z_2 = -0.7896 x_1 + 0.6135 x_2$$

② Find the principal components of the following data.

x_1	x_2	x_3	<u>sd:</u>
90	60	90	
90	90	30	
60	60	60	
60	60	90	
30	30	30	

$$x = \begin{bmatrix} 90 & 60 & 90 \\ 90 & 90 & 30 \\ 60 & 60 & 60 \\ 60 & 60 & 90 \\ 30 & 30 & 30 \end{bmatrix}$$

$$x - \mu = \begin{bmatrix} 24 & 0 & 30 \\ 24 & 30 & -30 \\ -6 & 0 & 0 \\ -6 & 0 & 30 \\ -36 & -30 & -30 \end{bmatrix}$$

$$\mu = [66 \ 60 \ 60]$$

Variance covariance matrix is given by

$$C = \frac{(x - \mu)^T(x - \mu)}{n}$$

$$C = \begin{bmatrix} 504 & 360 & 180 \\ 360 & 3160 & 0 \\ 180 & 0 & 720 \end{bmatrix}$$

$$C = \begin{bmatrix} 504-d & 360 & 180 \\ 360 & 3160-d & 0 \\ 180 & 0 & 720-d \end{bmatrix}$$

To find eigen values, the characteristic eqn $|A - \lambda I| = 0$

$$\Rightarrow \begin{vmatrix} 504-d & 360 & 180 \\ 360 & 3160-d & 0 \\ 180 & 0 & 720-d \end{vmatrix} = 0$$

$$\Rightarrow (504-d)((3160-d)(720-d)) - 360(360)(720-d) + 180(-180)(3160-d) = 0$$

$$|A| = \begin{vmatrix} 504 & 360 & 180 \\ 360 & 360 & 0 \\ 180 & 0 & 720 \end{vmatrix}$$

$$\bar{A}^{-1} = \begin{bmatrix} 0.0101 & -0.0101 & 0.0025 \\ -0.0101 & 0.0128 & 0.0025 \\ -0.0025 & 0.0025 & 0.0202 \end{bmatrix}$$

$$= 2566080$$

$$\boxed{\lambda^3 - \text{trace}(A) \cdot \lambda^2 + \text{trace}(\text{adj}(A)) \cdot \lambda - \det(A) = 0}$$

$$\therefore \text{adj}(A) = \bar{A}^{-1}|A|$$

$$\lambda^3 - 1584\lambda^2 + [20431][2566080] \cdot \lambda - 2566080 = 0$$

641263.391

$$\lambda_1 = 911.027, \lambda_2 = 628.13, \lambda_3 = 44.842$$

are the eigen values.

How many principal components to be Retained?

Principal Component	Variance Explained	Cumulative proportion of total variance
Z_1	$\lambda_1 = 911.027$	$\frac{911.027}{1584} \times 100 = 57.5\%$
Z_2	$\lambda_2 = 628.13$	$\frac{911.027 + 628.13}{1584} \times 100 = 97.17\%$
Z_3	$\lambda_3 = 44.84$	$\frac{911.07 + 628.13 + 44.842}{1584} \times 100 = 100\%$

Note:- If the threshold has fixed as 95%, we retain the first two principal components.

cont. in ^{lab} CS observation books

continuation:

To find the eigen vector when $\lambda_1 = 911.02$

$$\begin{pmatrix} 504 - 911.02 & 360 & 180 \\ 360 & 360 - 911.02 & 0 \\ 180 & 0 & 720 - 911.02 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0$$

$$-407.02x_1 + 360x_2 + 180x_3 = 0 \rightarrow ①$$

$$360x_1 - 551.02x_2 = 0 \rightarrow ②$$

$$180x_1 - 191.02x_3 = 0 \rightarrow ③$$

$$\begin{array}{cccc|c} x_1 & x_2 & x_3 & x_1 \\ \hline -407.02 & 360 & 180 & -407.02 \\ 360 & -551.02 & 0 & 360 \end{array}$$

$$\frac{x_1}{0 + (180)(551.02)} = \frac{x_2}{-(180)(360)} = \frac{x_3}{(360)(360) - (551.02)(407.02)}$$

$$\frac{x_1}{99183.6} = \frac{x_2}{+64800} = \frac{x_3}{-94676.16}$$

Normalizing the value

$$= \frac{99183.6}{\sqrt{9837386509 + 4199040000 + 8963575272}}$$

$$= \frac{99183.6}{151657.51} = 0.6539$$

$0.6539, 0.4292, 0.6210$ is the

eigen vector when $\lambda_1 = 911.02$

To find the eigen vector for $\lambda_2 = 628.13$

$$\begin{pmatrix} 504 - 628.13 & 360 & 180 \\ 360 & 360 - 628.13 & 0 \\ 180 & 0 & 720 - 628.13 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0$$

$$-121.13x_1 + 360x_2 + 180x_3 = 0 \rightarrow ①$$

$$360x_1 - 268.13x_2 = 0 \rightarrow ②$$

$$180x_1 - 91.87x_3 = 0 \rightarrow ③$$

$$\begin{array}{cccc|c} x_1 & x_2 & x_3 & x_1 \\ \hline -121.13 & 360 & 180 & -121.13 \\ 360 & -268.13 & 0 & 360 \\ \hline x_1 & +x_2 & & x_3 \\ \hline 32478.58 - 129600 & -64800 & & -97121.42 \\ 48263.4 & & & \end{array}$$

Normalizing the Value

$$= \frac{48263.4}{\sqrt{2329355780 + 419904000 + 9432570223}}$$

$$= \frac{48263.4}{126336.71} = 0.3820$$

0.3820, 0.5129, -0.7687 is the eigen vector when $\lambda_2 = 628.13$

Now, the principal components equations are

$$z_1 = 0.6539x_1 + 0.4278x_2 + 0.6243x_3$$

$$z_2 = 0.3820x_1 + 0.5164x_2 - 0.7644x_3$$

03/09/22

Q3. Find the principal components as follows.

$$\begin{array}{c}
 \underline{x_1} & \underline{x_2} & \underline{x_3} \\
 = & = & = \\
 7 & 4 & 3 \\
 4 & 1 & 4 \\
 4 & 3 & 5 \\
 6 & 6 & 1 \\
 8 & 5 & 7 \\
 8 & 2 & 9 \\
 5 & 5 & 3 \\
 9 & 4 & 8 \\
 7 & 2 & 5 \\
 8 & 2 & 2
 \end{array}
 \quad x = \begin{bmatrix} 7 & 4 & 3 \\ 4 & 1 & 8 \\ 6 & 3 & 5 \\ 8 & 6 & 1 \\ 8 & 5 & 7 \\ 7 & 2 & 9 \\ 5 & 3 & 3 \\ 9 & 5 & 8 \\ 7 & 4 & 5 \\ 8 & 2 & 2 \end{bmatrix}$$

$$U = [6.9 \ 3.5 \ 5.1]$$

$$x - U = \begin{bmatrix} 0.1 & 0.5 & -2.1 \\ -2.9 & -2.5 & 2.9 \\ -0.9 & -0.5 & 1.5 \\ 1.1 & 2.5 & -4.1 \\ 1.1 & 1.5 & 1.9 \\ 0.1 & -1.5 & 3.9 \\ -1.9 & -0.5 & -2.1 \\ 2.1 & 1.5 & 2.9 \\ 0.1 & 0.5 & -0.1 \\ 1.1 & -1.5 & -3.1 \end{bmatrix}$$

$$C = \frac{(x - U)^T(x - U)}{n}$$

$$C = \begin{bmatrix} 2.09 & 1.45 & -0.39 \\ 1.45 & 2.25 & -1.15 \\ -0.39 & -1.15 & 7.09 \end{bmatrix}$$

To find the eigen values of Variance Covariance matrix:

$|C - \lambda I| = 0$ is the characteristic equation

$$\begin{vmatrix} 2.09-\lambda & 1.45 & -0.39 \\ 1.45 & 2.25-\lambda & -1.15 \\ -0.39 & -1.15 & 7.09-\lambda \end{vmatrix} = 0$$

$$\lambda^3 - \text{trace}(A) \cdot \lambda^2 + \text{trace}(|A| |A^{-1}|) \lambda - |A| = 0$$

↳ characteristic equation for 3×3 matrix.

$$\lambda^3 - 11.43\lambda^2 + 31.896\lambda - 16.6284 = 0$$

$$\lambda_1 = 7.4465 \quad \lambda_2 = 3.3085 \quad \lambda_3 = +0.6749$$

To find the eigen vector,

$$\begin{pmatrix} 2.09 - 7.4465 & 1.45 & -0.39 \\ 1.45 & 2.25 - 7.4465 & -1.15 \\ -0.39 & -1.15 & 7.09 - \lambda \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0$$

$$-5.3565x_1 + 1.45x_2 - 0.39x_3 = 0 \rightarrow ①$$

$$1.45x_1 - 5.1965x_2 - 1.15x_3 = 0 \rightarrow ②$$

$$-0.39x_1 - 1.15x_2 - 0.3572x_3 = 0 \rightarrow ③$$

$$\begin{array}{cccc} x_1 & x_2 & x_3 & x_1 \\ \hline -5.3565 & 1.45 & -0.39 & -5.3565 \\ 1.45 & -5.1965 & -1.15 & 1.45 \\ -0.39 & -1.15 & -0.3572 & -0.39 \end{array}$$

$$\frac{x_1}{27.8335 - 2.1025} = \frac{x_2}{6.1599 + 0.5655} = \frac{x_3}{2.1025 - 27.8335}$$

Normalizing the value

$$= \frac{25.731}{\sqrt{662.0843 + 45.9310 + 662.161}}$$

$$= \frac{25.731}{\sqrt{1180.1153}} = -0.1375$$

-0.1375, -0.2504, 0.9583 is the eigen vector when $\lambda_1 = 7.4465$.

To find the eigen vector for $\lambda_2 = 3.3085$,

$$\begin{pmatrix} 2.09 - 3.3085 & 1.45 & -0.39 \\ 1.45 & 2.25 - 3.3085 & -1.15 \\ -0.39 & -1.15 & 7.09 - 3.3085 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = 0$$

$$-1.2185x_1 + 1.45x_2 - 0.39x_3 = 0 \rightarrow ①$$

$$1.45x_1 - 1.0585x_2 - 1.15x_3 = 0 \rightarrow ②$$

$$-0.39x_1 - 1.15x_2 + 3.7815x_3 = 0 \rightarrow ③$$

$$\begin{array}{cccc} \underline{x_1} & \underline{x_2} & \underline{x_3} & \underline{x_1} \\ -1.2185 & 1.45 & -0.39 & -1.2185 \\ 1.45 & -1.0585 & -1.15 & 1.45 \end{array}$$

\Rightarrow if 0.6990, if 0.6608, if 0.2730 are the eigen vector when $\lambda_2 = 3.3085$.

Q: Find the principal components of the following.

Screent Plot / Elbow Plot:-

There is always a question of how many components to retain. There is no definite answer to this question. Values to consider include the amount of total sample variance explain, the relative sizes of the eigen values. A component with less eigen value deemed unimportant and may indicate an

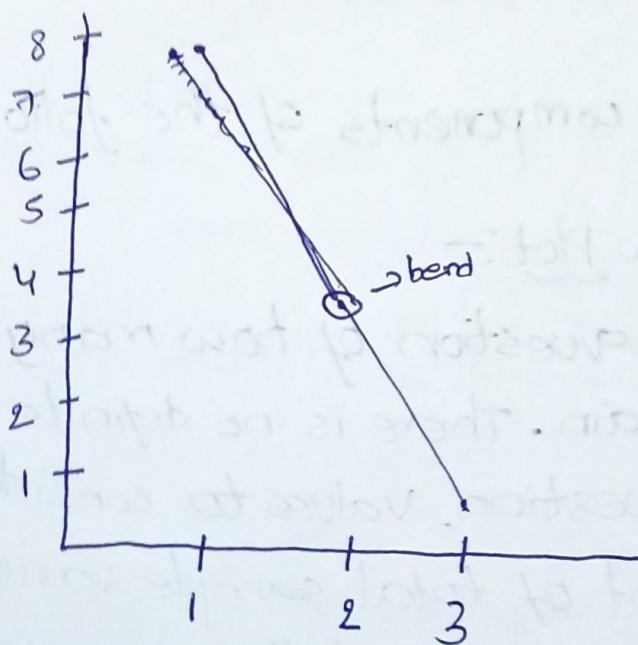
unsuspected linear dependency in the data.

An usual aid to determine an appropriate number of principal components is a scree plot/elbow plot with the eigen values ordered from largest to smallest, a scree plot is a plot of λ_i vs i^{th} magnitude of eigen value and its number. To determine the appropriate number of components we look for an elbow (bend) in the scree plot.

Scree plot for Q3:-

The 3 eigen values are

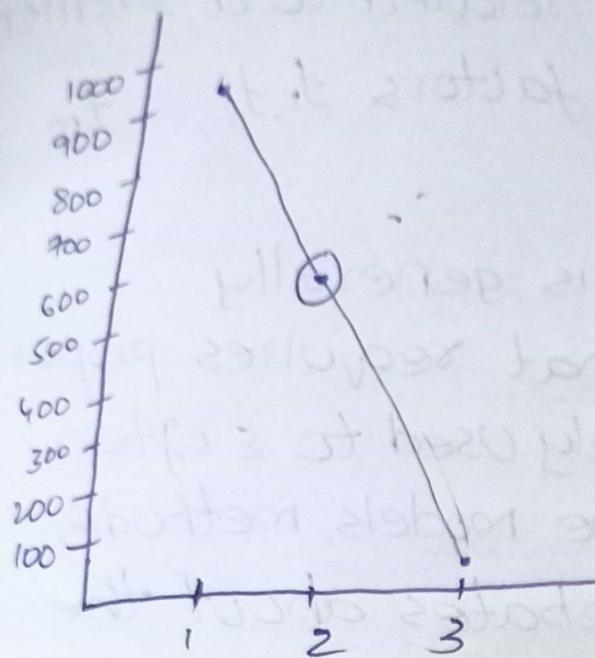
$$\lambda_1 = 7.74465, \lambda_2 = 3.3085, \lambda_3 = 0.6749$$



sreee plot for Q2:-

The 3 eigen values are

$$\lambda_1 = 911.02, \lambda_2 = 688.13, \lambda_3 = 44.84$$



Q4:- Find the principal components for the following data.

$\underline{x_1}$	$\underline{x_2}$	$\underline{x_3}$	$\underline{x_4}$	$\underline{x_5}$	$\underline{x_6}$	$\underline{x_7}$	$\underline{x_8}$	$\underline{x_9}$
16	16	13	18	16	15	14	16	16
18	19	15	16	18	18	18	17	19
17	17	14	14	17	17	20	14	15
17	17	17	16	18	18	16	20	14
16	15	17	17	18	18	19	16	19
15	17	16	17	18	18	15	19	16
17	16	16	18	18	18	17	15	18
20	18	16	20	15	15	19	14	17
14	16	18	17	19	19	18	17	18
16	16	15	19	18	18	18	15	14
18	19	16	14	14	14	17	16	13
19	15	15	18	16	16	18	19	17

Factor Analysis:- (Classification of Groups)

Factor analysis is a method for investigating whether a number of variables of interest x_1, x_2, \dots, x_p are linearly related to a smaller number of unobservable factors f_1, f_2, \dots, f_m ($m < p$).

Factor Analysis is generally exploratory/descriptive that requires proper judgements. It is widely used to & often controversial because the models, methods, are so flexible that debates about the interpretations that occur.

Notations for Factor Analysis:-

Let all the independent variables be denoted by

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \quad u = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_p \end{pmatrix}$$

let the population mean vector is denoted by
Here $E(x_i) = u_i$.

Consider m unobservable common factors f_1, f_2, \dots, f_m denoted by $f = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{pmatrix}$

Factor Analysis Model:-

The factor model can be thought of as a series of multiple regressions predicting each of the observable variables x_i from the values of unobservable common factors f_j .

$$x_1 = u_1 + l_{11}f_1 + l_{12}f_2 + \dots + l_{1m}f_m + \epsilon_1$$

$$x_2 = u_2 + l_{21}f_1 + l_{22}f_2 + \dots + l_{2m}f_m + \epsilon_2$$

:

$$x_p = u_p + l_{p1}f_p + l_{p2}f_2 + \dots + l_{pm}f_m + \epsilon_p$$

$$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_p \end{pmatrix} + \begin{pmatrix} l_{11} & l_{12} & \dots & l_{1m} \\ l_{21} & l_{22} & \dots & l_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ l_{p1} & l_{p2} & \dots & l_{pm} \end{pmatrix} \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_p \end{pmatrix}$$

$$x_{px_1} = u_{px_1} + l_{pxm} \cdot f_{mx_1} + \epsilon_{px_1}$$

$$X - U = Lf - \epsilon \rightarrow ①$$

is known as Factor model.

- If the factors in the F matrix are dependent then it is known as oblique factor analysis.
- If the factors in the F matrix are independent then it is known as orthogonal factor analysis.