

PRACTICALS ON STATISTICS

BY

TANUJIT CHAKRABORTY

Indian Statistical Institute

Mail : tanujitisi@gmail.com

PROBLEMS ON LINEAR ALGEBRA

- 1) (a) Determine k so that the set S is linearly independent in E^3
- $S = \{(1, 2, 1), (k, 3, 1), (2, k, 0)\}$
 - $S = \{(k, 1, 1), (1, k, 1), (1, 1, k)\}$
- (b) For what values of α will the vectors $(1, 5, 7), (4, 0, \alpha), (1, 0, 0)$ form a basis for E^3 ? [C.U. 2009]

ANS :- (i) Construct a matrix with the given vectors :

$$(a) \quad A = \begin{pmatrix} 1 & 2 & 1 \\ k & 3 & 1 \\ 2 & k & 0 \end{pmatrix}$$

As these vectors are linearly independent, so, $\text{rank}(A) = 3$.

$$\therefore |A| \neq 0$$

$$\Rightarrow \begin{vmatrix} 1 & 2 & 1 \\ k & 3 & 1 \\ 2 & k & 0 \end{vmatrix} \neq 0$$

$$\Rightarrow (k^2 - 6) - (k - 4) \neq 0$$

$$\Rightarrow (k^2 - k - 2) \neq 0$$

$$\Rightarrow (k-2)(k+1) \neq 0 \Rightarrow k \neq 2, -1.$$

(ii) Construct a matrix with the given vectors :

$$A = \begin{pmatrix} k & 1 & 1 \\ 1 & k & 1 \\ 1 & 1 & k \end{pmatrix}$$

As these vectors are linearly independent, so, $\text{rank}(A) = 3$

$$\therefore |A| \neq 0$$

$$\Rightarrow (k-1)^2(k+2) \neq 0$$

$$\Rightarrow k \neq -2, 1.$$

- (b) As the vectors $(1, 5, 7), (4, 0, \alpha), (1, 0, 0)$ are linearly independent, so we construct a matrix A with $\text{rank}(A) = 3$.

$$A = \begin{pmatrix} 1 & 5 & 7 \\ 4 & 0 & \alpha \\ 1 & 0 & 0 \end{pmatrix}$$

$$\Rightarrow |A| \neq 0$$

$$\Rightarrow 5\alpha \neq 0$$

$$\Rightarrow \alpha \neq 0$$

2) Extend the set $\{(2, 3, -1), (1, -2, -9)\}$ of vectors to a basis for E^3 . Also find an orthonormal basis for E^3 .

Ans:-

$(0, 0, 1) \in E^3$ and the set $\{(2, 3, -1), (1, -2, -9), (0, 0, 1)\}$ of vectors form a basis of E^3 as the vectors $(0, 0, 1), (2, 3, -1), (1, -2, -9)$ are LIN and they span in E^3 .

$$\tilde{v}_1 = (0, 0, 1)$$

$$\therefore \tilde{u}_1 = \frac{\tilde{v}_1}{\|\tilde{v}_1\|} = \frac{(0, 0, 1)}{\sqrt{0^2 + 0^2 + 1^2}} = (0, 0, 1)$$

$$\begin{aligned}\therefore \tilde{v}_2 &= (2, 3, -1) - \{(2, 3, -1)'(0, 0, 1)\}(0, 0, 1) \\ &= (2, 3, -1) + (0, 0, 1) \\ &= (2, 3, 0)\end{aligned}$$

$$\therefore \tilde{u}_2 = \frac{\tilde{v}_2}{\|\tilde{v}_2\|} = \frac{(2, 3, 0)}{\sqrt{4+9}} = \frac{1}{\sqrt{13}} (2, 3, 0)$$

$$\begin{aligned}\therefore \tilde{v}_3 &= (1, -2, -9) - (\tilde{u}_1' \tilde{v}_3) \tilde{u}_1 - (\tilde{u}_2' \tilde{v}_3) \tilde{u}_2 \\ &= (1, -2, -9) + 9(0, 0, 1) + \frac{1}{\sqrt{13}} \left(\frac{1}{\sqrt{13}} (2, 3, 0) \right) \\ &= (1, -2, -9) + (0, 0, 9) + \left(\frac{8}{13}, \frac{12}{13}, 0 \right)\end{aligned}$$

$$= \left(\frac{21}{13}, -\frac{14}{13}, 0 \right)$$

$$= \frac{7}{13} (3, -2, 0)$$

$$\begin{aligned}\therefore \tilde{u}_3 &= \frac{\tilde{v}_3}{\|\tilde{v}_3\|} = \frac{\frac{7}{13} (3, -2, 0)}{\sqrt{\left(\frac{21}{13}\right)^2 + \left(\frac{14}{13}\right)^2}} = \frac{7}{13} (3, -2, 0) \times \frac{\sqrt{13}}{7\sqrt{13}} \\ &= \frac{1}{\sqrt{13}} (3, -2, 0)\end{aligned}$$

We know, $\{\tilde{u}_1, \tilde{u}_2, \tilde{u}_3\}$ forms orthonormal basis for E^3 , $\{\tilde{v}_1, \tilde{v}_2, \tilde{v}_3\}$ forms orthogonal basis for E^3 .

In general, $\tilde{v}_k = \tilde{q}_k - \sum_{i=1}^{k-1} (u_i' \tilde{q}_k) u_i$ & $\tilde{u}_k = \frac{\tilde{v}_k}{\|\tilde{v}_k\|}$.

∴ the orthonormal basis for E^3 is

$$\left\{ (0, 0, 1), \frac{1}{\sqrt{13}} (2, 3, 0), \frac{1}{\sqrt{13}} (3, -2, 0) \right\}$$

3) Find the dimension of the vector space generated by the vectors:

$$\alpha_1' = \begin{pmatrix} 0 & 1 & 2 & 3 \end{pmatrix}$$

$$\alpha_2' = \begin{pmatrix} 2 & -1 & 5 & 1 \end{pmatrix}$$

$$\alpha_3' = \begin{pmatrix} 4 & 0 & 6 & 1 \end{pmatrix}$$

$$\alpha_4' = \begin{pmatrix} 0 & -2 & 4 & 7 \end{pmatrix}$$

Find a vector in the space orthogonal to the vector space spanned by $\alpha_1', \alpha_2', \alpha_3', \alpha_4'$. [C.U. 2001]

Ans:-

$$A = \begin{pmatrix} 0 & 1 & 2 & 3 \\ 2 & -1 & 5 & 1 \\ 4 & 0 & 6 & 1 \\ 0 & -2 & 4 & 7 \end{pmatrix}$$

Now, to find $\text{rank}(A)$, we will reduce this matrix into its echelon form, and the solution of $E\vec{x} = \vec{0}$ is a vector which is orthogonal to the vector space spanned by $\alpha_1, \alpha_2, \alpha_3, \alpha_4$.

$$\begin{array}{c} A = \begin{pmatrix} 0 & 1 & 2 & 3 \\ 2 & -1 & 5 & 1 \\ 4 & 0 & 6 & 1 \\ 0 & -2 & 4 & 7 \end{pmatrix} \xrightarrow{R_2 \leftrightarrow R_1} \begin{pmatrix} 2 & -1 & 5 & 1 \\ 0 & 1 & 2 & 3 \\ 4 & 0 & 6 & 1 \\ 0 & -2 & 4 & 7 \end{pmatrix} \\ \xrightarrow{R_3' = R_3 - 2R_1} \begin{pmatrix} 2 & -1 & 5 & 1 \\ 0 & 1 & 2 & 3 \\ 0 & 2 & -4 & 7 \\ 0 & -2 & 4 & 7 \end{pmatrix} \xrightarrow{R_4' = R_3 + R_4} \begin{pmatrix} 2 & -1 & 5 & 4 \\ 0 & 1 & 2 & 3 \\ 0 & 2 & -4 & 7 \\ 0 & 0 & 0 & 0 \end{pmatrix} \\ \xrightarrow{R_1' = R_1 + R_2} \begin{pmatrix} 2 & 0 & 7 & 7 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & -8 & -13 \\ 0 & 0 & 0 & 0 \end{pmatrix} \xrightarrow{R_3' = \frac{R_3}{-8}} \begin{pmatrix} 1 & 0 & 7/2 & 7/2 \\ 0 & 1 & 2 & 3 \\ 0 & 0 & 1 & 13/8 \\ 0 & 0 & 0 & 0 \end{pmatrix} = E \end{array}$$

$$\therefore \text{Rank}(A) = 3.$$

$$\therefore \dim(S_4) = 3.$$

Now, $E \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vec{x}_3 \\ \vec{x}_4 \end{pmatrix} = \vec{0}$ gives —

$$\Rightarrow \vec{x}_1 + \frac{7}{2} \vec{x}_3 + \frac{7}{2} \vec{x}_4 = \vec{0}$$

$$\vec{x}_2 + 2 \vec{x}_3 + 3 \vec{x}_4 = \vec{0}$$

$$\vec{x}_3 + \frac{13}{8} \vec{x}_4 = \vec{0}$$

$$\text{Let, } \vec{x}_4 = t, \text{ then } \vec{x}_3 = -\frac{13}{8}t, \vec{x}_2 = \frac{1}{4}t, \vec{x}_1 = \frac{35}{18}t.$$

$$\therefore \vec{x} = t \left(\frac{35}{18}, \frac{1}{4}, -\frac{13}{8}, 1 \right), t \in \mathbb{R}.$$

4) The following vectors from a spanning set for a vector space S_4 :

$$\alpha_1' = (1.000 \quad 0.313 \quad 0.280 \quad 0.156)$$

$$\alpha_2' = (0.333 \quad 1.313 \quad 0.628 \quad 0.315)$$

$$\alpha_3' = (0.309 \quad 0.553 \quad 0.480 \quad 0.165)$$

$$\alpha_4' = (1.024 \quad 1.073 \quad 0.428 \quad 0.306)$$

i) Find the $\dim(S_4)$?

ii) Find a basis for S_4 ?

iii) If $\dim(S_4) < 4$, find a basis for $O(S_4)$? [C.U.1996]

Ans:- Dimension of the vector space (S) spaned by α_i' 's, $i=1(1)4$,
= the rank of A = $\begin{pmatrix} \alpha_1' \\ \alpha_2' \\ \alpha_3' \\ \alpha_4' \end{pmatrix}$

Now,

$$A = \begin{pmatrix} 1.000 & 0.313 & 0.280 & 0.156 \\ 0.333 & 1.313 & 0.628 & 0.315 \\ 0.309 & 0.553 & 0.480 & 0.165 \\ 1.024 & 1.073 & 0.428 & 0.306 \end{pmatrix}$$

$$R_4' = R_3 + R_4 - R_2 - R_1$$

$$R_2' \rightarrow R_2 - 0.333R_1 \quad \begin{pmatrix} 1.000 & 0.313 & 0.280 & 0.156 \\ 0 & 1.2088 & 0.5848 & 0.2630 \\ 0.309 & 0.4563 & 0.3985 & 0.1168 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$R_3' = R_3 - 0.309R_1$$

$$R_2' = R_2 / 1.2088 \quad \begin{pmatrix} 1.000 & 0.313 & 0.280 & 0.156 \\ 0 & 1.000 & 0.4424 & 0.2176 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$$R_3' = R_3 - 0.4563R_2$$

$$R_3' = R_3 / 0.1916 \quad \begin{pmatrix} 1.000 & 0.313 & 0.280 & 0.156 \\ 0 & 1.000 & 0.4424 & 0.2176 \\ 0 & 0 & 1.000 & 0.0913 \\ 0 & 0 & 0 & 0 \end{pmatrix} = E$$

$\therefore \text{Rank}(A) = \text{No. of non-null rows in it's echelon form}$

$$= 3$$

$$\therefore \dim(S) = 3$$

ii) The basis of S is $\{(1, 0.00, 0.318, 0.280, 0.156), (0, 1.000, 0.4424, 0.2176), (0, 0, 1.000, 0.0913)\}$

iii) $E \begin{pmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \alpha_4 \end{pmatrix} = 0$

$$\Rightarrow 1.000\alpha_1 + 0.318\alpha_2 + 0.280\alpha_3 + 0.156\alpha_4 = 0 \\ 1.000\alpha_2 + 0.4424\alpha_3 + 0.2176\alpha_4 = 0 \\ 1.000\alpha_3 + 0.0918\alpha_4 = 0$$

Let, $\alpha_4 = t, t \in \mathbb{R}$.

$$\left\{ \begin{array}{l} \therefore \alpha_3 = -0.0913t \\ \therefore \alpha_2 = -0.1772t \\ \therefore \alpha_1 = -0.0749t \end{array} \right.$$

$$\therefore O(S) = t(-0.0749, -0.1772, -0.0913, 1), t \in \mathbb{R}.$$

5) S and T are subspaces of V_4 given by

$$S = \{(x_1, x_2, x_3, x_4) : 2x_1 + x_2 + 3x_3 + x_4 = 0\}$$

$$T = \{(x_1, x_2, x_3, x_4) : x_1 + 2x_2 + x_3 + 3x_4 = 0\}$$

Find a basis and the dimension of (i) $S \cap T$, (ii) $S+T$.

Ans:- Hence, $S \cap T = \{(x_1, x_2, x_3, x_4) : \begin{matrix} 2x_1 + x_2 + 3x_3 + x_4 = 0 \\ x_1 + 2x_2 + x_3 + 3x_4 = 0 \end{matrix}\}$

Let, $\tilde{x} \in S \cap T$, then

$$2x_1 + x_2 + 3x_3 + x_4 = 0 \Rightarrow 2x_1 + x_2 = -3x_3 - x_4 \quad \text{--- (1)}$$

$$x_1 + 2x_2 + x_3 + 3x_4 = 0 \Rightarrow x_1 + 2x_2 = -x_3 - 3x_4 \quad \text{--- (2)}$$

$$\text{Solving (1) \& (2), we get, } x_1 = \frac{-5x_3 - x_4}{3}$$

$$x_2 = \frac{x_3 - 5x_4}{3}$$

$$\text{Hence, } \tilde{x} = \left(\frac{-5x_3 - x_4}{3}, \frac{x_3 - 5x_4}{3}, x_3, x_4 \right)$$

$$= x_3 \left(-\frac{5}{3}, \frac{1}{3}, 1, 0 \right) + x_4 \left(\frac{1}{3}, -\frac{5}{3}, 0, 1 \right)$$

Hence, $\left\{ \left(-\frac{5}{3}, \frac{1}{3}, 1, 0 \right), \left(\frac{1}{3}, -\frac{5}{3}, 0, 1 \right) \right\}$ forms a basis.

$$\therefore \dim(S \cap T) = 2.$$

Now, $\underline{x} \in S$

$$\begin{aligned}\therefore \underline{x} &= (x_1, x_2, x_3, x_4) ; 2x_1 + x_2 + 3x_3 + x_4 = 0 \\ &= (x_1, x_2, x_3, -2x_1 - x_2 - 3x_3) \\ &= x_1(1, 0, 0, -2) + x_2(0, 1, 0, -1) + x_3(0, 0, 1, -3)\end{aligned}$$

Clearly, $\{(1, 0, 0, -2), (0, 1, 0, -1), (0, 0, 1, -3)\}$ forms a basis for S .

$$\therefore \dim(S) = 3.$$

Now, $\underline{x} \in T$

$$\begin{aligned}\therefore \underline{x} &= (x_1, x_2, x_3, x_4) ; x_1 + 2x_2 + x_3 + 3x_4 = 0 \\ &= (-2x_2 - x_3 - 3x_4, x_2, x_3, x_4) \\ &= x_2(-2, 1, 0, 0) + x_3(-1, 0, 1, 0) + x_4(-3, 0, 0, 1)\end{aligned}$$

Clearly, $\{(-2, 1, 0, 0), (-1, 0, 1, 0), (-3, 0, 0, 1)\}$ forms a basis for T .

$$\therefore \dim(T) = 3.$$

$$\begin{aligned}\therefore \dim(S+T) &= \dim(S) + \dim(T) - \dim(S \cap T) \\ &= 3 + 3 - 2 = 4\end{aligned}$$

Clearly, $S+T \subseteq Y_4$

And $\dim(S+T) = \dim(Y_4)$

$$\Rightarrow S+T = Y_4$$

$\Rightarrow \{\underline{r}_1, \underline{r}_2, \underline{r}_3, \underline{r}_4\}$ is a basis of $(S+T)$.

6) If $U = \text{Span}\{(1, 2, 1), (2, 1, 3)\}$, $W = \text{Span}\{(1, 0, 0), (0, 1, 0)\}$
 Show that U and W are subspaces of \mathbb{R}^3 . Determine
 a basis and dimension of $U \cap W$ and $U + W$.

ANS:- Let $\underline{x} \in U \cap W$

$$\begin{aligned}\therefore \underline{x} &= l_1(1, 2, 1) + l_2(2, 1, 3) \\ &= l_1(1, 0, 0) + l_2(0, 1, 0)\end{aligned}$$

$$\therefore l_1 + 2l_2 = l_1$$

$$2l_1 + l_2 = l_2$$

$$l_1 + 3l_2 = 0$$

$$\therefore l_1 = -3l_2$$

$$\begin{aligned}\therefore \underline{x} &= l_2 \left\{ (-3)(1, 2, 1) + (2, 1, 3) \right\} = l_2(-1, -5, 0) \\ &= t(1, 5, 0), t \in \mathbb{R}\end{aligned}$$

$$\therefore U \cap W = \left\{ t(1, 5, 0) : t \in \mathbb{R} \right\}$$

$$\therefore \text{Basis of } U \cap W = \{(1, 5, 0)\}$$

$$\therefore \dim(U \cap W) = 1.$$

$(1, 2, 1), (2, 1, 3) \in U$ and they are linearly independent,
 so therefore they form a basis of U .

$$\therefore \dim(U) = 2$$

$$\text{Similarly, } \dim(W) = 2$$

$$\begin{aligned}\therefore \dim(U + W) &= \dim(U) + \dim(W) - \dim(U \cap W) \\ &= 2 + 2 - 1 = 3.\end{aligned}$$

Clearly, $U + W \subseteq \mathbb{R}^3$

$$\therefore \dim(U + W) = \dim(\mathbb{R}^3)$$

$$\therefore U + W = \mathbb{R}^3$$

$\therefore \{\underline{e}_1, \underline{e}_2, \underline{e}_3\}$ is a basis of $(U + W)$.

7) Given $S_1 = \{(1, 2, 3), (0, 1, 2), (3, 2, 1)\}$ and

$S_2 = \{(1, -2, 3), (-1, 1, -1), (1, -3, 4)\}$, Determine the
 dimension and a basis for

i) $[S_1] \cap [S_2]$, ii) $[S_1] + [S_2]$,

where $[S]$ denotes the span of S .

[C.U. 2004]

$$i) S_1 \cap S_2 = \left[\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 3 & 2 & 1 \\ 1 & -2 & 3 \\ -1 & 1 & -2 \\ 1 & -3 & 4 \end{array} \right] \sim \left[\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 5 & -5 \\ 0 & -1 & 1 \\ 0 & -2 & 2 \\ 0 & -5 & 1 \end{array} \right]$$

$R_3' = R_3 + 3R_5$
 $R_4' = R_4 + R_5$
 $R_5' = R_5 + R_6$
 $R_6' = R_6 - R_1$

$$\sim \left[\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \\ 0 & -5 & 1 \end{array} \right]$$

$R_3' = R_3/5$
 $R_5' = R_4 - 2R_3$

$$\sim \left[\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & -5 & 1 \end{array} \right]$$

$R_3' = R_3/3$
 $R_4 \leftrightarrow R_6$

$$\sim \left[\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{array} \right]$$

$R_2' = R_2 - 2R_3$
 $R_4' = R_4 + 5R_2 - R_3$

Hence, ~~$\dim(S_1 \cap S_2)$~~ $\dim(S_1 \cap S_2) = 3$

Now, $\text{dis}(S) = \dim(S_1) + \dim(S_2) - \dim(S_1 \cap S_2)$

$$S_1 = \left(\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 3 & 2 & 1 \end{array} \right) \xrightarrow{\substack{R_3 \rightarrow R_3 - 3R_1 \\ R_3' = R_3 + 1R_2}} \left(\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & -4 & -8 \end{array} \right)$$

$$\sim \left(\begin{array}{ccc} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 0 \end{array} \right)$$

$\therefore \dim(S_1) = 2$

$$S_2 = \left(\begin{array}{ccc} 1 & -2 & 3 \\ -1 & 1 & -2 \\ 1 & -3 & 4 \end{array} \right) \xrightarrow{\substack{R_2' \rightarrow R_1 + R_2 \\ R_3' \rightarrow R_3 - R_1}} \left(\begin{array}{ccc} 1 & -2 & 3 \\ 0 & -1 & 1 \\ 0 & -1 & 1 \end{array} \right)$$

$$\xrightarrow{\substack{R_3' \rightarrow R_3 - R_2 \\ R_2' \rightarrow R_2/4}} \left(\begin{array}{ccc} 1 & -2 & 3 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{array} \right)$$

$\dim(S_2) = 2$

$$\therefore \dim(S_1 + S_2) = \text{dis}(S) + \dim(S_2) - \dim(S_1 \cap S_2)$$

$$= 2 + 2 - 3$$

$$= 1.$$

8) Construct an orthogonal matrix $A^{4 \times 4}$ whose first row is $(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$. Find A^{-1} . The matrix B is obtained by replacing the third column of A by (2x3rd column of A). Using A^{-1} , find B^{-1} . [C.U. 2002]

Ans:- Note that $\{\underline{e}_1, \underline{e}_2, \underline{e}_3, \underline{e}_4\}$ form a basis for E^4 .

$$\text{Now, } (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}) = \frac{1}{2} \underline{e}_1 + \frac{1}{2} \underline{e}_2 + \frac{1}{2} \underline{e}_3 + \frac{1}{2} \underline{e}_4$$

$\therefore \{(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2}), \underline{e}_2, \underline{e}_3, \underline{e}_4\}$ forms a basis for E^4 .

$$\text{Taking, } \underline{v}_1 = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$$

$$\therefore \underline{u}_1 = \frac{\underline{v}_1}{\|\underline{v}_1\|} = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$$

$$\underline{v}_2 = (0, 1, 0, 0) - \{(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})(0, 1, 0, 0)\} \frac{1}{2}(1, 1, 1, 1)$$

$$= (0, 1, 0, 0) - \frac{1}{4}(1, 1, 1, 1)$$

$$= (-\frac{1}{4}, \frac{3}{4}, -\frac{1}{4}, -\frac{1}{4})$$

$$\therefore \underline{u}_2 = \frac{\underline{v}_2}{\|\underline{v}_2\|} = \frac{(-\frac{1}{4}, \frac{3}{4}, -\frac{1}{4}, -\frac{1}{4})}{\sqrt{\frac{1}{16} + \frac{9}{16} + \frac{1}{16} + \frac{1}{16}}} = -\frac{1}{\sqrt{12}}(1, -3, 1, 1)$$

$$\underline{v}_3 = \underline{e}_3 - \{\underline{e}_3 \cdot \underline{u}_1\}\underline{u}_1 - \{\underline{e}_3 \cdot \underline{u}_2\}\underline{u}_2$$

$$= (0, 0, 1, 0) - \{(0, 0, 1, 0)(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})\}(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$$

$$- \{(0, 0, 1, 0)(-\frac{1}{4}, \frac{3}{4}, -\frac{1}{4}, -\frac{1}{4})\}(-\frac{1}{\sqrt{12}}, \frac{3}{\sqrt{12}}, -\frac{1}{\sqrt{12}}, \frac{1}{\sqrt{12}}) \}(-\frac{1}{\sqrt{12}}, \frac{3}{\sqrt{12}}, \frac{1}{\sqrt{12}}, \frac{1}{\sqrt{12}})$$

$$= (-\frac{1}{3}, 0, \frac{2}{3}, -\frac{1}{3})$$

$$\therefore \underline{u}_3 = \frac{\underline{v}_3}{\|\underline{v}_3\|} = \frac{(-\frac{1}{3}, 0, \frac{2}{3}, -\frac{1}{3})}{\sqrt{\frac{1}{9} + 0 + \frac{4}{9} + \frac{1}{9}}} = \frac{1}{\sqrt{6}}(-1, 0, 2, -1).$$

$$\underline{v}_4 = (0, 0, 0, 1) - \{(0, 0, 0, 1)(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})\}(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})$$

$$- \{(0, 0, 0, 1)(-\frac{1}{\sqrt{12}}, 1, -3, 1, 1)\}(-\frac{1}{\sqrt{12}})(1, -3, 1, 1)$$

$$- \{(0, 0, 0, 1)(\frac{1}{\sqrt{6}})(-1, 0, 2, -1)\}(\frac{1}{\sqrt{6}})(-1, 0, 2, -1)$$

$$= (-\frac{1}{2}, 0, 0, \frac{1}{2})$$

$$\therefore \underline{u}_4 = \frac{\underline{v}_4}{\|\underline{v}_4\|} = \frac{(-\frac{1}{2}, 0, 0, \frac{1}{2})}{\sqrt{\frac{1}{4} + 0 + 0 + \frac{1}{4}}} = \frac{1}{\sqrt{2}}(-1, 0, 0, 1)$$

Hence, $\{u_1, u_2, u_3, u_4\}$ is a set of 4 orthonormal vectors.
 Hence, $A = \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix}$ is an orthonormal matrix with first

$$\text{row } u_1 = \left(\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2} \right)^T.$$

$$A = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ -\frac{1}{\sqrt{2}} & \frac{3}{\sqrt{2}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{6}} & 0 & \frac{2}{\sqrt{6}} & -\frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{2}} & 0 & 0 & \frac{1}{\sqrt{2}} \end{bmatrix}$$

As A is an orthogonal matrix, so $A' = A^{-1}$.

$$\therefore A^{-1} = A' = \begin{bmatrix} \frac{1}{2} & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{2} & \frac{3}{\sqrt{2}} & 0 & 0 \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{2}{\sqrt{6}} & 0 \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{2}} \end{bmatrix}$$

$$\text{Hence, } B = AE_3(2), \text{ where } E_3(2) = \begin{pmatrix} 1 & & & \\ 0 & 0 & 0 & \\ 0 & 0 & 0 & \\ 0 & 0 & 0 & -1 \end{pmatrix}$$

$$\text{Now, } B^{-1} = [A E_3(2)]^{-1}$$

$$= E_3^{-1}(2) A^{-1}$$

$$= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \frac{1}{2} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{2} & \frac{3}{\sqrt{2}} & 0 & 0 \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{2}{\sqrt{6}} & 0 \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & 0 \end{pmatrix}$$

$$= \begin{pmatrix} \frac{1}{2} & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} & -\frac{1}{\sqrt{2}} \\ \frac{1}{2} & \frac{3}{\sqrt{2}} & 0 & 0 \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{2}{\sqrt{6}} & 0 \\ \frac{1}{2} & -\frac{1}{\sqrt{2}} & \frac{1}{\sqrt{6}} & 0 \end{pmatrix}$$

Q) Let A and B be defined by

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ -1 & 1 & 1 \end{pmatrix}, \quad B = \begin{pmatrix} 0 & 2 & 1 \\ 0 & 4 & 2 \\ 0 & -2 & -1 \end{pmatrix}$$

(i) Show that the column space of B is a subspace of A.

(ii) Find a matrix C such that $AC = B$. [C.U. 2000]

Ans:— (i) Column space of A = $C(A) = \left\{ \lambda_1 \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} + \lambda_2 \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix} + \lambda_3 \begin{pmatrix} 1 \\ 2 \\ -3 \end{pmatrix}; \lambda_1, \lambda_2, \lambda_3 \in \mathbb{R} \right\}$

Column space of B = $C(B) = \left\{ \lambda \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix}; \lambda \in \mathbb{R} \right\} \quad \lambda_1, \lambda_2, \lambda_3 \in \mathbb{R}$

Now, for any vector $x \in C(B)$

$$\text{we have, } x = \lambda \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} = \lambda_1 \begin{pmatrix} 1 \\ 2 \\ -1 \end{pmatrix} + \lambda_2 \begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix} + \lambda_3 \begin{pmatrix} 1 \\ 2 \\ -3 \end{pmatrix}, \text{ where } \lambda_1 = \lambda, \lambda_2 = \lambda_3 = 0.$$

which belongs to $C(A)$.

$$\therefore x \in C(A) \Rightarrow x \in C(B).$$

$$\therefore C(B) \subseteq C(A)$$

(ii) Hence, note that, $r(A)=2$, i.e. $r(A) < 3$,
since 1st & 2nd rows are LD.

$\Rightarrow A$ is singular.

$\Rightarrow A^{-1}$ does not exist.

[If A^{-1} exists, then $AC = B \Rightarrow C = A^{-1}B$ and it is unique]

Hence, choice of C is not unique.

Note that,

$$(A) \begin{pmatrix} 1 & 1 & 1 \\ 2 & 2 & 2 \\ -1 & 1 & 1 \end{pmatrix} (C) \begin{pmatrix} 0 & 2 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} = (B) \begin{pmatrix} 0 & 2 & 1 \\ 0 & 4 & 2 \\ 0 & -2 & -1 \end{pmatrix}$$

Alternative choice of C = $\begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{pmatrix}$.

- 10) Let $A = \begin{pmatrix} 1 & 2 & 3 \\ 3 & 5 & 7 \\ 9 & 16 & 23 \end{pmatrix}$. Suppose $\text{r}(A) = r$.
 i) Find r . ii) Write A as the sum of n matrices each of rank unity. iii) find an orthonormal basis of the row-space of A .
 [C.U. 2008]

ANS:-

Note that, $(9, 16, 23) = 3(1, 2, 3) + 2(3, 5, 7)$

i) $A = \begin{pmatrix} 1 & 2 & 3 \\ 3 & 5 & 7 \\ 9 & 16 & 23 \end{pmatrix}_{3 \times 1 + 2 \times 1 \quad 3 \times 2 + 2 \times 5 \quad 3 \times 3 + 2 \times 7}$ have two LN rows.
 clearly, $\text{rank}(A) = 2$, as there

ii) $A = \begin{pmatrix} 1 & 2 & 3 \\ 0 & 0 & 0 \\ 3 & 5 & 7 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 \\ 3 & 5 & 7 \\ 9 & 16 & 23 \end{pmatrix}$
 $= \begin{pmatrix} 1 & 2 & 3 \\ 0 & 0 & 0 \\ 3 & 5 & 7 \end{pmatrix} + \begin{pmatrix} 0 & 0 & 0 \\ 3 & 5 & 7 \\ 3 \times 2 & 5 \times 2 & 7 \times 2 \end{pmatrix}$
 $= A_1 + A_2$, where $\text{Rank}(A_i) = 1, \forall i = 1, 2$.

iii) $\{(1, 2, 3), (3, 5, 7)\}$ form a basis of the row space of A .

$$\therefore v_1 = (1, 2, 3)$$

$$\therefore u_1 = \frac{v_1}{\|v_1\|} = \frac{(1, 2, 3)}{\sqrt{1+4+9}} = \frac{1}{\sqrt{14}} (1, 2, 3)$$

$$\therefore v_2 = (3, 5, 7) - \left\{ (3, 5, 7) \frac{1}{\sqrt{14}} (1, 2, 3) \right\} \frac{1}{\sqrt{14}} (1, 2, 3)
= \left(\frac{4}{7}, \frac{1}{7}, -\frac{2}{7} \right)$$

$$\therefore u_2 = \frac{v_2}{\|v_2\|} = \frac{7}{\sqrt{21}} \left(\frac{4}{7}, \frac{1}{7}, -\frac{2}{7} \right)
= \frac{1}{\sqrt{21}} (4, 1, -2)$$

1) Suppose $A = \begin{pmatrix} 7 & -12 & 4 \\ 3 & -2 & 5 \\ -6 & 11 & 8 \end{pmatrix}$.

(a) Find A^{-1} .

(b) Find a matrix B such that $AB = \begin{pmatrix} 8 & -12 & 4 \\ 3 & -1 & 5 \\ -6 & 11 & 9 \end{pmatrix}$.

(c) Write A as a sum of symmetric and a skew-symmetric matrices.

(d) Let V be a vector space generated by the first two column vectors of A . Then write the third column vectors of A as the sum of the two non-zero vectors such that one is a member of V and the other is orthogonal to V .

[C.U. 1998]

ANS:-

(d) $V = \left\{ l_1 \begin{pmatrix} 7 \\ 3 \\ -6 \end{pmatrix} + l_2 \begin{pmatrix} -12 \\ -2 \\ 11 \end{pmatrix} : l_1, l_2 \in \mathbb{R} \right\}$.

Let, $\begin{pmatrix} 4 \\ 5 \\ 8 \end{pmatrix} = \underline{x} + \underline{y}$ where $\underline{x} \in V$ and \underline{y} is \perp to V .

$$\text{Now, } \underline{x} \in V, \underline{x} = l_1 \begin{pmatrix} 7 \\ 3 \\ -6 \end{pmatrix} + l_2 \begin{pmatrix} -12 \\ -2 \\ 11 \end{pmatrix}$$

$$\text{Hence, } \underline{y} \text{ is } \perp \text{ to } V \Rightarrow 7y_1 + 3y_2 - 6y_3 = 0 \\ -12y_1 - 2y_2 + 11y_3 = 0$$

$$\text{Let, } y_3 = t, \text{ now, } y_1 = \frac{21t}{22}, y_2 = -\frac{5t}{22}; \underline{y} = t \begin{pmatrix} 21/22 \\ -5/22 \\ 1 \end{pmatrix} \\ \therefore \underline{x} = l_1 \begin{pmatrix} 7 \\ 3 \\ -6 \end{pmatrix} + l_2 \begin{pmatrix} -12 \\ -2 \\ 11 \end{pmatrix}; l_1, l_2 \in \mathbb{R} = 22t \begin{pmatrix} 21 \\ -5 \\ 22 \end{pmatrix} \\ = l_3 \begin{pmatrix} 21 \\ -5 \\ 22 \end{pmatrix} \\ \therefore \begin{pmatrix} 4 \\ 5 \\ 8 \end{pmatrix} = l_1 \begin{pmatrix} 7 \\ 3 \\ -6 \end{pmatrix} + l_2 \begin{pmatrix} -12 \\ -2 \\ 11 \end{pmatrix} + l_3 \begin{pmatrix} 21 \\ -5 \\ 22 \end{pmatrix}; \\ l_1, l_2, l_3 \in \mathbb{R}.$$

12) (a) Suppose $A = \begin{pmatrix} 1 & 3 & 6 \\ 5 & 9 & 8 \\ 2 & 7 & 1 \end{pmatrix}$, Find two non-singular matrices P and Q such that $PAQ = I_3$. [C.U.1999]

(b) Find inverse of $A = \begin{pmatrix} 1 & -1 & 0 \\ 2 & 3 & -1 \\ -1 & 2 & 0 \end{pmatrix}$, by Pivotal Condensation or sweeping-out method.

(c) Obtain the fully reduced normal form of $A = \begin{pmatrix} 1 & 0 & 1 \\ 2 & 1 & 3 \\ 1 & 1 & 2 \end{pmatrix}$

(d) Express $A = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 2 \end{pmatrix}$ as the product of elementary matrices.

ANS:-

(i) By problem, $PAQ = I_3$

$\therefore \text{Rank}(PAQ) = \text{Rank}(I_3)$

$\therefore \text{rank}(A) = 3$, as P and Q are non-singular matrices and rank is not altered by non-singular matrix multiplication.

$\Rightarrow A$ is non-singular.

$\Rightarrow A^{-1}$ exists.

$\therefore A^{-1}AI_3 = I_3$

$\Rightarrow PAQ = I_3$, where $P = A^{-1}$, $Q = I_3$.

$$(ii). (A | I_3) = \left(\begin{array}{ccc|ccc} 1 & -1 & 0 & 1 & 0 & 0 \\ 2 & 3 & -1 & 0 & 1 & 0 \\ -1 & 2 & 0 & 0 & 0 & 1 \end{array} \right)$$

R

[Row operations]

~~~~~

~~~~~

~~~~~

$$\sim \left( \begin{array}{ccc|ccc} 1 & 0 & 0 & 2 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 7 & -1 & 5 \end{array} \right)$$

$$\therefore A^{-1} = \begin{pmatrix} 2 & 0 & 1 \\ 1 & 0 & 1 \\ 7 & -1 & 5 \end{pmatrix}$$

$$(c) [I_3 \mid A^{3 \times 3} \mid I_3]$$

$$= \left[ \begin{array}{ccc|ccc|c} 1 & 0 & 0 & 1 & 0 & 1 & & I_3 \\ 0 & 1 & 0 & 2 & 1 & 3 & & \\ 0 & 0 & 1 & 1 & 1 & 2 & & \end{array} \right]$$

$$\begin{array}{l} R'_2 \rightarrow R_2 - 2R_1 \\ R'_3 \rightarrow R_3 - R_1 \end{array} \quad \left[ \begin{array}{ccc|ccc|c} 1 & 0 & 0 & 1 & 0 & 1 & & I_3 \\ -2 & 1 & 0 & 0 & 1 & 1 & & \\ -1 & 0 & 1 & 0 & 1 & 1 & & \end{array} \right]$$

$$\begin{array}{l} R'_3 \rightarrow R_3 - R_2 \end{array} \quad \left[ \begin{array}{ccc|ccc|c} 1 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ -2 & 1 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right]$$

$$\begin{array}{l} C'_3 \rightarrow C_3 - C_1 - C_2 \end{array} \quad \left[ \begin{array}{ccc|ccc|c} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & -1 \\ -2 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & -1 \\ 1 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right]$$

$$= \left[ \begin{array}{c|cc|c} P & I_2 & 0 & Q \\ \hline 0 & 0 & 0 & \end{array} \right]$$

$$\text{Hence, } P = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ 1 & -1 & 1 \end{pmatrix}, Q = \begin{pmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \\ 0 & 0 & 1 \end{pmatrix}$$

$$(d) [I_3 \mid A^{3 \times 3} \mid I_3]$$

$$= \left[ \begin{array}{ccc|ccc|c} 1 & 0 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 & 2 & 0 & 0 & 1 \end{array} \right]$$

~

$$\sim \sim \sim \left[ \begin{array}{ccc|ccc|c} 1 & 0 & 0 & 1 & 0 & 0 & 2 & 0 & -1 \\ -1 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 \\ 1 & -1 & 1 & 0 & 0 & 1 & -1 & 0 & 1 \end{array} \right]; P = \begin{pmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 1 & -1 & 1 \end{pmatrix}, Q = \begin{pmatrix} 2 & 0 & -1 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix}$$

$$\therefore A = P^{-1} I_3 Q^{-1} = P^{-1} Q^{-1}$$

$$\left. \begin{array}{l} \left( \begin{array}{ccc} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 1 & -1 & 1 \end{array} \right) = \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) P \\ \Rightarrow \left( \begin{array}{ccc} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) = \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{array} \right) P \quad [R'_3 \rightarrow R_3 + R_2] \\ \Rightarrow \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) = \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{array} \right) P \quad [R'_2 \rightarrow R_2 + R_1] \\ \therefore P^{-1} = \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \end{array} \right) \end{array} \right| \quad \left. \begin{array}{l} \left( \begin{array}{ccc} 2 & 0 & -1 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{array} \right) = \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) Q \\ \Rightarrow \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ -1 & 0 & 1 \end{array} \right) = \left( \begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) Q \\ \Rightarrow \left( \begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right) = \left( \begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 2 \end{array} \right) Q \\ \Rightarrow Q^{-1} = \left( \begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 2 \end{array} \right) \end{array} \right|$$

12) (a) Solve the following system of linear equations using a numerical method:

$$x_1 - 2x_2 + 3x_3 + 4x_4 = 4.5$$

$$3x_1 - x_2 + 2x_3 + 5x_4 = 9.5$$

$$2x_1 + 4x_2 + 5x_3 + x_4 = 15$$

$$4x_1 + 2x_2 - x_3 + 3x_4 = 12$$

[C.U. 2007]

ANS:-

$$\left( \begin{array}{cccc|c} 1 & -2 & 3 & 4 & 4.5 \\ 3 & -1 & 2 & 5 & 9.5 \\ 2 & 4 & 5 & 1 & 15 \\ 4 & 2 & -1 & 3 & 12 \end{array} \right)$$

$$R_2' = R_2 - 3R_1$$

$$\tilde{R}_3' = R_3 - 2R_1$$

$$R_4' = R_4 - 4R_1$$

$$\left( \begin{array}{cccc|c} 1 & -2 & 3 & 4 & 4.5 \\ 0 & 5 & -7 & -7 & 9.5 \\ 0 & 0 & 1 & -1 & 15 \\ 0 & 0 & 0 & 0 & 12 \end{array} \right)$$

$$\sim \left( \begin{array}{cccc|c} 1 & -2 & 3 & 4 & 4.5 \\ 0 & 5 & -7 & -7 & 9.5 \\ 0 & 0 & 1 & -1 & 15 \\ 0 & 0 & 0 & 0 & 12 \end{array} \right)$$

$$\therefore \left( \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array} \right) = \left( \begin{array}{c} 1.5 \\ 1.5 \\ 15 \\ 12 \end{array} \right)$$

$$\therefore x_1$$

$$x_2$$

$$x_3$$

$$x_4$$

(b) For what value of  $k$  the planes

$$x - 4y + 5z = k$$

$$x - y + 2z = 3 \quad (\text{i}) \text{ intersect in a line?}$$

$$2x + y + z = 0, (\text{ii}) \text{ intersect in a point?}$$

Ans:-

Augmented matrix =  $[A : b]$

$$= \left[ \begin{array}{ccc|c} 1 & -4 & 5 & k \\ 1 & -1 & 2 & 3 \\ 2 & 1 & 1 & 0 \end{array} \right]$$

$$R_2' \rightarrow R_2 - R_1$$

$$\sim \left[ \begin{array}{ccc|c} 1 & -4 & 5 & k \\ 0 & 3 & -3 & 3-k \\ 0 & 9 & -9 & -2k \end{array} \right]$$

$$R_3' \rightarrow R_3 - 3R_2$$

$$\sim \left[ \begin{array}{ccc|c} 1 & -4 & 5 & k \\ 0 & 1 & -1 & 3-k/3 \\ 0 & 0 & 0 & k-9 \end{array} \right]$$

$$R_2' \rightarrow R_2/3$$

$$\sim \left[ \begin{array}{ccc|c} 1 & -4 & 5 & k \\ 0 & 1 & -1 & 3-k/3 \\ 0 & 0 & 0 & k-9 \end{array} \right]$$

(i) If  $k=9$ , then  $\text{Rank}(A : b) = 2 = \text{Rank}(A)$

$\Rightarrow$  the system is consistent and has infinitely many solutions and intersect in a line.

(ii)  $\text{R}(A : b) = \text{R}(A) = 3$  is not satisfied by no value of  $k$ .

13) (i) Identify the definiteness of the following quadratic forms:

$$(i) x_1^2 + 2x_2^2 + 3x_3^2 + 2x_1x_2 + 4x_2x_3 + 2x_1x_3 .$$

$$(ii) 2x_1^2 + 2x_2^2 + 5x_3^2 - 4x_1x_2 - 2x_1x_3 + 2x_2x_3 ,$$

Ans:- (i)  $A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 2 & 3 \end{pmatrix}; a_{11} = 1 > 0, \begin{vmatrix} 1 & 1 \\ 1 & 2 \end{vmatrix} = 1 > 0,$

$$\begin{vmatrix} 1 & 1 & 1 \\ 1 & 2 & 2 \\ 1 & 2 & 3 \end{vmatrix} = 1 > 0$$

$\therefore A$  is a p.d. matrix.

(ii)  $A = \begin{pmatrix} 2 & -2 & -1 \\ -2 & 2 & 1 \\ -1 & 1 & 5 \end{pmatrix}$

$$a_{11} = 2 > 0, \begin{vmatrix} 2 & -2 \\ -2 & 2 \end{vmatrix} = 0, \begin{vmatrix} 2 & -2 & -1 \\ -2 & 2 & 1 \\ -1 & 1 & 5 \end{vmatrix} = 0$$

$\therefore A$  is a p.s.d. matrix.

14) S.I.T. A =  $\begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 7 \\ 3 & 7 & 11 \end{pmatrix}$  is p.d. and find the matrix B such that  $A = BB^T$ .

Ans:-  $a_{11} = 1 > 0, \begin{vmatrix} 1 & 2 \\ 2 & 5 \end{vmatrix} = 1 > 0, \begin{vmatrix} 1 & 2 & 3 \\ 2 & 5 & 7 \\ 3 & 7 & 11 \end{vmatrix} = 170$

∴ All the principal minors of A is positive.  
∴ A is p.d. matrix.

$$\left( \begin{array}{ccc|ccc|ccc} 1 & 0 & 0 & 1 & 2 & 3 & 1 & 0 & 0 \\ 0 & 1 & 0 & 2 & 5 & 7 & 0 & 1 & 0 \\ 0 & 0 & 1 & 3 & 7 & 11 & 0 & 0 & 1 \end{array} \right)$$

$$R_2' \rightarrow R_2 - 2R_1 \sim C_2' \rightarrow C_2 - 2C_1 \left( \begin{array}{ccc|ccc|ccc} 1 & 0 & 0 & 1 & 0 & 3 & 1 & -2 & 0 \\ -2 & 1 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 3 & 1 & 11 & 0 & 0 & 1 \end{array} \right)$$

$$R_3' \rightarrow R_3 - 3R_1 \sim C_3' \rightarrow C_3 - 3C_1 \left( \begin{array}{ccc|ccc|ccc} 1 & 0 & 0 & 1 & 0 & 0 & 1 & -2 & -3 \\ -2 & 1 & 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ -3 & 0 & 1 & 0 & 1 & 2 & 0 & 0 & 1 \end{array} \right)$$

$$R_3' \rightarrow R_3 - R_2 \sim C_3' \rightarrow C_3 - C_2 \left( \begin{array}{ccc|ccc|ccc} 1 & 0 & 0 & 1 & 0 & 0 & 1 & -2 & -1 \\ -2 & 1 & 0 & 0 & 1 & 0 & 0 & 1 & -1 \\ -1 & -1 & 1 & 0 & 0 & 1 & 0 & 0 & 1 \end{array} \right)$$

$$\therefore CAc' = I_3 \Rightarrow A = c^{-1} I_3 (c^{-1})' \\ = (c^{-1})(c^{-1})' \\ = BB^T$$

where  $B = c^{-1} = \begin{pmatrix} 1 & 0 & 0 \\ -2 & 1 & 0 \\ -1 & -1 & 1 \end{pmatrix}^{-1}$   
 $= \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 1 & 1 \end{pmatrix}$

(ii) Solve the equations:

$$5x_1 - 2x_2 - x_3 = 42$$

$$5x_2 - 2x_3 - x_1 + 3x_2 = 0$$

$$5x_3 - 2x_1 - x_2 = 32$$

To express  $x_i$ 's in terms of  $z$ , for what value of  $z$ , do  $x_1, x_2, x_3$  further satisfy  $x_1 + x_2 + x_3 = 12$

[C.U. 1991]

Ans:- Augmented matrix is  $[A|b]$

$$\begin{pmatrix} 5 & -2 & -1 & | & 42 \\ -1 & 5 & -2 & | & -32 \\ -2 & -1 & 5 & | & 32 \end{pmatrix}$$

$$\begin{matrix} R_2 \leftrightarrow R_1 \\ \sim \end{matrix} \begin{pmatrix} -1 & 5 & -2 & | & -32 \\ 5 & -2 & -1 & | & 42 \\ -2 & -1 & 5 & | & 32 \end{pmatrix}$$

$$\begin{matrix} R_2' \leftrightarrow R_2 + 5R_1 \\ R_3' \rightarrow R_3 - 2R_1 \end{matrix} \begin{pmatrix} -1 & 5 & -2 & | & -32 \\ 0 & 23 & -11 & | & -112 \\ 0 & -11 & 9 & | & 92 \end{pmatrix}$$

$$\begin{matrix} R_1' \rightarrow R_1(-1) \\ R_2' \rightarrow R_2/23 \\ R_3' \rightarrow R_3/11 \end{matrix} \begin{pmatrix} 1 & -5 & 2 & | & 32 \\ 0 & 1 & -11/23 & | & -112/23 \\ 0 & 1 & -9/11 & | & -92/11 \end{pmatrix}$$

$$\begin{matrix} R_3' \rightarrow R_3 - R_2 \end{matrix} \begin{pmatrix} 1 & -5 & 2 & | & 32 \\ 0 & 1 & -11/23 & | & -11/23 \\ 0 & 0 & -86/253 & | & -\frac{86}{253}z \end{pmatrix}$$

$$\begin{matrix} R_3' = R_3 / -\frac{86}{253} \end{matrix} \begin{pmatrix} 1 & -5 & 2 & | & 32 \\ 0 & 1 & -11/23 & | & -11/23 \\ 0 & 0 & 1 & | & z \end{pmatrix}$$

$$\therefore x_3 = t, x_2 = -\frac{11}{23}t,$$

$$\Rightarrow x_2 = 0$$

$$\Rightarrow x_1 = t$$

$\therefore (x_1, x_2, x_3) = t(1, 0, 1)$  is a solution.

$$\text{Again } x_1 + x_2 + x_3 = 1 \Rightarrow t + 0 + t = 1 \Rightarrow t = 1/2$$

15) i) Determine the null space of  $A = \begin{bmatrix} 1 & 1 & -1 & 2 \\ 2 & 2 & -3 & 1 \\ -1 & -1 & 0 & -5 \end{bmatrix}$ .

Also find  $\dim(N(A))$  and  $\text{rank}(A)$ .

Ans:- To find  $\dim(N(A))$ , we use,  
 $\dim(N(A)) = 4 - \text{rank}(A)$ .

We know,  $\{x : Ax = 0\} \subset N(A)$

Hence,

$$\begin{pmatrix} 1 & 1 & -1 & 2 \\ 2 & 2 & -3 & 1 \\ -1 & -1 & 0 & -5 \end{pmatrix} \sim \begin{pmatrix} 1 & 1 & -1 & 2 \\ 0 & 0 & -1 & -3 \\ 0 & 0 & -1 & -3 \end{pmatrix} \sim \begin{pmatrix} 1 & 1 & -1 & 2 \\ 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 0 \end{pmatrix} \sim \begin{pmatrix} 1 & 1 & 0 & 5 \\ 0 & 0 & 1 & 3 \\ 0 & 0 & 0 & 0 \end{pmatrix} = 1$$

Rank(A) = 2.

$$Hx = 0$$

$$\Rightarrow x_1 + x_2 + 5x_4 = 0$$

$$\Rightarrow x_3 + 3x_4 = 0 \Rightarrow x_3 = -3x_4$$

$$\therefore x_4 = t, x_3 = -3t, x_2 = -3t, x_1 = -2t$$

$$[N(A)] = t \begin{pmatrix} 1 \\ -3 \\ -3 \\ -2 \end{pmatrix}$$

$$\therefore \dim(N(A)) = 4 - 2 = 2.$$

(ii) Find an orthogonal matrix which diagonalizes  
 $A = \begin{pmatrix} 6 & -2 & 2 \\ -2 & 3 & -1 \\ 2 & -1 & 3 \end{pmatrix}$ . Also, find  $A^8$ .

ANS:- The characteristic equation of  $A$  is

$$\begin{aligned} 0 &= |A - \lambda I_3| = \begin{vmatrix} 6-\lambda & -2 & 2 \\ -2 & 3-\lambda & -1 \\ 2 & -1 & 3-\lambda \end{vmatrix} \\ &= \begin{vmatrix} 6-\lambda & -2 & 2 \\ 0 & 2-\lambda & 2-\lambda \\ 2 & -1 & 3-\lambda \end{vmatrix} \quad R_2' \rightarrow R_2 + R_3 \\ &= (2-\lambda) \begin{vmatrix} 6-\lambda & -2 & 2 \\ 0 & 1 & 1 \\ 2 & -1 & 3-\lambda \end{vmatrix} \\ &= (2-\lambda) \left\{ (6-\lambda)(3-\lambda+1) + 2(-2-2) \right\} \\ &= (2-\lambda)(\lambda-2)(\lambda-8) \end{aligned}$$

$\therefore \lambda = 2, 2, 8$  are the eigen values of  $A$ .

For  $\lambda = 8$ ,

$$(A - 8I_3)\vec{x} = 0$$

$$\Rightarrow \begin{pmatrix} -2 & -2 & 2 \\ -2 & -5 & -1 \\ 2 & -1 & -5 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\Rightarrow \begin{cases} x_1 + x_2 - x_3 = 0 \\ 2x_1 + 5x_2 + x_3 = 0 \\ 2x_1 - x_2 - 5x_3 = 0 \end{cases}$$

$$\text{Let, } x_3 = t, \Rightarrow x_1 + x_2 = t.$$

$$\Rightarrow x_2 = t - x_1$$

$$\therefore x_2 = -t, \& x_1 = 2t.$$

$$\therefore \vec{x} = t \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix}, \quad x_3 = \frac{\begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix}}{\sqrt{4+1+1}} = \frac{1}{\sqrt{6}} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix}$$

$$\text{For } \lambda = 2, (A - 2I_3)\vec{x} = 0$$

$$\Rightarrow \begin{pmatrix} 4 & -2 & 2 \\ -2 & 1 & -1 \\ 2 & -1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

$$\Rightarrow 2x_1 - x_2 + x_3 = 0$$

$$x_2 = t_1, \quad x_3 = t_2, \quad x_1 = \frac{t_1 - t_2}{2}, \quad \vec{x} = t_1 \begin{pmatrix} \frac{1}{2} \\ 1 \\ 0 \end{pmatrix} + t_2 \begin{pmatrix} 1/2 \\ 0 \\ 1 \end{pmatrix}$$

$\vec{x}_1 = \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}, \quad \vec{x}_2 = \begin{pmatrix} -1 \\ 0 \\ 2 \end{pmatrix}$  are L.I.N. eigen vectors when  $\lambda = 2$ .

$$\tilde{u}_1 = u_1 = \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}$$

$$u_1 = \frac{\tilde{u}_1}{\|\tilde{u}_1\|} = \frac{\begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}}{\sqrt{5}} = \frac{1}{\sqrt{5}} \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix}$$

$$\begin{aligned} \tilde{u}_2 &= u_2 - (u_1 \cdot u_2) u_1 \\ &= \begin{pmatrix} -4/5 \\ 2/5 \\ 2 \end{pmatrix} \end{aligned}$$

$$u_2 = \frac{\tilde{u}_2}{\|\tilde{u}_2\|} = \frac{1}{\sqrt{30}} \begin{pmatrix} -2 \\ 1 \\ 5 \end{pmatrix}$$

$$\tilde{u}_3 = u_3 - (u_3 \cdot u_1) u_1 - (u_3 \cdot u_2) u_2$$

$$\therefore u_3 = \frac{\tilde{u}_3}{\|\tilde{u}_3\|} = \frac{1}{\sqrt{6}} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix} = \frac{1}{\sqrt{6}} \begin{pmatrix} 2 \\ -1 \\ 1 \end{pmatrix}$$

$$\begin{aligned} \therefore Q &= (u_1 \ u_2 \ u_3) \\ &= \left( \begin{array}{ccc} \frac{1}{\sqrt{5}} & -\frac{2}{\sqrt{30}} & \frac{2}{\sqrt{6}} \\ \frac{2}{\sqrt{5}} & \frac{1}{\sqrt{30}} & \frac{1}{\sqrt{6}} \\ 0 & \frac{5}{\sqrt{30}} & \frac{1}{\sqrt{6}} \end{array} \right) \text{ is orthogonal matrix.} \end{aligned}$$

$$\therefore Q^T A Q = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 8 \end{pmatrix}$$

$\lambda_i^m$  is the eigen value of the matrix  $A^m$ .

$$\text{So, } Q^T A^8 Q = \begin{pmatrix} 2^8 & 0 & 0 \\ 0 & 2^8 & 0 \\ 0 & 0 & 8^8 \end{pmatrix}$$

$$\Rightarrow Q^T A^8 = Q' \begin{pmatrix} 256 & 0 & 0 \\ 0 & 256 & 0 \\ 0 & 0 & 16777216 \end{pmatrix} Q$$

$$\text{as } Q^{-1} = Q'$$

## PRACTICAL PROBLEMS ON DESCRIPTIVE STATISTICS

Q.No.1. [CU'03] The sound intensity levels measured in decibels at 50 construction sites are

68, 63, 59, 77, 60, 57, 63, 62, 64, 73, 63, 70, 71, 65, 68, 67, 67, 62, 56, 61, 69, 64, 58, 73, 68, 66, 65, 64, 68, 67, 69, 68, 67, 70, 69, 61, 65, 62, 68, 62, 67, 64, 82, 86, 65, 69, 68, 70, 66.

- (i) Prepare a frequency table for considering 6 classes of equal width, (ii) Calculate the exact mean and variance of sample from the raw data. (iii) Estimate the mean and variance from the grouped frequency distribution by assuming each value is equal to the mid-point of the class to which it belongs.
- (iv) Compare the answers in (ii) and (iii) and in this connection mention about the hole of Sheppard's correction.

Q.No.2. [CU'96] A Chemical compound contains 12.5% of iron was given to two technicians A and B for chemical analysis. A made of 15 determinations, and B 10 determinations, of the percentage of iron. Their results are given in the following table

| Determinations by A |       |       | Determinations by B |       |
|---------------------|-------|-------|---------------------|-------|
| 12.46               | 12.11 | 12.47 | 12.23               | 12.11 |
| 11.80               | 12.44 | 12.11 | 11.91               | 12.45 |
| 12.40               | 11.85 | 12.56 | 12.45               | 12.39 |
| 11.95               | 12.12 | 12.65 | 12.22               | 12.37 |
| 12.77               | 12.43 | 12.72 | 12.05               | 12.65 |

- (i) Find separately for A and B various measures of central tendency and dispersion. Also, find their respective coefficients of variation.
- (ii) Based on the above measures prepare a report comparing the accuracy and consistency of the two technicians.

Q.No.3. The following tables relate to the weights of new born babies recorded at the different clinics. One of the clinics is located in a locality where the average family income is also homogeneous throughout the locality where as the other clinic caters to the same category but also has a considerable number of citizens from a significant lower income group. Examine the locations, dispersions and the shapes of the two distributions and interpret your findings.

| Weights in (kgs.) | Clinic I | Clinic II |
|-------------------|----------|-----------|
| 0 - 1             | 0        | 0         |
| 2 - 3             | 1        | 43        |
| 3 - 4             | 60       | 230       |
| 4 - 5             | 304      | 372       |
| 5 - 6             | 318      | 320       |
| 6 - 7             | 56       | 54        |
| > 7               | 3        | 2         |
|                   | 0        | 0         |

[CU'1995]

Q. No. 4. [CU'2001] The scores in English of 250 candidates appearing in an examination have

$$m_1 = 39.7213, \delta = 9.8894$$

$$\frac{m_3}{m^{3/2}} = -0.1182, b_2 = 2.9719.$$

It is later found onscrutining that score 61 has been wrongly recorded as 51. Obtain the correct values of  $b_1$  and  $b_2$ .

Q. No. 5. Particulars relating to the monthly wage distributions of two manufacturing firms are given below:

| <u>Measures</u> | <u>Firm A</u>  | <u>Firm B</u> |
|-----------------|----------------|---------------|
| Mean wage       | Rs. 1477       | 1495          |
| Median wage     | Rs. 1389       | 1354          |
| Modal wage      | Rs. 1350       | 1312          |
| Quartiles       | Rs. 1278, 1422 | 1262, 1435    |
| S.D.            | Rs. 87         | Rs. 99        |

Compare the two distributions w.r.t. the characteristic central tendency, dispersion and skewness, kurtosis.

Q. No. 6. The weights (in grams) of 25 indicator housing used on gauges are as follows:

|       |       |       |       |         |
|-------|-------|-------|-------|---------|
| 102.0 | 106.3 | 106.6 | 108.8 | 109.7   |
| 106.1 | 105.9 | 106.7 | 106.8 | 110.2   |
| 101.7 | 106.6 | 106.3 | 110.2 | 109.9   |
| 102.0 | 105.8 | 106.1 | 106.7 | 108.7.3 |
| 102.0 | 106.8 | 110.0 | 107.9 | 109.3   |

- (a) Construct an ordered stem-leaf display using integers as the stems and tenths as the leaves.
- (b) Find the five-number summary of the data and draw a box-plot.
- (c) Are there any suspected outliers?

Q.No.7. The following table gives the yearly corn yield ( $X$ ) in bushels per acre, in six Corn Belt states (Iowa, Illinois, Nebraska, Missouri, Indiana, Ohio) and rainfall measurements in inches ( $Y$ ) in six states from 1915 to 1927.

| Year: | 1915 | 1916 | 1917 | 1918 | 1919 | 1920 | 1921 | 1922 | 1923 | 1924 | 1925 | 1926 | 1927 |
|-------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| $X :$ | 33.3 | 29.7 | 35.0 | 29.9 | 35.2 | 38.3 | 35.2 | 35.5 | 36.7 | 26.8 | 38.0 | 31.7 | 32.6 |
| $Y :$ | 16.5 | 9.3  | 9.4  | 9.0  | 9.5  | 11.6 | 12.1 | 8.0  | 10.7 | 13.9 | 11.3 | 11.6 | 10.9 |

Fit a linear regression model:  $X = \alpha + \beta Y + \epsilon$  to the data by the method of least squares, making the usual assumptions.

Q.No.8. [CU'1997] In an experiment on certain fertilizers were applied at various levels (in appropriate units) with resulting yields (in appropriate unit) as follows:

|                           |      |      |      |      |      |      |      |
|---------------------------|------|------|------|------|------|------|------|
| Fertilizer level ( $x$ ): | 0    | 5    | 10   | 15   | 20   | 25   | 30   |
| Yield ( $y$ ):            | 27.1 | 32.1 | 35.0 | 36.2 | 36.9 | 36.1 | 35.2 |

- (i) Fit an appropriate polynomial to the given data.
- (ii) Obtain a suitable measure of association between  $x$  and  $y$  and comment
- (iii) Obtain the optimum fertilizer level which maximizes yield.

Q.No.9. [CU'2001] The following data relate to the height ( $x$ ) and weight ( $y$ ) of 15 students

|     |      |      |      |      |      |      |      |      |      |      |      |      |      |      |      |
|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| $x$ | 5'6" | 5'8" | 5'8" | 5'9" | 5'6" | 5'6" | 5'9" | 5'3" | 5'6" | 5'6" | 5'9" | 5'6" | 5'3" | 5'9" | 5'3" |
| $y$ | 55   | 52   | 60   | 61   | 57   | 59   | 60   | 51   | 58   | 50   | 59   | 55   | 44   | 65   | 47   |

Compute  $r_{yz}$  and  $e_{yz}$ . Is the regression linear? If not, compute a measure of deviation from linearity.

Q.No.10. Two supervisors ranked 12 workers working under them in order of efficiency as follows. Compute Spearman's rank correlation coefficient between the two rankings. Also compute Kendall's  $\gamma$ .

| Worker :        | A | B | C | D | E | F | G | H | I | J  | K  | L  |
|-----------------|---|---|---|---|---|---|---|---|---|----|----|----|
| Supervisor I :  | 5 | 6 | 1 | 2 | 3 | 8 | 8 | 4 | 7 | 11 | 10 | 12 |
| Supervisor II : | 5 | 5 | 1 | 1 | 1 | 9 | 7 | 4 | 8 | 10 | 12 | 10 |

Q.No.11. [CU'2008]

The following table gives data on income in thousand dollars( $x$ ), the number of families ( $N$ ) at income  $x$  and the number of families owing a house ( $n$ ).

|       |    |    |     |    |    |    |    |    |
|-------|----|----|-----|----|----|----|----|----|
| $x$ : | 10 | 13 | 15  | 20 | 25 | 30 | 35 | 40 |
| $N$ : | 60 | 80 | 100 | 70 | 65 | 50 | 40 | 25 |
| $n$ : | 18 | 28 | 45  | 36 | 39 | 33 | 30 | 20 |

Suggest an appropriate regression equation to explain the effect of income on owning a house. Also estimate the parameters of this equation using the above data and predict there from the proportion of families of income 32 thousands dollars who own a house.

Solutions: PROBLEMS ON DESCRIPTIVE STATS.

1. Minimum value = 56, Maximum value = 86

$$\text{Range} = 30.$$

(i) Considering 6 classes of width 5:

| <del>class-limit</del><br>Class-boundary | Frequency<br>( $f_i$ ) | Tally marks | cls. mark<br>( $x_i'$ ) |
|------------------------------------------|------------------------|-------------|-------------------------|
| 56 - 60.5                                | 2                      |             | 55.5                    |
| 60.5 - 64.5                              | 16                     |             | 61.5                    |
| 64.5 - 68.5                              | 25                     |             | 67.5                    |
| 68.5 - 72.5                              | 4                      |             | 70.5                    |
| 72.5 - 76.5                              | 2                      |             | 74.5                    |
| 76.5 - 80.5                              | 1                      | /           | 79.5                    |
| <b>TOTAL</b>                             | <b>50</b>              |             |                         |

(ii) From the raw data,

$$\text{mean}(\bar{x}) = \frac{1}{50} \sum_{i=1}^{50} x_i$$

$$\text{Variance}(s^2) = \frac{1}{50} \sum_{i=1}^{50} x_i^2 - \bar{x}^2$$

(iii) from the grouped data,

$$\text{mean}(\bar{x}') = \frac{1}{\sum f_i} \sum_{i=1}^6 x'_i f_i$$

$$\text{Variance} (s'^2) = \frac{1}{50} \sum_{i=1}^6 x'_i f_i - \bar{x}'^2$$

(iv) The mean and variance computed from the raw data and from grouped data are different.

In general, frequency distribution, we assume that all the values in a class are equal to the mid-point which is true if values are uniformly distributed over the class. So mean and variance computed from the raw data and grouped data are different due to the error due to grouping. The correction for the errors due to grouping is given by Sheppard's correction.

$$m_1' (\text{corrected}) = m_1'$$

$$m_2' (\text{corrected}) = m_2 - \frac{c^4}{12}, c = \text{class width}$$

$$2. (i) \text{ Exact mean of } A = \frac{1}{15} \sum_{i=1}^{15} x_{iA} =$$

$$\text{Exact mean of } B = \frac{1}{10} \sum_{i=1}^{10} x_{iB} =$$

$$\text{Measure of Dispersion, } SD(A) = \sqrt{\frac{1}{15} \sum_{i=1}^{15} x_{iA}^2 - \bar{x}_A^2}$$

$$SD(B) = \sqrt{\frac{1}{10} \sum_{i=1}^{10} x_{iB}^2 - \bar{x}_B^2} =$$

$$RMSD_A (12.5) = \sqrt{\frac{1}{15} \sum_{i=1}^{15} (x_{iA} - 12.5)^2}$$

$$RMSD_B (12.5) = \sqrt{\frac{1}{10} \sum_{i=1}^{10} (x_{iB} - 12.5)^2}$$

$$\text{Co-efficient of variation, } CV(A) = \frac{SD(A)}{\bar{x}_A} \times 100\% =$$

$$CV(B) = \frac{SD(B)}{\bar{x}_B} \times 100\% =$$

(ii) As a measure of accuracy, we use RMSD about 12.5.  
The smaller the RMSD, is more accuracy.

As a measure of consistency, we use C.V. The smaller C.V. is more consistency in Data set.

3) The frequency distribution of clinic-I is highly positively skewed and that of clinic-II is near about symmetric.  
[To get the measure of distribution we may draw histogram of the frequency distribution].

i) As a measure of location, we may use median.

$\therefore$  Median of Clinic I is

$\therefore$  Median of Clinic II is

(ii) For Dispersion, we use  $Q.D. = \frac{Q_3 - Q_1}{2}$ .

$\therefore$  Q.D. of Clinic I is

& Q.D. of Clinic II is

(iii) For shape, we use,  $S_K = \frac{Q_3 + Q_1 - 2Q_2}{Q_3 - Q_1}$

$\therefore$  Skewness of Clinic I is

$\therefore$  Skewness of Clinic II is

iv) For peakedness, we use,  $K_p = \frac{Q_3 - Q_1}{2(P_{90} - P_{10})}$

$\therefore$  Kurtosis of Clinic I is

$\therefore$  Kurtosis of Clinic II is

The smaller the value of  $K_p$ , the higher the kurtosis.

| Weights in<br>(Kgs) | Clinic I | Clinic II | $\leq I$ | $\leq II$ | $P_{90}^I = 4 + \frac{667.8 - 365}{318} \times 1$ |
|---------------------|----------|-----------|----------|-----------|---------------------------------------------------|
| 0-1                 | 0        | 0         | 0        | 0         | $= 4.95$                                          |
| 1-2                 | 1        | 43        | 1        | 43        | $P_{10}^I = 3 + \frac{74.2 - 61}{304} \times 1$   |
| 2-3                 | 60       | 230       | 61       | 278       | $= 3.04$                                          |
| 3-4                 | 304      | 372       | 365      | 645       | $Q_1^I = 3 + \frac{185.5 - 61}{304} \times 1$     |
| 4-5                 | 318      | 320       | 683      | 965       | $= 3.41$                                          |
| 5-6                 | 56       | 54        | 739      | 1019      | $Q_2^I = 4 + \frac{371 - 365}{318} \times 1$      |
| 6-7                 | 8        | 2         | 742      | 1021      | $= 4.02$                                          |

$$Q_3^I = 4 + \frac{556.5 - 365}{318} \times 1 \\ = 4.60$$

$$n = 250$$

$$m_1' (\text{incorrected}) = 39.7213$$

$$\Rightarrow \frac{1}{250} \sum_{i=1}^{50} x_i (\text{incorrected}) = 39.7213$$

$$\Rightarrow \sum_{i=1}^{50} x_i (\text{incorrected}) = 9930.325$$

$$\therefore \sum_i x_i (\text{connected}) = (9930.325 - 51 + 61) \\ = 9940.325$$

$$\therefore m_1' (\text{connected}) = 39.7613.$$

$$\delta (\text{incorrected}) = \sqrt{\frac{1}{250} \sum x_i^2 (\text{incorrect}) - \bar{x}^2 (\text{incorrect})} \\ = 9.8894$$

$$\therefore \sum x_i^2 (\text{incorrected}) = 418895.4765$$

$$\therefore \sum x_i^2 (\text{connected}) = 418895.4765 - 51^2 + 61^2 \\ = 420015.4765$$

$$\therefore \delta (\text{connected}) = \sqrt{\frac{1}{250} \sum x_i^2 (\text{connected}) - \bar{x}^2 (\text{connected})} \\ = 9.9549$$

$$b_1 (\text{incorrected}) = \frac{m_3 (\text{incorrected})}{m_2^{3/2} (\text{incorrected})} = -0.1182$$

$$\Rightarrow m_3 (\text{incorrected}) = -114.32134.$$

5) For Firm A: Mean > Median > Mode and the distribution is +ve skewed.

For Firm B: Mean > Median > Mode & the distn. is +ve skewed.

(For skewed distn. measures should be based on quantiles)

Measures:-

(i) Location: Median ( $Q_2$ )

(a) Median of Firm A is 1389

(b) " " " B is 1354

(ii) Dispersion: Q.D. =  $\frac{Q_3 - Q_1}{2}$ .

(a) Q.D. of Firm A is  $\frac{1422 - 1278}{2} = 72$

(b) Q.D. of Firm B is  $\frac{1435 - 1262}{2} = 86.5$

(iii) Skewness:  $S_K = \frac{Q_3 + Q_1 - 2Q_2}{Q_3 - Q_1}$

(a)  $S_K$  of firm A is  $(1422 + 1278 - 2 \times 1350) / 144 = 0$

(b)  $S_K$  of firm B is  $(1262 + 1435 - 2 \times 1312) / 173 = 0.42$

(iv) Kurtosis:  $K_p = \frac{Q_3 - Q_1}{2 \times S.D.}$

(a)  $K_p$  of firm A is  $(1422 - 1278) / 2 \times 87 = 0.8276$

(b)  $K_p$  of firm B is  $(1435 - 1262) / 2 \times 99 = 0.8739$

[ $(Q_3 - Q_1)$  represents the length within which we have central 50% value. The smaller the length  $(Q_3 - Q_1)/2$ , the higher the kurtosis. To set a unit free measure, we divide  $\frac{Q_3 - Q_1}{2}$  by S.D.]

6) (a) Consider the integer as stems and decimals as leaves.

| Stem | Leaf        |
|------|-------------|
| 101  | 7           |
| 102  | 0 0 0       |
| 105  | 8 9         |
| 106  | 1 3 3       |
| 107  | 6 6 7 7 8 8 |
| 107  | 3 7 9       |
| 108  | 8           |
| 109  | 1 3 9       |
| 110  | 0 2 2       |

Comment: The frequency distribution is located at 106 (median or mode) & is +ve skewed.

(b) Minimum value = 101.7, maximum value = 110.2

$$Q_1 = \frac{25}{4}^{\text{th}} \text{ ordered value} = 6^{\text{th}} \text{ value} + \frac{1}{4}(7^{\text{th}} - 6^{\text{th}})$$

$$= 105.9 + \frac{1}{4}(0.2) = 105.95 \text{ gram}$$

$$Q_2 = \frac{25}{2}^{\text{th}} \text{ ordered value} = 12^{\text{th}} \text{ ordered value}$$

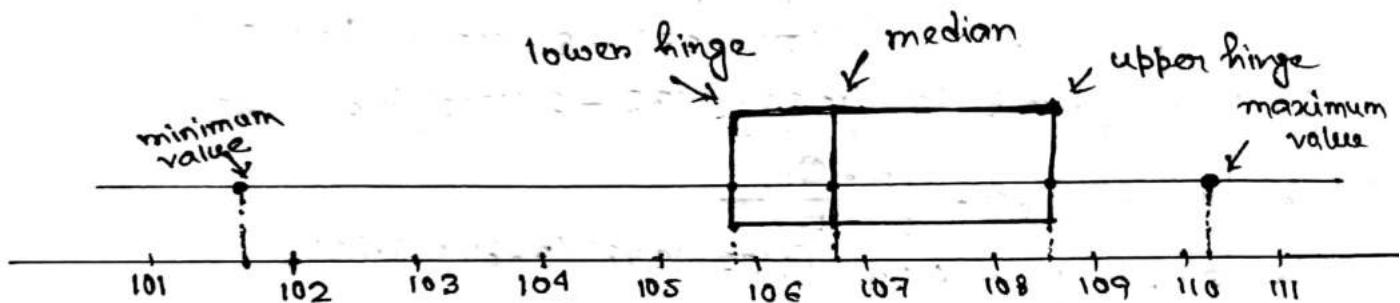
$$+ \frac{1}{2}(13^{\text{th}} - 12^{\text{th}}) \text{ ordered value}$$

$$= 108.7 + \frac{1}{2} \times 0 = 108.7 \text{ gm.}$$

$$Q_3 = \frac{25 \times 3}{4}^{\text{th}} \text{ ordered value} = 18^{\text{th}} \text{ ordered value}$$

$$+ \frac{3}{4}(19^{\text{th}} - 18^{\text{th}}) \text{ ordered value.}$$

$$= 107.9 + \frac{3}{4}(0.9) = 108.575 \text{ gm.}$$



Comment: — The frequency distribution is located at (near about) 107. The H-shaped spread is small. Hence, the dispersion of the data is not high. The distance between upper hinge and median is greater than that of lower hinge and median. The distribution is very skewed. The length of H-spread is small w.r.t. the range, the kurtosis is high.

(c) The data values that are smaller than  $L_H - \frac{3}{2}H$  or greater than  $L_U + \frac{3}{2}H$  are called outliers.

$$L_H - \frac{3}{2}H = 105.95 - \frac{3}{2}(108.575 - 105.95)$$

$$= 102.0125$$

$$L_U + \frac{3}{2}H = 108.575 + \frac{3}{2}(108.575 - 105.95)$$

$$= 112.5125$$

In the data, 101.7, 102.0, 102.0, 102.0 are the outliers of the data set.

### 7) Linear Regression model:

$$x = \alpha + \beta y + \epsilon$$

The constant  $\alpha, \beta$  are determined by minimising error sum of squares or residual sum of squares.

$$S = \sum_{i=1}^n (x_i - \alpha - \beta y_i)^2$$

The normal equations are:

$$\sum_{i=1}^n x_i = n\alpha + \beta \sum_{i=1}^n y_i$$

$$\sum_{i=1}^n x_i y_i = \alpha \sum_{i=1}^n y_i + \beta \sum_{i=1}^n y_i^2$$

From data,  $n=13$ ,  $\sum x_i = 437.9$ ,  $\sum y_i = 143.8$ ,

$$\sum x_i y_i = 4419.35, \sum y_i^2 = 1543.67,$$

The normal equations are:

$$437.9 = \beta \alpha + (\beta \times 143.8) \quad \text{--- ①}$$

$$4419.35 = 143.8 \alpha + 1543.67 \beta \quad \text{--- ②}$$

$$\text{①} \times 143.8 - \text{②} \times 13, \text{ we get}$$

$$\beta = 11.34.$$

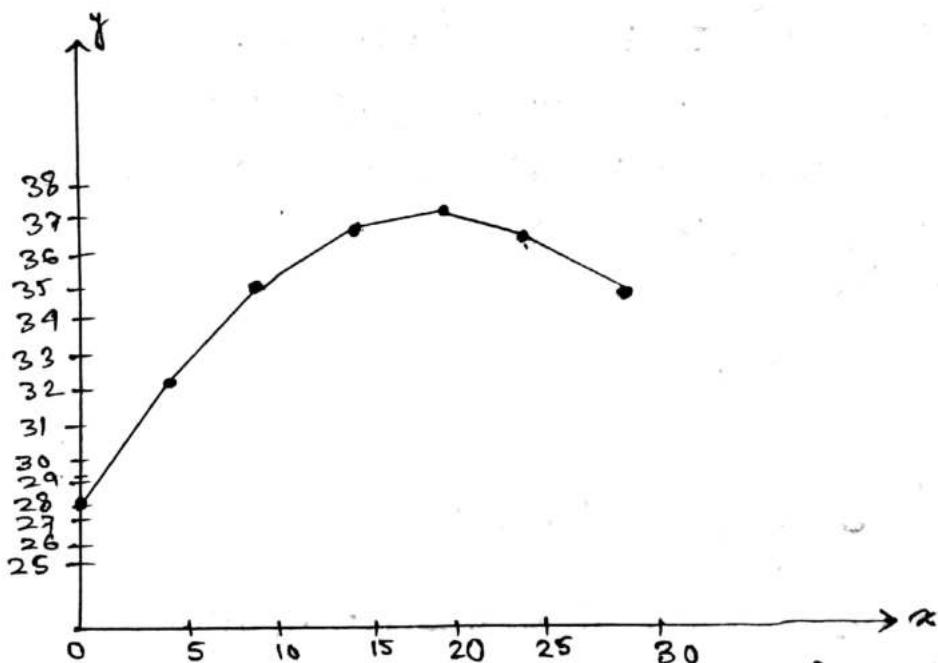
$$\therefore \alpha = -91.36$$

$$\therefore \hat{\alpha} = -91.36, \hat{\beta} = 11.34$$

Hence, the predicting formula is —

$$x = -91.36 + 11.34 y.$$

8)



From scatter plot, it is clear that the relationship between  $x$  &  $y$  is approximately a quadratic.

(a) A measure of association between  $x$  and  $y$  is the measure of usefulness of the 2nd degree polynomial (least square) regression as a predicting formula, i.e.

$$R^2 = \frac{\sum (Y_{pi} - \bar{Y}_p)^2}{\sum (y_i - \bar{y})^2} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (x_i - \bar{x})^2}$$

| $x$ | $y$  | $u_i = \frac{x_i - 15}{5}$ | $u_i^2$ | $u_i^3$ | $u_i^4$ | $u_i y_i$ | $u_i^2 y_i$ | $\hat{y}_i$ |
|-----|------|----------------------------|---------|---------|---------|-----------|-------------|-------------|
| 0   | 27.1 | -3                         |         |         |         |           |             |             |
| 5   | 32.1 | -2                         |         |         |         |           |             |             |
| 10  | 35.0 | -1                         |         |         |         |           |             |             |
| 15  | 36.2 | 0                          |         |         |         |           |             |             |
| 20  | 36.9 | 1                          |         |         |         |           |             |             |
| 25  | 36.1 | 2                          |         |         |         |           |             |             |
| 30  | 35.2 | 3                          |         |         |         |           |             |             |

$$\text{Let, } Y = a_0 + a_1 u + a_2 u^2$$

the constants  $a_0, a_1, a_2$  are determined by minimizing

$$\sum_{i=1}^n (y_i - a_0 - a_1 u_i - a_2 u_i^2)$$

∴ Normal equations are:-

$$\sum y_i = n a_0 + a_1 \sum u_i + a_2 \sum u_i^2$$

$$\sum u_i y_i = a_0 \sum u_i + a_1 \sum u_i^2 + a_2 \sum u_i^3$$

$$\sum u_i^2 y_i = a_0 \sum u_i^2 + a_1 \sum u_i^3 + a_2 \sum u_i^4$$

Here,  $n=7$ ,  $\sum u_i = 0$ ,  $\sum u_i^2 = 28$ , etc.

Now, the total variability  $= \sum y_i^2 - n \bar{y}^2$

& unexplained variability (on Residual sum of squares (RSS))

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \sum (y_i - \hat{y}_i) e_i^2$$

$$= \sum y_i^2 - \sum \hat{y}_i^2$$

$$= \sum y_i^2 - (y_i - a_0 - a_1 u_i - a_2 u_i^2)^2$$

$$= \sum y_i^2 - a_0^2 \sum y_i^2 - a_1^2 \sum u_i y_i - a_2^2 \sum u_i^2 y_i$$

$$\therefore R^2 = 1 - \frac{\text{Unexplained variability}}{\text{Total variability}}$$

(iii)  $y = \hat{a}_0 + \hat{a}_1 u + \hat{a}_2 u^2, \quad u = \frac{x-15}{5}$

$$\frac{dy}{du} = \hat{a}_1 + 2\hat{a}_2 u$$

$$\Rightarrow u = -\frac{\hat{a}_1}{2\hat{a}_2}, \text{ since } \frac{dy}{dx} = 0$$

$$\& \frac{d^2y}{du^2} = 2\hat{a}_2 < 0$$

Hence,  $y$  is maximum when  $u_0 = -\frac{\hat{a}_1}{2\hat{a}_2}$   
 $\Rightarrow x = 15 + 5u_0$

9)

| $x_i$                   | $y_{ij}$               | $\bar{y}_i$ | $n_i$ | $n_i \bar{y}_i$ |
|-------------------------|------------------------|-------------|-------|-----------------|
| $x_1 = 5'3''$<br>= 63'' | 52, 60, 61, 44, 47     | 50.8        | 5     | 129.08          |
| $x_2 = 5'6''$<br>= 66'' | 55, 57, 59, 58, 50, 55 | 54.8        | 6     | 180.18          |
| $x_3 = 5'9''$<br>= 69'' | 61, 60, 58, 65         | 61          | 4     | 148.84          |

Now,  $\tilde{e}_{yx} = \frac{\sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2}{\sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y})^2} = \frac{\sum_{i=1}^k n_i \bar{y}_i^2 - n \bar{y}^2}{\sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 - n \bar{y}^2}$ , where

$$\bar{y} = \frac{\sum_{i=1}^k n_i \bar{y}_i}{\sum_{i=1}^k n_i}, \quad \bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}, \quad i=1(1)k$$

$$\bar{x} = \frac{1}{\sum_{i=1}^k n_i} \sum_{i=1}^k n_i x_i, \quad \bar{x} = 65.8, \quad \bar{y} = \frac{\sum_{i=1}^k \bar{y}_i}{n} = 55.13$$

$$\therefore \hat{e}_{yx} = \frac{45805 - 45590}{46039 - 45590} = 0.4788$$

$$\therefore e_{yx} = 0.692$$

$$\hat{s}_{yx} = \frac{\sum_{i=1}^n n_i x_i \bar{y}_i - n \bar{x} \bar{y}}{\sqrt{\sum_i n_i x_i - n \bar{x}^2} \sqrt{\sum_j y_{ij} - n \bar{y}^2}}$$

$$= \frac{54538.8 - 54413.31}{8.97 \times 21.20}$$

$$= 0.1658$$

$$\therefore s_{yx} = 0.2585$$

the regression equation is linear iff  $\hat{e}_{yx} = \hat{s}_{yx}$ . A measure of deviation of the (true) regression equation from linearity is  $d_{yx} = \hat{e}_{yx} - \hat{s}_{yx} = 0.4788 - 0.1658 = 0.313$ .

10) Both series of efficiency have some tie. So, we replace the tied rank by their average of those.

| Workers | A   | B   | C | D | E | F   | G   | H | I | J    | K  | L    |
|---------|-----|-----|---|---|---|-----|-----|---|---|------|----|------|
| I       | 5   | 6   | 1 | 2 | 3 | 8.5 | 8.5 | 4 | 7 | 11   | 10 | 12   |
| II      | 5.5 | 5.5 | 2 | 2 | 2 | 9   | 7   | 4 | 8 | 10.5 | 12 | 10.5 |

$$\frac{n-1}{12} - \frac{T_u + T_v}{2} - \frac{1}{2n} \sum d_{ij}$$

$$s_R = \sqrt{\frac{\frac{n-1}{12} - T_u}{\frac{n-1}{12} - T_v}}$$

11)

~~Explanatory variable ( $x$ ): income..~~  
~~Response variable ( $y$ ): owning a house.~~

Hence  $y$  is Binary.

Hence, a logistic regression may be appropriate.

In logistic regression,

$$\hat{y}_x = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

$$\Rightarrow \ln\left(\frac{\hat{y}_x}{1 - \hat{y}_x}\right) = \alpha + \beta x.$$

$$\Rightarrow u_x = \alpha + \beta x$$

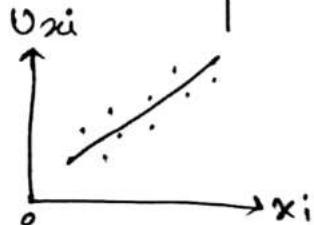
$\Rightarrow x$  &  $u_x$  are linearly related,

| $x$        | $y$ -values             | Array total<br>$= \sum_{j=1}^{N_i} y_{ij}$ | $\bar{y}_{xi} = \frac{y_i}{N_i}$ | $u_{xi}$ |
|------------|-------------------------|--------------------------------------------|----------------------------------|----------|
| $x_1 = 10$ | $y_{11} \dots y_{1N_1}$ |                                            |                                  |          |
| $x_2 = 13$ | $y_{21} \dots y_{2N_2}$ |                                            |                                  |          |
| $x_3 = 15$ | $\dots$                 |                                            |                                  |          |
| $\vdots$   |                         |                                            |                                  |          |
| $x_8 = 40$ | $y_{81} \dots y_{8N_8}$ |                                            |                                  |          |

Plot the points  $(x_i, u_{xi})$ ,  $i=1(1)8$ ,

From graph the points

$(x_i, u_{xi})$  lie near a straight line.



Hence logistic regression is appropriate.  
Hence,  $\alpha$  &  $\beta$  can be determined by minimising  $\sum w_i (u_{xi} - (\alpha + \beta x_i))^2$

$$\text{where } w_i = N_i \bar{y}_{xi} (1 - \bar{y}_{xi})$$

$$\text{Normal eqn. } \sum w_i u_{xi} = \alpha \sum w_i + \beta \sum x_i w_i$$

$$\sum w_i x_i u_{xi} = \alpha \sum x_i w_i + \beta \sum x_i^2 w_i$$

| $U_{xi}$ | $w_i$ | $x_i w_i$ | $x_i \check{w}_i$ | $w_i U_{xi}$ | $w_i x_i U_{xi}$ |
|----------|-------|-----------|-------------------|--------------|------------------|
|          |       |           |                   |              |                  |

fitted logistic regression is

$$Y = \frac{\alpha + \beta x}{1 + \gamma \hat{\alpha} + \hat{\beta} x}$$

## PROBLEMS ON CATEGORICAL DATA ANALYSIS

1. In one phase of a study regarding the effectiveness of several drugs on fast-operating nausea 167 patients were assigned at random, 30 to Drug-P, 67 to Drug-C and the remaining 70 to Placebo. The no. of patients suffering from severe, moderate, slight or no nausea is shown below:

|         | SEVERE | MODERATE | SLIGHT | NO NAUSEA | TOTAL |
|---------|--------|----------|--------|-----------|-------|
| PLACEBO | 8      | 8        | 19     | 35        | 70    |
| DRUG-P  | 2      | 3        | 5      | 20        | 30    |
| DRUG-C  | 3      | 4        | 15     | 45        | 67    |
| TOTAL   | 13     | 15       | 39     | 100       | 167   |

calculate a suitable measure of association.

2. A drug, supposed to have some effect in curing diabetes was treated on 100 patients in a certain hospital and their records were compared with 100 other patients not treated with drug. Study the efficacy of the drug in curing diabetes.

|           | Cured | Not cured |
|-----------|-------|-----------|
| Treated   | 54    | 46        |
| Untreated | 24    | 76        |

3. The following table gives the data on the results of a visual test and a Balance test performed on 413 male college students.

|              | Left eyed | Ambiocular | Right-eyed | TOTAL |
|--------------|-----------|------------|------------|-------|
| Left handed  | 48        | 25         | 52         | 125   |
| Ambidextrous | 32        | 13         | 25         | 70    |
| Right handed | 94        | 33         | 91         | 218   |
| TOTAL        | 174       | 71         | 168        | 413   |

Compute at least two measures of association and comment.

4. A study was conducted to find out what people of different educational level feel about the role of the caste in life.

| <del>Edu. level</del>    | <del>Role</del> | No Role | little Role | Large Role | Complete Role |
|--------------------------|-----------------|---------|-------------|------------|---------------|
| Upto 8th<br>(Grp-I)      |                 | 2       | 8           | 25         | 36            |
| CIS-9 - H.S.<br>(Grp-II) |                 | 10      | 15          | 14         | 12            |
| Graduate<br>(Grp-III)    |                 | 29      | 24          | 2          | 1             |

- (a) Compute a suitable measure to find if there is any association between educational level and the individual's perception towards the role of region or caste in social life. Interpret your finding.
- (b) Merge the columns 'No role' and 'little role' and rename it as 'minor role', similarly merge, the last two columns 'Large role' and 'complete Role' and rename it as 'Major Role'. Find the odds ratio of
- i) Group I with Group II
  - ii) Group III const Group I, const. the Minor and Major roles and interpret your findings.

## PROBLEMS ON CATEGORICAL DATA ANALYSIS

Solutions:-

$$\begin{aligned}
 1. \text{ We now calculate } & \sum_i \sum_j \frac{(n_{ij})^2}{n_{inj}} \\
 & = \frac{8^2}{13 \times 70} + \frac{2^2}{13 \times 30} + \frac{3^2}{13 \times 67} + \frac{8^2}{70 \times 15} + \frac{3^2}{15 \times 30} \\
 & \quad + \frac{4^2}{15 \times 67} + \frac{19^2}{39 \times 70} + \frac{5^2}{39 \times 30} + \frac{15^2}{39 \times 67} \\
 & \quad + \frac{35^2}{100 \times 70} + \frac{20^2}{100 \times 30} + \frac{45^2}{100 \times 67} \\
 & = 1.0381
 \end{aligned}$$

Mean-square contingency:-

$$\begin{aligned}
 \chi^2 &= N \left\{ \sum_{i=1}^8 \sum_{j=1}^4 \frac{(n_{ij})^2}{n_{inj}} - 1 \right\} \\
 &= 167 (1.0381 - 1) \\
 &= 6.3827
 \end{aligned}$$

Karl Pearson's coefficient of mean-square contingency :-

$$C = \sqrt{\frac{\chi^2}{N + \chi^2}} = 0.1916$$

Tschuprow's coefficient:-

$$\begin{aligned}
 T^2 &= \frac{\chi^2}{N \sqrt{(s-1)(t-1)}} , s=4, t=3 \\
 &= 0.1247
 \end{aligned}$$

2.

|                        | Cured (B)         | Not cured ( $\beta$ ) | Total |
|------------------------|-------------------|-----------------------|-------|
| Treated (A)            | $(AB) = 54$       | $(A\beta) = 46$       | 100   |
| Untreated ( $\alpha$ ) | $(\alpha B) = 24$ | $(\alpha\beta) = 76$  | 100   |

A measure of association between A and B in a  $2 \times 2$  table is:

(i) Yule's coefficient of association:-

$$Q = \frac{(AB)(\alpha\beta) - (A\beta)(\alpha B)}{(AB)(\alpha\beta) + (A\beta)(\alpha B)}$$

$$= 0.5760$$

(ii) Yule's coefficient of colligation:-

$$Y = \frac{\sqrt{(AB)(\alpha\beta)} - \sqrt{(A\beta)(\alpha B)}}{\sqrt{(AB)(\alpha\beta)} + \sqrt{(A\beta)(\alpha B)}}$$

$$= 0.3169$$

Hence, A and B have moderate positive association.  
i.e. Drug has moderate effect on Diabetes.

3. Here both the characters 'Balance Test' (A) and 'Visual Test' (B) are in normal scale.

|       |  | B <sub>1</sub>       | B <sub>2</sub>        | B <sub>3</sub>       | TOTAL                 |                       |
|-------|--|----------------------|-----------------------|----------------------|-----------------------|-----------------------|
|       |  | f <sub>11</sub> = 48 | f <sub>12</sub> = 25  | f <sub>13</sub> = 52 | f <sub>10</sub> = 125 |                       |
|       |  | A <sub>1</sub>       |                       |                      |                       |                       |
|       |  | A <sub>2</sub>       | f <sub>21</sub> = 32  | f <sub>22</sub> = 13 | f <sub>23</sub> = 25  | f <sub>20</sub> = 70  |
|       |  | A <sub>3</sub>       | f <sub>31</sub> = 94  | f <sub>32</sub> = 33 | f <sub>33</sub> = 91  | f <sub>30</sub> = 218 |
| TOTAL |  |                      | f <sub>01</sub> = 174 | f <sub>02</sub> = 71 | f <sub>03</sub> = 168 | 413                   |

Here, we shall use measures based on  $\chi^2$ .

$$\chi^2 = N \left\{ \sum_{i=1}^{n_s} \sum_{j=1}^{s_j} \frac{f_{ij}^2}{f_{0j} f_{i0}} - 1 \right\}$$

$$= 2.3738$$

(i) Karl Pearson's coefficient :-  $C_{AB} = \sqrt{\frac{\chi^2}{N + \chi^2}}$

$$= \sqrt{\frac{2.3738}{413 + 2.3738}}$$

$$= 0.0756 .$$

(ii) Tschuprow's coefficient :-  $T_{AB} = \sqrt{\frac{\chi^2}{N \sqrt{(n-1)(s-1)}}}$

$$= 0.05361 .$$

4. (a) Here both the characters  $\rightarrow$  educational level ( $X$ ) and Role of religion ( $Y$ ) are in ordinal scale.  
 We shall use measures based on the no. of concordant pairs ( $C$ ) and no. of discordant pairs ( $D$ ).

A pair  $(i, j)$  of individuals with scores  $(X_i, Y_i)$  and  $(X_j, Y_j)$  is in concordance if

$$\{X_i > X_j, Y_i > Y_j\} \text{ or } \{X_i < X_j, Y_i < Y_j\}$$

and discordant if

$$\{X_i > X_j, Y_i < Y_j\} \text{ or } \{X_i < X_j, Y_i > Y_j\}.$$

a pair  $(i, j)$  is on tie const.  $X$  if  $X_i = X_j$ .

Here,  $C = \text{No. of concordant pairs}$

$$= 2(15+14+12+24+2+1) + 8(14+12+2+1) + 25(2+1) \\ + 36(0) + 10(24+12+1) + 15(2+1) + 14 \times 1 \\ = 1022.$$

$D = \text{No. of discordant pairs}$

$$= 8(10+29) + 25(10+15+29+24) + 36(10+15+14+29+24+2) + 15 \times 29 + 14(29+24) \\ + 12(29+24+2) \\ = 7483.$$

Goodman-Kruskal ( $\gamma$ ) measure of association:-

$$\gamma = \frac{C - D}{C + D} = -0.7597,$$

So, we can say educational level and the individual's perception towards the role of religion is negatively associated.

|         | Minor Role    | Major Role    |
|---------|---------------|---------------|
| Gro I   | $f_{11} = 10$ | $f_{12} = 61$ |
| Gro II  | $f_{21} = 25$ | $f_{22} = 26$ |
| Gro III | $f_{31} = 53$ | $f_{32} = 3$  |

(i) Sample odds ratio of Gro.I with Gro.II =  $\frac{f_{11} f_{22}}{f_{21} f_{12}}$

Hence the success in 'minor role' in Gro.I is less likely than Gro.II, i.e., Hence, more people feel that religion has 'Major Role' from Gro.II compare to Gro.I.

(ii) Sample odds ratio of Gro.III w.r.t. Gro.I is

$$\hat{\theta} = \frac{f_{11} f_{22}}{f_{21} f_{12}} = 107.76$$

|         | Minor         | Major         |
|---------|---------------|---------------|
| Gro III | $f_{11} = 53$ | $f_{12} = 3$  |
| Gro I   | $f_{21} = 10$ | $f_{22} = 61$ |

The success 'minor role' in Gro.III is more likely than Gro.I.



# PRACTICAL PROBLEMS FROM VITAL STATISTICS

Q.1. Consider the following data set for two countries [CU'92, 09]

Country I

| <u>Age Group</u> | <u>Population size<br/>(in '000) Px</u> | <u>No. of Males<br/>(in '000) m Px</u> | <u>No. of<br/>Deaths<br/>Dx</u> | <u>Number of<br/>deaths in males<br/>m Dx</u> |
|------------------|-----------------------------------------|----------------------------------------|---------------------------------|-----------------------------------------------|
| < 1              | 100                                     | 50                                     | 750                             | 400                                           |
| 1-4              | 400                                     | 200                                    | 3000                            | 1600                                          |
| 5-20             | 1500                                    | 750                                    | 12000                           | 6500                                          |
| 21-100           | 4000                                    | 1500                                   | 31000                           | 12000                                         |

Country II

| <u>Age Group</u> | <u>Population size<br/>(in '000)</u> | <u>No. of Males<br/>(in '000)</u> | <u>No. of<br/>Deaths</u> | <u>Number of<br/>deaths in males</u> |
|------------------|--------------------------------------|-----------------------------------|--------------------------|--------------------------------------|
| < 1              | 15                                   | 7.5                               | 12                       | 6                                    |
| 1-4              | 60                                   | 30                                | 45                       | 24                                   |
| 5-20             | 200                                  | 100                               | 160                      | 85                                   |
| 21-100           | 440                                  | 170                               | 350                      | 200                                  |

- (i) Compute the CDR for both the countries.
- (ii) Compute the ASDR for both the countries separately for male, female and total population.
- (iii) In order to compare the mortality situation of the two countries propose a good measure and calculate its value for both the countries.

Q.2: [CU'05]

Fill in the blanks of the following table which are marked with question marks:

| <u>Age(x)</u> | <u><math>l_x</math></u> | <u><math>d_x</math></u> | <u><math>q_x</math></u> | <u><math>p_x</math></u> | <u><math>l_{x+1}</math></u> | <u><math>T_x</math></u> | <u><math>r_x</math></u> |
|---------------|-------------------------|-------------------------|-------------------------|-------------------------|-----------------------------|-------------------------|-------------------------|
| 20            | 6,93,435                | ?                       | ?                       | ?                       | ?                           | 35081126                | ?                       |
| 21            | 6,90,673                | -                       | -                       | -                       | -                           | -                       | ?                       |

Q.3. The number of persons dying at age 75 is 476 and the complete expectation of life at 75 and 76 years are 3.92 and 3.66 years. Find the numbers living at ages 75 and 76.

Q.4. [CU'08,95] From the following data relating to a particular community compute the GRR and the NRR. Interpret your results.

| <u>Age of mother</u> | <u>No. of Women</u> | <u>No. of Births</u> | <u>Survival factor</u> |
|----------------------|---------------------|----------------------|------------------------|
| 15-19                | 9000                | 140                  | 0.920                  |
| 20-24                | 9200                | 1312                 | 0.914                  |
| 25-29                | 8900                | 1067                 | 0.908                  |
| 30-34                | 8600                | 771                  | 0.891                  |
| 35-39                | 8400                | 468                  | 0.879                  |
| 40-44                | 8500                | 160                  | 0.869                  |

Assume that 48.7% of the total are female births. Survival factor gives the Rate of survival from Birth to the mid point of the corresponding age-group.

[CU'07]

Q.5. The following table provides data of female population and the number of live births in different age groups in the USA in 2004.

| <u>Age Group (years)</u> | <u>Female Population ('000)</u> | <u>Live Births ('000)</u> |
|--------------------------|---------------------------------|---------------------------|
| 14-19                    | 20,724                          | 422                       |
| 20-24                    | 20,973                          | 1034                      |
| 25-29                    | 19,555                          | 1104                      |
| 30-34                    | 20,467                          | 966                       |
| 35-39                    | 21,050                          | 476                       |
| 40-44                    | 23,055                          | 104                       |
| 45-49                    | 22,121                          | 6                         |

Given that in 2004, the total population of the USA was 293,657,000 and the sex-ratio at birth was 105 male births per 100 female births. Compare, for the year 2004,

- (a) the CBR
- (b) the GFR
- (c) the ASFR
- (d) the TFR
- (e) the GRR.

Q.6. In 1951, the total number of live births in W.B. was estimated [CU'06] as 399680. The table below records the number of females in the child bearing age intervals and the survival rates for WB and also the number of female live births in India during the year mentioned. Can you gather any idea about the GRR and NRR for WB from the given data?

| <u>Age</u> | <u>Female Population in WB ('00)</u> | <u>Survival rate (100000)</u> | <u>Total Number of female live-births in India (10<sup>6</sup>)</u> |
|------------|--------------------------------------|-------------------------------|---------------------------------------------------------------------|
| 15-19      | 12652                                | 59753                         | 4632                                                                |
| 20-24      | 11403                                | 56924                         | 14493                                                               |
| 25-29      | 10001                                | 54032                         | 14058                                                               |
| 30-34      | 8346                                 | 50282                         | 8329                                                                |
| 35-39      | 6847                                 | 45921                         | 4036                                                                |
| 40-44      | 5696                                 | 40256                         | 2158                                                                |
| 45-49      | 4728                                 | 36205                         | 689                                                                 |

Q.7. [CU/1997] The following informations are obtained from a sample survey

| Age-group | Number of women | Proportion of Marriage | NO. of births | Proportion surviving from birth to mid point of age group among married women |
|-----------|-----------------|------------------------|---------------|-------------------------------------------------------------------------------|
| 15-19     | 16592           | 0.82                   | 2692          | 0.902                                                                         |
| 20-24     | 19137           | 0.83                   | 4272          | 0.891                                                                         |
| 25-29     | 10800           | 0.84                   | 2179          | 0.878                                                                         |
| 30-34     | 4990            | 0.85                   | 790           | 0.865                                                                         |
| 35-39     | 2463            | 0.74                   | 203           | 0.849                                                                         |
| 40-44     | 928             | 0.69                   | 47            | 0.830                                                                         |

(a) Calculate the GIRR and NRR, assuming that there was no illegitimate birth and sex ratio at birth is 1000:942 in favour of male.

(b) Determine the probability of a female children dying (after marriage) at 32 age.

Q.8. The following table gives the female ASFR for 1993 and the life-table for females of India with cohort

$$l_0 = 1,000.$$

Age of Years

| Age of Years | ASFR   |
|--------------|--------|
| 15-19        | 0.0696 |
| 20-24        | 0.2346 |
| 25-29        | 0.1897 |
| 30-34        | 0.1143 |
| 35-39        | 0.0611 |
| 40-44        | 0.0285 |
| 45-49        | 0.0101 |

Female life-table  
stationary popn.

|      |
|------|
| 4180 |
| 4123 |
| 4063 |
| 4001 |
| 3939 |
| 3860 |
| 3763 |

Compute TFR, GIRR and NRR for India assuming sex ratio at birth is 1.05:1.

Q.10. Consider a population on July 1, 1985 equal to 10,00,000 and growing at 2% per year as a continuous instantaneous rate, the crude rate of natural increase for 1995 was 3%, and the crude rate con 1%. Determine the number of live births in 1995.

Q.11. [CU'07] For a stationary population with radix,  $l_0 = 1,00,000$ , out of the children born in 1980, number of deceased was 20,000 and the number of deceased in 1981 was 5000. Given the ASDRs of this population for the following ages (l.b.d)

| $x$   | 0     | 1    | 2     | 3     | 4     |
|-------|-------|------|-------|-------|-------|
| $m_x$ | .1073 | .021 | .0023 | .0017 | .0012 |

(a) Compute complete expectation of life at age 4, given the same at birth is 65.89 years.

(b) What is the chance that two newborn babies will survive 4 years after their birth?

Q.9. [CU'03] The following table gives the ten decennial census popln. of two countries, say, A and B.

| Year :                           | 1901  | 1911  | 1921  | 1931  | 1941  | 1951  | 1961  | 1971  |
|----------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|
| Popln. of Country A (in million) | 283.3 | 252.0 | 251.2 | 278.9 | 318.5 | 361.6 | 439.1 | 547.0 |
| Popln. of Country B (in million) | 5.3   | 7.2   | 9.6   | 12.9  | 17.1  | 23.2  | 31.4  | 38.5  |

| Year :                     | 1981  | 1991  |
|----------------------------|-------|-------|
| Popln. of C-A (in million) | 653.8 | 823.8 |
| Popln. of C-B (in million) | 50.2  | 62.9  |

Justify by some graphical procedure which of these countries population growth follows approximately Logistic model and comment on it. Fit a logistic model to the population that fits as well.

### PROBLEMS ON VITAL STATS.

7) (iii) To compare the mortality situation of country I and country II, we consider STDR based on age specific death rates,

$$STDR_I = \frac{\sum_m m_{lx}^m p_x^S + \sum_f m_{fx}^f p_x^S}{\sum_m m_{lx}^m p_x^S + \sum_f m_{fx}^f p_x^S}$$

Hence, we select  $\frac{m}{n} p_x^S = \frac{m p_x^I + m p_x^{II}}{2}$ ,  $\frac{f}{n} p_x^S = \frac{f p_x^I + f p_x^{II}}{2}$

2)  $d_x = l_x - l_{x+1}$   
 $\therefore d_{20} = 6,93,435 - 6,90,673 = 2762$

$$\therefore q_{rx} = \frac{d_x}{l_x} \Rightarrow q_{20} = \frac{d_{20}}{l_{20}} = 3.98 \times 10^{-3}$$

$$\therefore p_x = 1 - q_{rx} \Rightarrow p_{20} = 1 - q_{20} = 0.996$$

$l_{rx} = \frac{l_x + l_{x+1}}{2}$ , assuming that deaths are uniformly distributed.

$$\therefore l_{20} = \frac{1}{2}(l_{20} + l_{21}) = 692054.$$

$$e_x^o = \frac{T_x}{l_x} = \frac{T_{20}}{l_{20}} = e_{20}^o = 50.59$$

$$\therefore T_x = l_x + l_{x+1}$$

$$e_{21}^o = \frac{T_{21}}{l_{21}}$$

$$\Rightarrow T_{x+1} = T_x - l_x$$

$$\Rightarrow T_{21} = T_{20} - l_{20} = 34389072,$$

3&gt;

$$e_{75}^o = 3.92, e_{76}^o = 3.66$$

$$d_{75} = 476$$

$$\Rightarrow l_{75} = -l_{76} = 476$$

$$\therefore e_{75}^o \approx e_{75} + \frac{1}{2}$$

$$\Rightarrow e_{75} = e_{75}^o - \frac{1}{2} = 3.92 - 5 = 3.42$$

$$\therefore e_{76} = 3.66 - 5 = 3.16$$

$$\begin{aligned} \therefore \frac{e_{76}}{e_{75}} &= \frac{\left( \sum_{t=1}^{\infty} l_{76+t} \right) / l_{76}}{\left( \sum_{t=1}^{\infty} l_{75+t} \right) / l_{75}} \\ &= \frac{l_{75}}{l_{76}} \times \frac{\sum_{t=1}^{\infty} l_{76+t}}{l_{76} + \sum_{t=1}^{\infty} l_{76+t}} \\ &= \frac{l_{75}}{l_{76}} \times \frac{1}{\frac{l_{76}}{\sum_{t=1}^{\infty} l_{76+t}} + 1} \\ &= \frac{l_{75}}{l_{76}} \times \frac{1}{1 + \frac{1}{e_{76}}} \end{aligned}$$

$$\Rightarrow \frac{l_{76}}{l_{75}} = \frac{e_{75}}{1 + e_{76}}$$

$$\Rightarrow l_{76} = l_{75} \times \left( \frac{e_{75}}{1 + e_{76}} \right)$$

$$= (476 + l_{76}) \times \frac{e_{75}}{1 + e_{76}}$$

$$\Rightarrow l_{76} \left( 1 - \frac{e_{75}}{1 + e_{76}} \right) = 476 \times \frac{e_{75}}{1 + e_{76}}$$

$$\Rightarrow l_{76} = \frac{476 \times e_{75}}{1 + e_{76} - e_{75}} = 2200$$

$$\therefore l_{75} = 476 + 2200 = 2676.$$

4. Age-specific fertility rate of age group (20-25) is

$$f_{ix} = \frac{f_B z}{f_P x}$$

$$\therefore GRR = 5 \times \sum_{x=1}^5 \frac{f_x B_x}{f_p x} = 5 \times \sum_{x=1}^5 f_x$$

$$NRR = 5 \sum_{\text{fix}} f_i \times \frac{s}{f_i} R_i$$

cohere if  $R_x$  is the survivor factor per year in the age group  $(x, x+5)$ .

$$\text{Hence, } \frac{\frac{f}{5}Bx}{\frac{1}{5}Bx} = 48.7\% = \frac{48.7}{100} \times x.$$

$$C_{IRR} = 5 \times \sum_{x=5}^{\infty} \frac{f_i x}{i} = 5 \times 2153679 = 10768$$

$$\therefore NRR = 5 \times \frac{1}{5} \sum_{x=5}^{\infty} ix \times \frac{1}{5} R_x = 5 \times 1.9419 = 0.9707$$

"GRR = 1.0768" indicates the no. 1.0768 of daughters would be born on the average, to each of a group of females beginning life together, supposing none of them died before reaching the end of the reproductive period and all of them experience, throughout the reproductive period.

"GRRR = 9707" implies a group of 10000 females is expected to be replaced by 9707 females in the next generation under the given rates of fertility and mortality and the population will show a tendency to decrease. Hence the population will ultimately decrease and will ultimately die out, unless the fertility and mortality change.

$$5) \quad P = 293,657,000$$

$$\frac{m_B}{f_B} = \frac{105}{100}$$

$$(a) CBR = \frac{B}{P} \times 1000$$

$$= \frac{4112,000}{293,657,000} \approx 14$$

$$(b) GFR = \frac{B}{\sum_x f_s P_x} \times 1000$$

$$(c) ASFR = \frac{5Bx}{5f_p x} \times 1000 = 5ix$$

$$(d) TFR = 5 \times \sum_x 5ix = 5 \times \text{sum of ASFR's.}$$

$$(e) GIRR = \left( TFR \times \frac{f_B}{B} \right) ; \quad GRR = \frac{100}{205} \times \frac{TFR}{1000}$$

%

6.  $B = 399680$   
 To compute GIRR & NRR of W.B., we should have information regarding  $\frac{f}{5}Bx$ .

$$\text{Assuming, } \frac{\frac{f}{5}Bx}{fB} = \frac{\frac{f}{5}Bx^I}{fB^I} \forall x$$

$$\Rightarrow \frac{f}{5}Bx = \frac{\frac{f}{5}Bx^I}{fB^I} \times fB$$

$$\text{Assume, } \frac{fB}{B} = \frac{100}{205} \Rightarrow fB = \frac{100}{205} \times B$$

$$= \frac{100}{205} \times 399680 \\ = 194966$$

| Age          | Female Popn.<br>$\frac{f}{5}Px$ | Survival<br>Rate ( $\frac{f}{5}Rx$ ) | Total no. of<br>female live<br>births in India<br>$\frac{f}{5}Bx^I$ | $\frac{f}{5}Bx$ | $\frac{f}{5}ix$ | $\frac{f}{5}ix \times \frac{f}{5}Rx$ |
|--------------|---------------------------------|--------------------------------------|---------------------------------------------------------------------|-----------------|-----------------|--------------------------------------|
| 15-19        | 1265200                         | .59753                               | 463200                                                              | 18680           | .01976          | .00882                               |
| 20-24        | 1140300                         | .56924                               | 1444300                                                             | 58248           | .05108          | .02908                               |
| 25-29        | 1000100                         | .54082                               | 1405800                                                             | 56693           | .05869          | .03063                               |
| 30-34        | 886600                          | .50352                               | 832900                                                              | 83589           | .04015          | .02022                               |
| 35-39        | 684700                          | .45921                               | 403600                                                              | 16276           | .02377          | .01091                               |
| 40-44        | 589600                          | .40356                               | 215800                                                              | 8703            | .01528          | .00617                               |
| 45-49        | 472800                          | .36205                               | 68900                                                               | 2779            | .00588          | .00213                               |
| <b>TOTAL</b> | <b>5969800</b>                  |                                      | <b>4834500</b>                                                      | <b>194966</b>   | <b>.20761</b>   | <b>.10797</b>                        |

$$\frac{f}{5}Bx = \frac{\frac{f}{5}Bx^I}{fB^I} \times 194966, \quad \frac{f}{5}ix = \frac{\frac{f}{5}Bx}{\frac{f}{5}Px}$$

$$\therefore GIRR = 5 \times \sum_x \frac{f}{5}ix \\ = 5 \times 0.20761 = 1.03805$$

$$\therefore NRR = 5 \times \sum_x \frac{f}{5}ix \times \frac{f}{5}Rx \\ = 5 \times 0.10797 \\ = 0.53985,$$

7)   
 (a) there was no illegitimate birth. We consider  
 $\frac{f_{Px}}{S} = \text{No. of women in age group } [x, x+5] \times$   
 proportion of marriage.

$$\frac{f_B}{mB} = \frac{942}{1000} \Rightarrow \frac{f_B}{B} = \frac{942}{1942}$$

$$\therefore GRR = S \times \frac{f_B}{B} \times \sum \frac{s_{Bx}}{f_{Px}} = S \times \frac{f_B}{B} \times \sum s_{ix}$$

| Age Group    | No. of women | Proportion of marriage | $\frac{f_{Px}}{S}$ | $s_{Bx}$ | $s_{ix}$       | $\frac{f_{Rx}}{S}$ | $\frac{f_{SRx} \times s_{ix}}{S}$ |
|--------------|--------------|------------------------|--------------------|----------|----------------|--------------------|-----------------------------------|
| 15-19        | 16592        | 0.82                   | 13605              | 2642     | 0.19419        | 0.902              | 0.17516                           |
| 20-24        | 19137        | 0.83                   | 15889              | 4272     | 0.26895        | 0.891              | 0.28963                           |
| 25-29        | 10860        | 0.84                   | 9122               | 2179     | 0.23887        | 0.878              | 0.20973                           |
| 30-34        | 4990         | 0.85                   | 4241               | 790      | 0.18628        | 0.865              | 0.16113                           |
| 35-39        | 2463         | 0.74                   | 1823               | 203      | 0.11138        | 0.849              | 0.09486                           |
| 40-44        | 925          | 0.64                   | 592                | 47       | 0.07939        | 0.830              | 0.06589                           |
| <b>TOTAL</b> | <b>—</b>     | <b>—</b>               | <b>—</b>           | <b>—</b> | <b>1.07906</b> | <b>—</b>           | <b>0.57328</b>                    |

$$\therefore GRR = S \times \frac{942}{1942} \times 1.07906 = 2.61708$$

$$\therefore NRR = S \times \frac{942}{1942} \times 0.57328 = 1.3909$$

(b) Probability of the female children dying after marriage at 32 age  
 $= \{1 - \text{Probability of surviving at age } 32\} \times$   
 $\{\text{Prob. of married}\}$

$$= (1 - 0.865) \times 0.85$$

$$= 0.11475$$

8)  $NRR = \sum \frac{s_{ix}}{S} \times \frac{f_{SLix}}{f_{L0}}$ , where  $f_{SLix} = f_{Lix} + f_{L_{x+1}} + \dots + f_{L_{x+4}}$

| Age Group    | ASFR ( $s_{ix}$ ) | $\frac{f_{Lix}}{S}$ | $\frac{f_{R0}}{f_{L0}[=100]}$ | $s_{ix} \times f_{R0}$ |
|--------------|-------------------|---------------------|-------------------------------|------------------------|
| 15-19        | 0.0696            | 4180                | 4.180                         | ·29023                 |
| 20-24        | 0.2346            | 4123                | 4.123                         | ·96726                 |
| 25-29        | 0.1897            | 4063                | 4.063                         | ·77075                 |
| 30-34        | 0.1143            | 4001                | 4.001                         | ·45931                 |
| 35-39        | 0.0611            | 3934                | 3.934                         | ·24037                 |
| 40-44        | 0.0289            | 3860                | 3.860                         | ·11001                 |
| 45-49        | 0.0101            | 3763                | 3.763                         | ·03801                 |
| <b>TOTAL</b> | <b>0.7079</b>     | <b>27924</b>        | <b>2.924</b>                  | <b>2.87399</b>         |

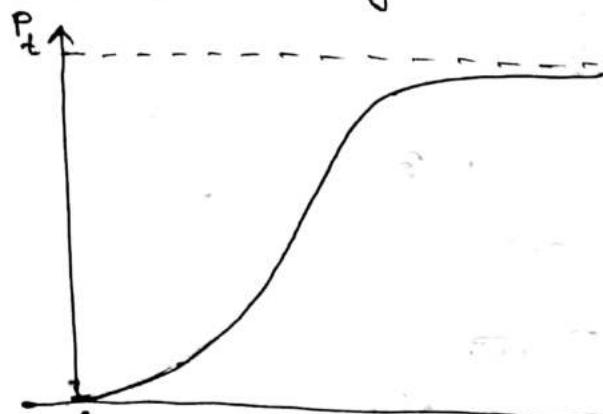
$$\therefore TFR = S \sum s_{ix} = S \times 0.7079 = 3.5395$$

$$\therefore GRR = \frac{f_B}{B} \times TFR = \frac{100}{205} \times 3.5395 = 1.7266$$

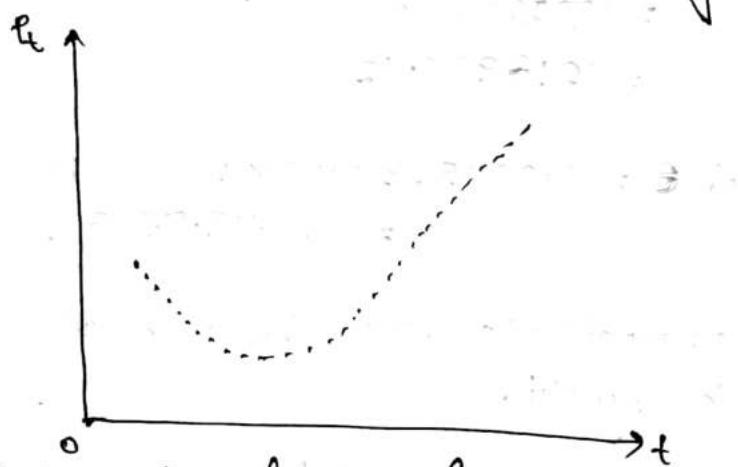
$$\therefore NRR = \frac{f_B}{B} \times \sum s_{ix} \times f_{R0} = \frac{100}{205} \times 2.87399 = 1.40192$$

9)

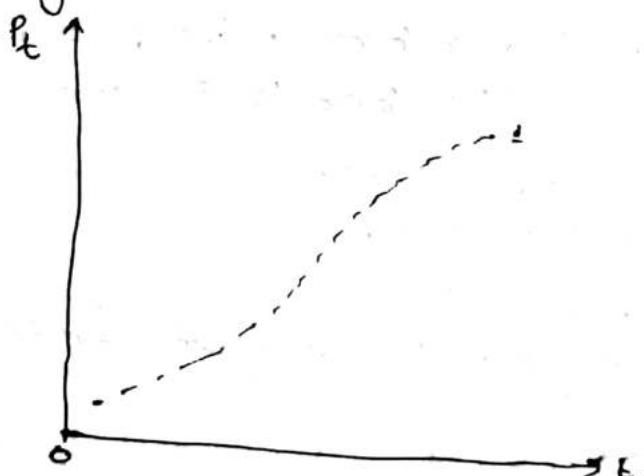
The shape of the Logistic curve is given by:



Plot  $(t, P_t)$  for both the population on graph paper.



Clearly the population figures for the country I is not following the logistic law.



Clearly, II population figures are following logistic law approximately.

10)  $CBR - CPR = \text{Crude rate of natural increase}$   
 $= \left( \frac{B}{P} - \frac{D}{P} \right) 100\%.$   
 $= 3\%.$

$$\Rightarrow \frac{B}{P} - \frac{D}{P} = 0.03 \quad \left[ \frac{D}{P} \times 100 = 1 \right].$$

$$\Rightarrow \frac{B}{P} = 0.04$$

$$P_0 = 10,000,00$$

$$P_{10} = P_0 (1 + 0.02)^{10}$$

$$= 10,000,000 (1 + 0.02)^{10}$$

$$= 1218994.42$$

$$\therefore B = 0.04 \times P_{10} = 0.04 \times 1218994.42$$

$$= 48759.7768$$

11) (a) To calculate the value  $a_0$ , we use the formula given by Kuczynski.

$$\hat{a}_0 = 1 - \left( 1 - \frac{D''}{B_0} \right) \left( 1 - \frac{D'''}{B_{-1} - D'} \right)$$

$B_{-1}$  = No. of children born in the preceding calendar year.

$B_0$  = No. of children born in the current calendar year.

$D'$  = No. of children born and deceased in the preceding calendar year,

$D''$  = No. of children born in the preceding calendar year & deceased in the current calendar year before reaching age 1 and

$D'''$  = No. of children born and deceased in the current calendar year.

$$B_{-1} = 1,00,000$$

$$B_0 = 1,00,000$$

$$D' = 20,000$$

$$D'' = 5,000$$

$$D''' = 25,000$$

$$\hat{a}_0 = 1 - \left( 1 - \frac{20000}{100000} \right) \left( 1 - \frac{5000}{80000} \right)$$

$$= 0.25$$

$$m_x = \frac{dx}{lx}$$

$$\Rightarrow m_x = \frac{dx}{lx - (1-\alpha_x)dx}$$

$$= \frac{dx}{1 - (1-\alpha_x)dx}$$

$$\begin{aligned}\alpha_1 &= 0.43 \\ \alpha_2 &= 0.45 \\ \alpha_3 &= 0.47 \\ \alpha_4 &= 0.49\end{aligned}$$

(b) Required probability =  $\left(\frac{\lambda_4}{\lambda_0}\right)^{\sim}$ .



## PROBLEMS ON INDEX NUMBER

FOR PRACTICAL EXAM

► The following table shows the quantities consumed and the values (price  $\times$  quantity) of 5 commodities for 3 successive years.

| Commodities | 1990     |       | 1991     |       | 1992     |       |
|-------------|----------|-------|----------|-------|----------|-------|
|             | Quantity | Value | Quantity | Value | Quantity | Value |
| I           | 50       | 350   | 60       | 420   | 70       | 490   |
| II          | 120      | 600   | 140      | 700   | 180      | 800   |
| III         | 30       | 330   | 20       | 200   | 15       | 225   |
| IV          | 20       | 360   | 15       | 300   | 10       | 220   |
| V           | 5        | 40    | 5        | 50    | 5        | 60    |

Calculate the price index numbers for 1992 taking 1990 as the base period adopting chain base formula and using Paasches formula at each stage. Also verify whether the circular test is satisfied by the Paasches formula or not on the basis of the above data. [C.U. 1996]

Soln. → The chain index number for 1992 cont. base 1990 is

$$P'_{90,92} = P_{90,91} \times P_{91,92}$$

$$= \frac{\sum p_{91} q_{91}}{\sum p_{90} q_{91}} \times \frac{\sum p_{92} q_{92}}{\sum p_{91} q_{92}}, \text{ by Paasches formula.}$$

Note that, Price ( $p$ ) =  $\frac{\text{Value}}{\text{Quantity}}$ .

For circular test,  $P_{90,91} \cdot P_{91,92} \cdot P_{92,90} = 1$ .

$$\Leftrightarrow P_{90,91} \cdot P_{91,92} = P_{90,92}$$

$$\Leftrightarrow P_{90,92} = P_{90,92}$$

i.e. chain index = fixed base index.

Hence,  $P_{90,92} = \frac{\sum p_{92} q_{92}}{\sum p_{90} q_{92}} =$

Hence,  $P'_{90,92} \neq P_{90,92}$

⇒ The circular test is not satisfied by the Paasches formula.

2. The following table shows the average per capita consumption of cereals and prices of cereals in rural India for four different periods  $T_1, T_2, T_3, T_4$ .

| Commodities | Consumptions (in kg) per month |       |       |       | Price (Rs. per kg) |       |       |       |
|-------------|--------------------------------|-------|-------|-------|--------------------|-------|-------|-------|
|             | $T_1$                          | $T_2$ | $T_3$ | $T_4$ | $T_1$              | $T_2$ | $T_3$ | $T_4$ |
| Rice        | 4.69                           | 3.94  | 3.58  | 4.28  | 5.37               | 5.75  | 6.15  | 6.30  |
| Wheat       | 5.51                           | 6.93  | 6.43  | 6.78  | 4.94               | 5.15  | 5.02  | 5.10  |
| Others      | 7.66                           | 8.19  | 7.75  | 7.71  | 3.32               | 2.94  | 3.25  | 3.52  |

Calculate the price index numbers of cereals for period  $T_4$  taking  $T_1$  as the base adopting chain-base formula and using any uniformly suitable formula at each stage. If after calculations it is found that consumption figures in  $T_4$  are in error by 5%, do you think that your index is going to be affected by this? Give reason for your answer. [C.U. 1992]

Soln. → If we use Laspeyres's formula, then

$$P_{14}' = P_{12} \cdot P_{23} \cdot P_{34}$$

$$= \frac{\sum p_2 q_1}{\sum p_1 q_1} \times \frac{\sum p_3 q_2}{\sum p_2 q_2} \times \frac{\sum p_4 q_3}{\sum p_3 q_3} = \dots$$

is independent of  $q_4$  (the consumption figure for  $T_4$ ). Hence, it will not affect the calculation of chain index.

If we use Paasche's formula, then

$$P_{14}' = P_{12} \cdot P_{23} \cdot P_{34}$$

$$= \frac{\sum p_2 q_2}{\sum p_1 q_2} \times \frac{\sum p_3 q_3}{\sum p_2 q_3} \times \frac{\sum p_4 q_4}{\sum p_3 q_4} = \dots$$

which depends on  $q_4$ .

Hence, all the consumption figures in  $T_4$  are in error by 5%, then the corrected figures are

$$q_{4i}' = q_{4i} \times \frac{105}{100} \text{ or } q_{4i} \times \frac{95}{100}$$

$$\text{and hence, } \frac{\sum p_4 i q_{4i}'}{\sum p_3 i q_{4i}} = \frac{\sum p_4 i \frac{105}{100} q_{4i}}{\sum p_3 i \frac{105}{100} q_{4i}} = \frac{\sum p_4 i q_{4i}}{\sum p_3 i q_{4i}}$$

Ultimately, it will not affect the calculation of the index.

3. The data below show the percentage increases in price of a few tabulated food items and weights attached to each of them. Calculate the index numbers for the food group.

|                                |      |       |     |      |     |        |      |      |        |
|--------------------------------|------|-------|-----|------|-----|--------|------|------|--------|
| Food items :                   | Rice | Wheat | Dal | Ghee | Oil | Spices | Milk | Fish | Others |
| Weights :                      | 30   | 11    | 8   | 6    | 5   | 3      | 7    | 9    | 19     |
| Percentage increase in price : | 180  | 202   | 115 | 212  | 175 | 517    | 260  | 428  | 332    |

Using the above Food Index and information given below.  
Calculate CLI.

| Group :  | Food | Clothing | Fuel and Light | Rents and Rates | Miscellaneous |
|----------|------|----------|----------------|-----------------|---------------|
| Index :  | —    | 310      | 220            | 150             | 300           |
| Weight : | 60   | 5        | 8              | 9               | 18            |

[C.U. 2000]

Soln. :  $\Rightarrow$  The percentage increase in price for the  $i$ th item

$$(\gamma) = \frac{p_{ii} - p_{oi}}{p_{oi}} \times 100$$

$$\Rightarrow \frac{p_{ii}}{p_{oi}} \times 100 = (100 + \gamma)$$

The food index is computed by

$$I_F = \frac{\sum \frac{p_{ii}}{p_{oi}} \times w_i}{\sum w_i} \times 100 = \frac{\sum \left( \frac{p_{ii}}{p_{oi}} \times 100 \right) \times w_i}{\sum w_i}$$

| food item | weights ( $w_i$ ) | $\frac{p_{ii}}{p_{oi}} \times 100$ | $\left( \frac{p_{ii}}{p_{oi}} \times 100 \right) w_i$        |
|-----------|-------------------|------------------------------------|--------------------------------------------------------------|
|           |                   |                                    |                                                              |
|           |                   |                                    |                                                              |
|           | $\sum w_i =$      |                                    | $\sum \left( \frac{p_{ii}}{p_{oi}} \times 100 \right) w_i =$ |

Hence,  $I_F = \frac{\text{Column 4 total}}{\text{Column 2 total}} = \dots$

$$\text{The CLI} = \frac{\sum I_i w_i}{\sum w_i}$$

4. The following table shows the group indices and the corresponding weights for the year 1995 with 1981 as the base year of a given community.

| <u>Group</u>   | <u>Group index</u> | <u>Weight</u> |
|----------------|--------------------|---------------|
| Food           | 212.45             | 65.3          |
| Clothing       | 328.06             | 4.8           |
| Fuel and light | 345.89             | 8.5           |
| Housenrent     | 173.841            | 7.6           |
| Miscellaneous  | 201.35             | 13.8          |

- (a) Find the CLI for the year 1995.  
 (b) What is the purchasing power in 1995 as compared to 1981.  
 (c) If Mr. Dasgupta's salary increased from Rs. 2400 in 1981 to Rs. 4950 in 1995, how has his economic status changed?  
 (d) The weights are proportional to the consumption expenditure of each group. Suppose Mr. Dasgupta has to maintain the same status for each of the first four groups and can only adjust his spending on Miscellaneous items to come to terms with his income changes.  
 find his spendings for each of the groups in 1995.  
 (e) What would be Mr. Dasgupta's weights for each of the groups in 1995? [C.U. 2001, 2009]

Soln. :→

$$(a) C.L.I = \frac{\sum I_i w_i}{\sum w_i} = 224.842$$

(b) the purchasing power of money in 1995 w.r.t.

$$1981 \text{ is } \frac{1}{(C.L.I/100)} = \frac{100}{C.L.I} = 0.44956.$$

(c) CLI for 1995 w.r.t. 1981 as base year is 224.842 means that if some individual spends Rs. 100 in 1981, in 1995 he/she has to spend Rs. 224.842 to maintain the same standard of living as she/he has in 1981.

For Mr. Dasgupta's expenditure on 1981 was Rs. 2400. In 1995 she has to spend  $Rs. 2400 \times 2.24842 = 5396.2$  to maintain the same standard of living as he has in 1981. But his salary in 1995 was Rs. 4950. She has to adjust  $Rs. (5396.2 - 4950) = 446.2$  i.e. his economic status falls in 1995 as compared to 1981.

| <u>Group</u>  | <u>Index</u> | <u>Weights</u> | <u>In 1981<br/>Expenditure</u>       | <u>In 1995 to expenditure<br/>to maintain the same<br/>standard of living</u> |
|---------------|--------------|----------------|--------------------------------------|-------------------------------------------------------------------------------|
| Food          | 212.45       | 65.3           | $2400 \times 65.3 / 100$<br>= 1567.8 | $1567.8 \times 212.45 / 100$<br>= 3329.51                                     |
| Clothing      | 328.06       | 4.8            | 115.2                                | 377.925                                                                       |
| Fuel & light  | 345.89       | 8.5            | 204                                  | 705.615                                                                       |
| Housenent     | 178.41       | 7.6            | 182.4                                | 316.299                                                                       |
| Miscellaneous | 201.85       | 13.8           | 326.4                                | 666.206                                                                       |
| TOTAL =       |              | 100            | 2400                                 | 5396.2                                                                        |

$$\text{In 1981, expenditure on a group} = \text{Income} \times \frac{w_i}{\sum w_i}$$

Hence his spendings on the first four groups remain same and the deficit of Rs. 446.27 will be adjusted on Miscellaneous group and the spending on Miscellaneous is Rs. (666.86 - 446.27)  
= 220.59 Rs. in 1995.

$$(e) \text{ In 1995, the weight of a group} = \frac{\text{Expenditure on the Group}}{\text{Salary (1995)}} \times 100$$

| <u>Group</u>  | <u>Expenditure</u> | <u>Weights</u>                      |
|---------------|--------------------|-------------------------------------|
| Food          | 3329.51            | $3329.51 \times 100 / 4950 = 67.2$  |
| Clothing      | 377.925            | $377.925 \times 100 / 4950 = \dots$ |
| Miscellaneous | 220.59             | $220.59 \times 100 / 4950 = \dots$  |
| TOTAL =       | 4950               | 100                                 |

— x —

- Q. 5. The following data relate to the urban middle-class people of a particular region in the years 2003 and 2005.

| Group         | % of total expenditure | Group Index (Base: 2000) |      |
|---------------|------------------------|--------------------------|------|
|               |                        | 2003                     | 2005 |
| Food          | 35                     | 118                      | 122  |
| Clothing      | 15                     | 112                      | 118  |
| Housing       | 20                     | 113                      | 115  |
| Transport     | 10                     | 112                      | 117  |
| Durable Goods | 8                      | 105                      | 110  |
| Others        | 12                     | 120                      | 125  |

(a) Compute the CLI for 2003 and 2005 with Base year 2000.

(b) If a family saved 15% of its monthly income in 2003, find the relative change in its average savings over the years 2003-2005, assuming that it maintained the same standard of living as in 2003. (C.U.- 2008)

Soln.

(a)

$$CLI_{2003} = \frac{\sum I_i w_i}{\sum w_i} = \frac{11470}{100} = 114.70$$

$$CLI_{2005} = \frac{\sum I_i w_i}{\sum w_i} = \frac{11890}{100} = 118.90$$

∴ Salary is constant, i.e., Rs. 100.

So, weights of the group = Expenditure of the group.

(b) Let Rs. 100 be the salary in 2003. Saving, Rs. = 15.

Expenditure = Rs. 80.5

| CL.I. | 2003   | 2005   |
|-------|--------|--------|
|       | 114.70 | 118.90 |
|       |        |        |

If some individual spent Rs. 114.70 in 2003, in 2005 he has spent Rs.  $\frac{118.90}{114.70} \times 80.5 = \cancel{80.5} = 88.11$  Rs.

In 2005, his salary is 100 Rs.

Exp. ~~Rs. 80.5~~ = 88.11 Rs.

$$\text{Saving} = 11.89 \text{ Rs.} \quad \therefore \text{Saving \%} = \frac{11.89}{100} \times 100 \\ = 11.89\%$$

6. The following table gives (with two missing values) the overall and groupwise CLI (with 2000 as the base year) with six different expenditure groups and their respective weights, for the urban middle class people of a particular city, in 2004 and 2005.

| Group         | Weight | Group Index |      |
|---------------|--------|-------------|------|
|               |        | 2004        | 2005 |
| Food          | 350    | 117         | 120  |
| Clothing      | 156    | 113         | 118  |
| Housing       | 187    | 118         | -    |
| Transport     | 108    | 112         | 117  |
| Durable Goods | 76     | 102         | 111  |
| Others        | 123    | 121         | 125  |

CLI       $\frac{\text{TOTAL} = 1000}{= 115.4} \rightarrow = 119.5$

(a) If Mr. X saved 20% of his salary in 2004, Determine the relative change in his average savings, relative to 2004, in 2005 if his salary increased by 10% and he maintained the same standard of living as in 2004. [C.O. 2007]

Soln. (a) Let Rs. 100 be the salary in 2004. Saving = Rs. 20

$$\text{Exp.} = \text{Rs. } 80,$$

| CLI . | 2004   2005 |       |
|-------|-------------|-------|
|       | 115.4       | 119.5 |

If X spends Rs. 115.4 in 2004, in 2005, he has spent.  
 $\text{Rs. } \frac{119.5}{115.4} \times 80 = 82.84 \text{ Rs.}$

In 2005, his salary is 110.

$$\text{Exp.} = \text{Rs. } 82.84.$$

$$\text{Saving} = \text{Rs. } 27.16$$

$$\text{Saving \%} = \frac{27.16}{110} \times 100 = 24.69\%$$

## PRACTICALS ON PROBABILITY DISTRIBUTION

1. (a) Given that  $\sum_{i=1}^{10} f_i = 500420$ ,  $\sum_{i=1}^{10} f_i = 329240$ ,  $\sum_{i=1}^{10} f_i = 175212$ ,  
 $f_{10} = 40365$ . Find  $f_1$ .  
(b)

1. (a) The no. of females in each of 100 queues of length 10 at a metro railway station in kolkata. The data are shown below:

| Count | 0 | 1 | 2 | 3  | 4  | 5  | 6  | 7 | 8 | 9 | 10 |
|-------|---|---|---|----|----|----|----|---|---|---|----|
| Freq  | 1 | 3 | 4 | 23 | 25 | 19 | 18 | 5 | 1 | 1 | 0  |

Propose an appropriate theoretical distn and fit it to the above data. Also comment on the fitting.

Solution:- Let  $X$ : 'No' of females in a queue, of length 10  
Consider "getting a female in queue" as a success.

Assuming probability of successes in each queue is  $p$  (constant).  
Then, Binomial model is appropriate to the given RV  $X$ .

Then PMF of  $X$  is

$$P(X) = \begin{cases} \binom{10}{x} p^x (1-p)^{10-x} & ; x=0(1)10 \\ 0 & ; \text{ow} \end{cases}$$

By method of moments,  $\mu'_1 = m'_1$ .

$$\Rightarrow 10p = \bar{x} = \sum_{i=1}^{10} xi / \sum_{i=1}^{10} f_i = 435/100 = 4.35.$$

$$\Rightarrow \hat{p} = \frac{\bar{x}}{10} = 0.435.$$

Hence the fitted distn is given by :

$$p(x) = \begin{cases} \binom{10}{x} (\hat{p})^x (1-\hat{p})^{10-x} & ; x=0(1)10 \\ 0 & ; \text{ow} \end{cases}$$

Comment on Fitting :- Expected frequency of the value 'X' is  
 $N \times P[X = x]$   
 $= 100 \times p(x)$ .

Consider getting a value ' $x$ ' as a success. Then  
 $f_x$  = the frequency of ' $x$ ' in  $N$  observations  
     = the no. of successes in  $N$  Bernoulli trials  
     ~  $\text{Bin}(N, P[X=x])$   
 $\therefore E(f_x) = N \times P[X=x]$

| $x$ | $p(x)$ | Exp. freq = $N \cdot p(x)$ | Observed freq. |
|-----|--------|----------------------------|----------------|
|     |        |                            |                |

(b) When the first proof of 200 pages of an encyclopaedia of 500 pages were read, the distn. of printing mistakes were found to be as shown in the table:

|                           |     |    |    |   |   |   |
|---------------------------|-----|----|----|---|---|---|
| No. of misprints on Page: | 0   | 1  | 2  | 3 | 4 | 5 |
| Frequency :               | 112 | 63 | 20 | 3 | 1 | 1 |

Fit a suitable probability distribution to the data.

Establish the total cost of correcting the first proof of the whole encyclopaedia by using the information given below:

No. of misprints on page: 0      1      2      3      4      5 or more  
 Cost of detection and  
 correction (dollars) per page: 0.10      0.18      0.23      0.29      0.34      0.36

Solution:- Let  $X$  denotes the number of misprints on a page. Getting a misprint on a page is a rare event. Hence  $X$  is expected to follow Poisson distribution.

Assume that  $X \sim P(\lambda)$ .

The PMF of  $X$  is

$$p(x) = \begin{cases} e^{-\lambda} \cdot \frac{\lambda^x}{x!}, & x=0,1,2,3,\dots \\ 0, & \text{ow} \end{cases}$$

cohere  $\lambda > 0$ .

By method of moments,  $\mu'_i = m'_i$

$$\Rightarrow \hat{\lambda} = \bar{x} = \frac{\sum_{i=0}^s x_i f_i}{\sum f_i} = \underline{\quad}$$

Hence the fitted distn. is

$$p(x) = \begin{cases} e^{-\hat{\lambda}} \cdot \frac{(\hat{\lambda})^x}{x!}, & x=0,1,2,\dots \\ 0, & \text{ow} \end{cases}$$

The expected frequency of the value ' $x$ ' is

$$N.P[X=x] = 5000 \times p(x)$$

Table showing expected freq. and cost of correction:-

| $x$   | $p(x)$ | Exp. freq. | Cost per page | cost of correction<br>$= \text{Exp. freq.} \times \text{Cost per page}$ |
|-------|--------|------------|---------------|-------------------------------------------------------------------------|
| 0     |        |            | 0.10          |                                                                         |
| 1     |        |            | 0.16          |                                                                         |
| 2     |        |            | 0.23          |                                                                         |
| 3     |        |            | 0.29          |                                                                         |
| 4     |        |            | 0.34          |                                                                         |
| $> 5$ |        | 1          | 0.36          |                                                                         |
|       |        | N          |               |                                                                         |

2. (a) A farmer sells bean seeds in packets of 100 and agrees to refund the price if the number of seeds germinating from a packet is less than 75. He knows from past experience that on an average 80% of the seeds germinate and it costs him Rs. 200 for a packet of seeds. How should he fix the price of a packet so as to ensure an average profit of 25%?

Solution:-  $X$ : The no. of seeds germinated in packet.

Here,  $X \sim \text{Bin}(100, p=0.8)$ .

Let  $x$  be the selling price.

Define,  $Z = \begin{cases} (x-2) & \text{if } x \geq 75 \\ -2 & \text{if } x < 75 \end{cases}$

$$\text{Now, } E(Z) = \frac{x}{4} = (x-2)P[X \geq 75] + (-2)P[X < 75]$$

$$= (x-2)[1 - P[X < 75]] - 2P[X < 75]$$

By CLT,  $\frac{x-80}{\sqrt{16}} \approx N(0,1)$

$$P[X \geq 75] = 1 - \Phi\left(\frac{5}{4}\right)$$

$$\Rightarrow P[X < 75] = \Phi\left(\frac{5}{4}\right)$$

(b) Let  $X \sim \text{Bin}(2, p)$ ,  $Y \sim \text{Pois}(\lambda=1)$  and  $X$  and  $Y$  are independently distributed. It is known that  $P[Y > X] = 1 - 2e^{-1}$ . What is the probability that  $X$  is strictly positive?

$$\begin{aligned} 1 - 2e^{-1} &= P[Y > X] = 1 - P[Y \leq X] \\ &= 1 - \sum_{x=0}^2 P[Y \leq x, X=x] \\ &= 1 - \left[ P[Y \leq 0]P[X=0] + P[Y \leq 1]P[X=1] \right. \\ &\quad \left. + P[Y \leq 2]P[X=2] \right] \\ &= 1 - \left[ e^{-1} \left( \frac{2}{0} \right) k^0 (1-p)^0 + \left\{ e^{-1} + \frac{e^{-1}}{1!} \right\} \left( \frac{2}{1} \right) \right. \\ &\quad \left. p^0 (1-p)^2 + \left\{ e^{-1} + \frac{e^{-1}}{1!} + \frac{e^{-1}}{2!} \right\} \left( \frac{2}{2} \right) p^2 \right] \end{aligned}$$

$$\Rightarrow 1 + 2p - \frac{p^2}{2} = 2$$

$$\Rightarrow p^2 - 4p + 2 = 0$$

$$\Rightarrow p = \frac{4 \pm \sqrt{4^2 - 4 \cdot 1 \cdot 2}}{2}$$

$$\Rightarrow p = 2 \pm \sqrt{2}$$

\*  $P[Y \leq 0] = P[Y=0]$   
 $P[Y \leq 1] = P[Y=0] + P[Y=1]$ .

$$\begin{aligned}
 \text{Now, required probability is } &= P[X > 0] \\
 &= 1 - P[X = 0] \\
 &= 1 - (1-p)^2 \\
 &= [1 - (\sqrt{2} - 1)^2] \\
 &= 2(\sqrt{2} - 1).
 \end{aligned}$$

- (c) Suppose that the number of accidents per week at an industrial plant is a Poisson random variable with mean four. Suppose also that the number of workers injured in different accidents are independent Poisson RVs with a common mean 2. Assume that the number of workers injured in each accident is independent of the number of accidents that occur. What are the mean and variance of the number of injured during a week?

Solution:- Let  $N$  denotes the no. of accidents per week and  $X_i$  denotes the no. of workers injured in the  $i^{\text{th}}$  accident.

Here  $N \sim P(\lambda = 4)$  and  $X_i \sim P(\lambda = 2)$ , independently.

Hence, the number of injured during a week is

$$S_N = X_1 + X_2 + \dots + X_N$$

$$E(S_N) = E(N) E(X_1) = 2 \times 4 = 8.$$

- (d) A sample of size 10 is drawn from a normal population with mean and variance both equal to  $\theta (> 0)$  and the first quartile equal to 2.65. Find the probability that
- (i) the first four observations are negative
  - (ii) four of the observations are negative.

Solution:- Let  $X_1, X_2, \dots, X_{10}$  be a r.v. from  $N(0, \theta)$ ;  $\theta > 0$ .

By definition of 1st quartile,

$$P[X_1 \leq Q_1] = 0.25$$

$$\Rightarrow P\left[\frac{X_1 - \theta}{\sqrt{\theta}} \leq \frac{Q_1 - \theta}{\sqrt{\theta}}\right] = 0.25$$

$$\Rightarrow \Phi\left(\frac{Q_1 - \theta}{\sqrt{\theta}}\right) = 0.25$$

$$\Rightarrow \Phi\left(-\frac{Q_1 - \theta}{\sqrt{\theta}}\right) = 0.75$$

From Biometrika,

| $x$  | $\Phi(x)$ |
|------|-----------|
| 0.67 | 0.74857   |
| 0.68 | 0.75175   |

$$\therefore -\frac{Q_1 - \theta}{\sqrt{\theta}} = 0.6745$$

$$\Rightarrow \frac{2.65 - \theta}{\sqrt{\theta}} = 0.6745$$

$$\Rightarrow \sqrt{\theta} = \frac{0.6745 + \sqrt{(0.6745)^2 + 4 \times 2.65}}{2}$$

as  $\sqrt{\theta} > 0$ .

$$\therefore \sqrt{\theta} = 1.99 \approx 2.$$

$$\therefore \theta = 4.$$

Hence,  $X \sim N(4, 2^2)$ .

(i) The probability that first four observations are negative,

$$P[X_1 < 0, X_2 < 0, X_3 < 0, X_4 < 0]$$

$$= [P[X_1 < 0]]^4, \text{ as } X_i \text{'s are i.i.d.}$$

$$= \left[ \Phi \left[ \frac{0 - \theta}{\sqrt{\theta}} \right] \right]^4$$

$$= \left[ \Phi(-\sqrt{\theta}) \right]^4$$

$$= \Phi^4(-2) = \underline{\hspace{2cm}}$$

(ii) Let  $Y$  denotes the no. of negative values in  $X_1, X_2, \dots, X_{10}$ .  
Then  $Y \sim \text{Bin}(10, p)$ , where,  $p = P[X_1 < 0] = \Phi(-2)$

$$\therefore \text{Required Probability} = P[Y=4] = \underline{\hspace{2cm}}$$

$$= \binom{10}{4} p^4 (1-p)^6 = \underline{\hspace{2cm}}$$

3. (i) Twenty five leaves were selected at random from each of six similar apple trees. The number of adult female European red mites of each leaf was counted, the resulting information is summarized in the table below:

|                         |    |    |    |    |   |   |   |   |
|-------------------------|----|----|----|----|---|---|---|---|
| No. of Mites per leaf : | 0  | 1  | 2  | 3  | 4 | 5 | 6 | 7 |
| Frequency :             | 70 | 38 | 17 | 10 | 9 | 3 | 2 | 1 |

(a) Fit a negative binomial distribution to the data.

(b) Plot the observed and fitted values on the same graph paper, and comment on the goodness of fit after visual inspection.

Solution:- (i) Let  $X$  denotes the no. of adult female European red mites on a leaf.

(a) Assume that  $X \sim NB(r, p)$

$$\text{The PMF of } X \text{ is } f(x) = \begin{cases} \binom{x+r-1}{x} p^r q^x & ; x=0,1,2,\dots \\ 0 & ; \text{ otherwise} \end{cases}$$

where  $0 < p < 1$  and  $p+q=1$  and  $r > 0$ .

By method of moments,

$$\mu_1' = \bar{x}, \mu_2 = s^2$$

$$\Rightarrow \frac{rq}{p} = \bar{x}, \frac{rq}{p^2} = s^2$$

$$\Rightarrow \hat{p} = \frac{\bar{x}}{s^2} \text{ and } \hat{r} = \bar{x} \cdot \frac{\hat{p}}{1-\hat{p}}$$

$$\text{From data, } \bar{x} = \frac{\sum_{i=0}^7 i f_i}{\sum f_i} = \text{_____}$$

$$s^2 = \left( \sum_{i=0}^7 i^2 f_i / \sum f_i \right) - \bar{x}^2 = \text{_____}$$

$$\hat{p} = \text{_____}$$

$$\hat{r} = \text{_____}$$

The fitted NB distribution is

$$f(x) = \begin{cases} \binom{x+\hat{r}-1}{x} (\hat{p})^{\hat{r}} (\hat{r})^x & ; x=0,1,2,\dots \\ 0 & ; \text{ otherwise} \end{cases}$$

(b)

| $x$      | $p(x)$           | Exp. freq.<br>$N \cdot p(x)$ | Obs. freq. |
|----------|------------------|------------------------------|------------|
| 0        |                  |                              |            |
| 1        |                  |                              |            |
| 2        |                  |                              |            |
| 3        |                  |                              |            |
| 4        |                  |                              |            |
| 5        |                  |                              |            |
| 6        |                  |                              |            |
| $\geq 7$ | (by subtraction) |                              |            |

Hence  $p(x+1) = \left(\frac{x+r}{x+1}\right) \hat{q} \cdot p(x)$ .

$$p(1) = r \cdot \hat{q} \cdot p(0)$$

$$\therefore p(0) = \hat{p}^r = \text{_____}$$

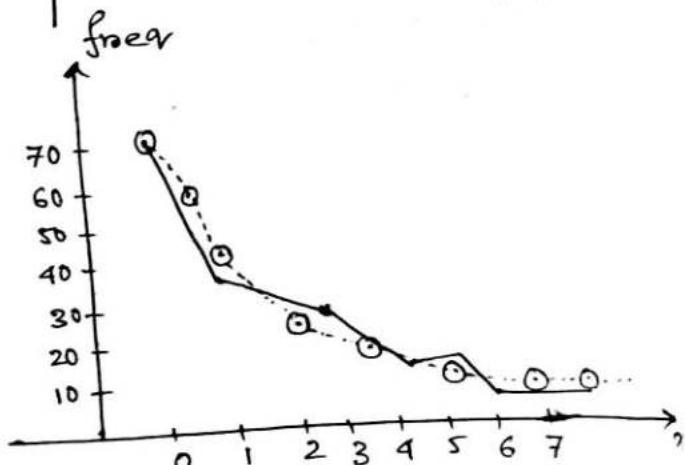
$$p(1) = \left(\frac{r}{1} \cdot \hat{q}\right) p(0) = \text{_____}$$

$$p(2) = \left(\frac{r+1}{2} \cdot \hat{q}\right) p(1) = \text{_____}$$

If the both frequency graphs are very near to each other throughout range, then the fitting is satisfactory.

Remark:- It is observed that the red mites on a leaf form a cluster, therefore it seems like that NB model is appropriate for  $X$ .

**Alt:-** For a poisson distribution, we should have  $\bar{x} \approx \delta^2$ . For a negative binomial distn., we should have  $\bar{x} \ll \delta^2$ .



██████████ fitted freq. graph

██████████ Observed freq. graph

(ii) The following table gives the freq. distn. of the number of albino children (in families of five children including at least one albino child) (Pearson's data)

| No. of albino in family | 1  | 2  | 3  | 4 | 5 | TOTAL |
|-------------------------|----|----|----|---|---|-------|
| No. of families         | 25 | 23 | 10 | 1 | 1 | 60    |

Fit an appropriate probability distribution to the data. Also, compare the observed and expected frequencies.

Solution:- Let  $X$  denotes the no. of albino children in families of five children. Then, considering "getting an albino children" as a success,  $X \sim \text{Bin}(n=5, p)$ .

$$\text{The PMF of } X \text{ is } p(x) = \begin{cases} \binom{5}{x} p^x q^{5-x}, & x=0(1)5 \\ 0 & \text{, otherwise} \end{cases}$$

Let  $Y$  denotes the no. of albino children in a family of five having at least one albino. The data is given on the R.V.  $Y$ .

Note that  $Y$  is the RV "  $X$  is truncated at  $x=0$ ".

$$\text{The PDF of } Y \text{ is } f(y) = \begin{cases} \frac{\binom{5}{y} p^y (1-p)^{5-y}}{1 - (1-p)^5}, & y=1(1)5 \\ 0 & \text{, otherwise} \end{cases}$$

By method of moments,  $E(Y) = \bar{y}$ ,

$$\Rightarrow \frac{5p}{1 - (1-p)^5} = \bar{y} \quad (*)$$

$$\bar{y} = \sum_{i=1}^5 i f_i / \sum_{i=1}^5 f_i = 1.833,$$

Solving (\*) by Iteration method:-

$$p = \frac{\bar{y}}{5} \left\{ 1 - (1-p)^5 \right\} = \phi(p), \text{ say}$$

$$\text{Trial root:- } p_0 = \frac{\bar{y}}{5} =$$

$$\text{Condition of Convergence:- } |\phi'(p_0)| = \text{_____} < 1.$$

Successive approximation,  $b_1 = f(b_0)$

$$\Phi_2 = \phi(\Phi_1)$$

14

$$\vec{p} = \underline{\hspace{2cm}}$$

∴ The fitted distribution is

$$f'(y) = \left\{ \begin{array}{l} \frac{\binom{s}{y} \hat{p}^y (1-\hat{p})^{s-y}}{1 - (1-\hat{p})^s}, y=1,2,\dots \\ 0, \text{ otherwise} \end{array} \right.$$

(iii) The table gives the frequency distribution of the number of dust nuclei in a small volume of air that fell onto a stage in a chamber containing moisture and filtered air. It is suspected that a number of zero counts were wrongly rejected on the ground that the apparatus was not working.

| No. of dust nuclei | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7 | 8 | Total |
|--------------------|----|----|----|----|----|----|----|---|---|-------|
| Frequency          | 23 | 56 | 86 | 95 | 73 | 40 | 17 | 5 | 3 | 400   |

fit an appropriate distribution to the frequency distribution, omitting the zero counts.

Solution:- Let  $X$  denotes the number of dust nuclei in a small volume of air.

Here, the poisson distribution is appropriate for  $X$ , let,  

$$X \sim P(\lambda).$$

Let,  $Y$  be the RV of  $X$  truncated at  $x=0$ .

Hence, the PMF of  $X$  is,

$$p(y) = \begin{cases} \frac{e^{-\lambda} \cdot \frac{\lambda^y}{y!}}{1 - e^{-\lambda}}, & \text{if } y = 1, 2, 3, \dots \\ 0, & \text{otherwise} \end{cases}$$

By method of moments,

$$E(Y) = \bar{y}$$
$$\Rightarrow \frac{\lambda}{1 - e^{-\lambda}} = \bar{y} = \frac{\sum_{i=1}^8 i \cdot f_i}{\sum_{i=1}^8 f_i} = \dots$$
$$\Rightarrow \lambda = \bar{y}(1 - e^{-\lambda}) = \phi(\lambda).$$

Trial root :-  $\lambda_0 = \bar{y} = \dots$

Condition of Convergence :-  $|\phi'(\lambda_0)| = \dots < 1.$

Successive improvements are :

$$\lambda_1 = \phi(\lambda_0)$$

$$\lambda_2 = \phi(\lambda_1)$$

$$\vdots$$
$$\hat{\lambda} = \dots, \text{ correct to two decimal places.}$$

Hence, fitted PMF of  $Y$  is,

$$p(y) = \begin{cases} e^{-\hat{\lambda}} \cdot \frac{(\hat{\lambda})^y}{y!}, & y = 1, 2, 3, \dots \\ 0, & \text{otherwise} \end{cases}$$

1. The life time ( $T$ ) in hours of an electron tube manufactured in a company is a r.v. with the d.f.

$$F(t) = 1 - e^{-t/\theta}, t > 0.$$

A sample of 337 electron tubes give the following frequency distribution of  $T$ .

|                  |      |        |         |         |         |         |         |         |
|------------------|------|--------|---------|---------|---------|---------|---------|---------|
| Life (in hours): | 0-50 | 50-100 | 100-150 | 150-200 | 200-250 | 250-300 | 300-350 | 350-400 |
| Frequency:       | 100  | 68     | 48      | 31      | 12      | 21      | 27      |         |

- (i) Estimate  $\theta$  from the data and compare the observed and expected frequencies.
- (ii) Estimate the probability that a supply of 20 tubes will not last more than 1900 hours, if they used one at a time successively.
- (iii) Find the number of electronic tubes likely to burn away within 80 hours of their life and also the average life of these tubes.

Solution:-

The DF of  $T$  is

$$F(t) = 1 - e^{-t/\theta}, t > 0$$

PDF of  $T$  is

$$f(t) = \frac{1}{\theta} e^{-t/\theta}, \text{ if } t > 0$$

- (i) By method of moments,

$$\begin{aligned} \mu'_1(T) &= \bar{t} \\ \Rightarrow \theta &= \bar{t} = \frac{\sum_{i=1}^7 t_i \cdot f_i}{\sum_{i=1}^7 f_i} = \hat{\theta} \end{aligned}$$

where,  $t_i, f_i$  are mid point and frequency of the  $i^{th}$  class.  
The fitted distribution is given by,

$$F(t) = 1 - e^{-t/\hat{\theta}}, t > 0.$$

$\therefore$  Expected frequency of the class interval  $(a, b)$  is

$$\begin{aligned} N.P[a < T < b] &= N \{ F(b) - F(a) \} \\ &= N \left\{ e^{-a/\hat{\theta}} - e^{-b/\hat{\theta}} \right\} \end{aligned}$$

## Computation of Expected Frequencies :-

| Class interval $(a, b)$ | $P[a < T < b]$ | Exp. Freq.<br>$N \cdot P[a < T < b]$ | Observed<br>Freq. |
|-------------------------|----------------|--------------------------------------|-------------------|
| 0 - 50                  |                |                                      |                   |
| 50 - 100                |                |                                      |                   |
| 100 - 150               |                |                                      |                   |
| ⋮                       |                |                                      |                   |
| ⋮                       |                |                                      |                   |

(ii) Let  $T_i$  denotes the life-time of the  $i^{th}$  tube,  $i = 1(1)20$ .  
 Here,  $T_i \stackrel{iid}{\sim} \text{Exp}(\theta)$ ,  $i = 1(1)20$ .

$$\text{Required Probability} = P\left[\sum_{i=1}^{20} T_i < 1900\right]$$

$$= P[T < 1900], T \sim \text{Gamma}(20, \theta).$$

$$= \int_0^{1900} \frac{e^{-t/\theta} \cdot t^{20-1}}{\theta^{20} \cdot \Gamma(20)} dt$$

$$= \int_0^{1900/\theta} \frac{e^{-x} \cdot x^{20-1}}{\Gamma(20)} dx \quad \left[ \because \frac{2T_i}{\theta} \stackrel{iid}{\sim} \chi^2_2, \right.$$

$$\Rightarrow \frac{2T}{\theta} = \sum_{i=1}^{20} \frac{2T_i}{\theta}$$

$$\Rightarrow \frac{2T}{\theta} \sim \chi^2_{40}$$

$$= \Gamma_{\frac{1900}{\theta}}(20)$$

$$\text{So, estimated probability} = \Gamma_{\frac{1900}{\theta}}(20).$$

[ Use Pearson Table for Incomplete Gamma ]

(iii) The no. of tubes which burns away within 80 hours.

$$= N \cdot P[T < 80]$$

$$= 337 \cdot \left\{ 1 - e^{-\frac{80}{\theta}} \right\}$$

$$= \underline{\quad}$$

The distribution of lifetime of a tube which burns away within 80 hours is

$$f'(t) = \begin{cases} \frac{1}{\theta} e^{-t/\theta}, & 0 < t < 80 \\ 0, & \text{otherwise} \end{cases}$$

$$E(T') = \int_0^{80} t \cdot \frac{1}{\theta} e^{-t/\theta} dt = \theta - 80 \left( \frac{e^{-80/\theta}}{1 - e^{-80/\theta}} \right).$$

$$\hat{E}(T') = \hat{\theta} - 80 \cdot \frac{e^{-80/\hat{\theta}}}{1 - e^{-80/\hat{\theta}}}.$$

5. (a) The following table gives the frequency distribution of the length of the left middle finger of 3000 criminals, obtained in the course of a study in clinical anthropometry.

|                                               |                                                      |
|-----------------------------------------------|------------------------------------------------------|
| Length (cm):<br>(Mid-point of class interval) | 9.5 9.6 10.1 10.4 10.7 10.0 11.3 11.6 11.9 12.2 12.5 |
| frequency :                                   | 1 4 24 67 193 417 575 691 509 306 131                |
| Length (cm):                                  | 12.8 13.1 13.4                                       |
| frequency :                                   | 63 16 3                                              |

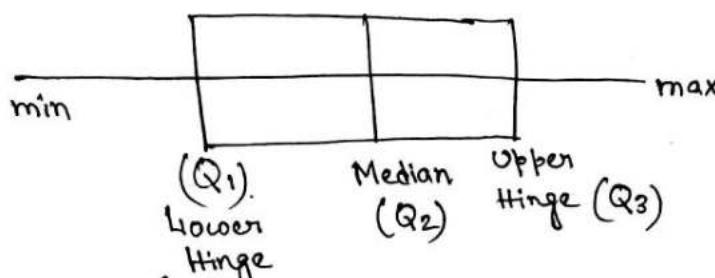
(i) Display the data graphically by means of box-plot.

(ii) Fit an appropriate probability distribution to the data.

(iii) Plot observed and estimated frequencies on a graph paper and comment on the fit.

Solution:-

(i)



(ii) Here the freq. distn. is more or less symmetric and the variable given continuous. Let  $X$  denotes the length of the left middle finger of a criminal.

Assume that  $X \sim N(\mu, \sigma^2)$ .  
The PDF of  $X$  is  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$ ;  $x \in \mathbb{R}$ .

By method of moments,

$$\mu_1' = m_1', \mu_2' = m_2.$$

$$\Rightarrow \begin{cases} \hat{\mu} = \bar{x} = \frac{\sum x_i f_i}{\sum f_i} = \dots \\ \hat{\sigma}^2 = s^2 = \frac{\sum x_i^2 f_i}{\sum f_i} = \dots = \bar{x}^2. \end{cases}$$

Hence, the fitted PDF of  $X$  is

$$f(x) = \frac{1}{\hat{\sigma}\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\hat{\mu}}{\hat{\sigma}})^2}; x \in \mathbb{R}.$$

| Cl. Mark | Class Boundaries | Exp. freq. of the class $(a, b)$ is<br>$= N.P[a < X < b] = N \left[ \Phi\left(\frac{b-\hat{\mu}}{\hat{\sigma}}\right) - \Phi\left(\frac{a-\hat{\mu}}{\hat{\sigma}}\right) \right]$ |
|----------|------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 9.5      | 9.35 - 9.65      |                                                                                                                                                                                    |
| 9.8      | 9.65 - 9.95      |                                                                                                                                                                                    |
| :        | :                |                                                                                                                                                                                    |

| Class boundaries | $\Phi\left(\frac{x-\hat{\mu}}{\hat{\sigma}}\right)$ | $\Phi\left(\frac{x-\hat{\mu}}{\hat{\sigma}}\right)$ | Exp. Freq<br>$N.A\Phi$ | Observed Freq. |
|------------------|-----------------------------------------------------|-----------------------------------------------------|------------------------|----------------|
| $-\infty$        | 0                                                   | -                                                   |                        |                |
| 9.35             | -                                                   |                                                     |                        |                |
| 9.65             | -                                                   |                                                     |                        |                |
| :                | :                                                   |                                                     |                        |                |
| 13.65            | -                                                   |                                                     |                        |                |
| $\infty$         | 1                                                   | ...                                                 |                        |                |

Comment:- If the graph of observed and expected frequencies are very close to each other then the fitting is good.

- (b) The following table gives the frequency distribution of annual salaries of individuals in rupees obtained from Indian Income Tax Returns (1995)

| <u>Statistics of Individual salaries Assessed, 1995</u> |                  |
|---------------------------------------------------------|------------------|
| <u>Annual Salary (rupees)</u>                           | <u>Frequency</u> |
| below 5000                                              | 27,000           |
| 5001 - 10000                                            | 60,000           |
| 10001 - 25000                                           | 90,000           |
| 25001 - 40000                                           | 26,032           |
| 40001 - 80000                                           | 13,000           |
| 80001 - 100000                                          | 60000            |
| 100001 - 200000                                         | 1178             |
| 200001 and above                                        | 312              |

Plot the data, using a suitable scale, and find out whether Pareto's law will give an adequate fit over the entire range of the observed distn., if not, try to fit the data in a suitable range.

Solution:- Let  $X$  denotes the salary of an individual.

Assume that the distribution of  $X$  is given. Then the DF of  $X$  is

$$F(x) = 1 - \left(\frac{x_0}{x}\right)^{\gamma} \text{ if } x > x_0$$

$$\Rightarrow 1 - F(x) = \left(\frac{x_0}{x}\right)^{\gamma}$$

$$\Rightarrow P[X \leq x] = \left(\frac{x_0}{x}\right)^{\gamma}$$

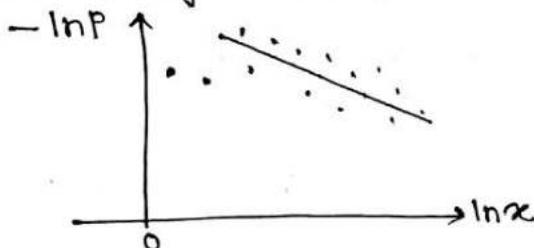
$$\Rightarrow \ln p = \ln x_0^{\gamma} - \gamma \ln x$$

$\Rightarrow \ln p = \alpha - \gamma \ln x$ , where,  $p$  is the proportion of individuals whose income is above  $x$ .

Note that,  $-\ln p = -\alpha + \gamma \ln x$ .

| Class | Freq. | $p = \frac{\text{c.f. of } \geq \text{ type}}{N}$ | $\ln p$ | $\ln x$ |
|-------|-------|---------------------------------------------------|---------|---------|
|       |       |                                                   |         |         |

$x$  = upper bound of the class.



From graph, we get — "Except the first two points, all the points are near about a straight line".  
Hence, the Pareto law is appropriate after including 3rd class.

Fitting of Pareto Distribution:— Here  $-\ln p = (-\alpha) + \ln x$   
 $\Rightarrow p' = -\alpha + \gamma x'$  ( $p' = \ln p$ ,  $x' = \ln x$ )

Hence  $(p', x')$  are linearly related.

The parameters  $\alpha, \gamma$  are estimated by method of least squares based on the values  $(p'_i, x'_i)$ , from 3rd class and above.

Normal Equations are —

$$\sum p'_i = N'(-\alpha) + \gamma \sum x'_i$$

$$\sum x'_i p'_i = -\alpha \sum x'_i + \gamma \sum x'^2_i$$

$$\Rightarrow \hat{\alpha} = \underline{\quad}, \hat{\gamma} = \underline{\quad}.$$

$$\text{Also } \hat{x}_0 = \underline{\quad}, \hat{\gamma} = \underline{\quad}.$$

Hence, the fitted Pareto law is given by the D.F.

$$F(x) = 1 - \left( \frac{\hat{x}_0}{x} \right)^{\hat{\gamma}}, \text{ if } x > x_0.$$

(c) The differences (in mins) of the actual arrival times and the scheduled arrival are tabulated below:

| Difference (in mins) | No. of days |
|----------------------|-------------|
| < 1                  | 12          |
| 4-7                  | 23          |
| 7-55                 | 141         |
| 55-90                | 14          |
| > 90                 | 11          |

- (i) fit a log-normal distribution to the data.  
(ii) Use the fitted distribution to estimate the probability that the train will arrive at 1hr, 15mins late on a typical day.

Solution:- Let  $X$  denotes the differences (in mins) of the actual arrival time and the scheduled arrival time.

(i) Here  $X \sim N(\mu, \sigma^2)$ .

Note that —  $P = P[X \leq e_{sp}] = P[\ln X \leq \ln e_{sp}]$

$$= P\left[\frac{\ln X - \mu}{\sigma} \leq \frac{\ln e_{sp} - \mu}{\sigma}\right]$$

$$= \Phi\left(\frac{\ln e_{sp} - \mu}{\sigma}\right)$$

$$= \Phi(z_p) \quad [\because \ln X \sim N(\mu, \sigma^2)]$$

$$\Rightarrow e_{sp} = e^{\mu + \sigma z_p}; \text{ where } \Phi(z_p) = p.$$

To protect the tails and to cover up the most of the values, we shall use  $P_{10}$  and  $P_{90}$ .

Now, equating the sample and population percentile, we get —

$$P_{10} = \hat{P}_{10} = 5.5 \Rightarrow e^{\mu + \sigma(z_{0.1})} = \hat{P}_{10}$$

$$P_{90} = \hat{P}_{90} = 72.5 \Rightarrow e^{\mu + \sigma(z_{0.9})} = \hat{P}_{90}$$

$$\Rightarrow \ln \hat{P}_{10} = \mu + \sigma(z_{0.1}) \quad \& \quad \ln \hat{P}_{90} = \mu + \sigma(z_{0.9})$$

$$\Rightarrow \hat{\mu} = \text{_____}, \quad \hat{\sigma} = \text{_____}.$$

(ii) Required Probability  $= P[74.5 < X < 75.5]$

$$= \Phi\left(\frac{\ln 74.5 - \hat{\mu}}{\hat{\sigma}}\right) - \Phi\left(\frac{\ln 75.5 - \hat{\mu}}{\hat{\sigma}}\right).$$

## PRACTICALS ON STATISTICAL INFERENCE I

### ESTIMATION

1.

The length of life recorded in hours for 10 electron tubes were:

980, 1020, 995, 1015, 990, 1030, 975, 950, 1050, 870

Assume that life-times are distributed in the form:

$$f(t) = \frac{1}{\theta} e^{-t/\theta}, \text{ if } t > 0, \text{ where } \theta > 0$$

- (i) Obtain an estimate of  $\theta$  and the estimated S.E. of this estimate.
- (ii) Estimate also the probability that an electron tube will survive at least 100 hours.
- (iii) Determine the lower confidence limit with confidence coefficient 0.95 to  $\theta$  and the probability of survival for 100 hours or more.

Solution:- (i) By method of moments :-

$$E(T) = \bar{t}$$

$$\Rightarrow \hat{\theta} = \bar{t} = \frac{1}{10} \sum_{i=1}^{10} t_i = \underline{\quad}$$

$$\text{S.E. of } \hat{\theta} = SE(\hat{\theta}) = \sqrt{\text{Var}(T)} = \sqrt{\frac{\theta^2}{n}} = \frac{\theta}{\sqrt{n}}$$

$$\therefore \text{Estimated S.E. is } \hat{SE}(\hat{\theta}) = \frac{\hat{\theta}}{\sqrt{n}} = \frac{\bar{t}}{\sqrt{10}} = \underline{\quad}$$

(ii) Probability that an electron tube will survive at least 100 hours

$$= P[T_i > 100]$$

$$= e^{-100/\hat{\theta}}$$

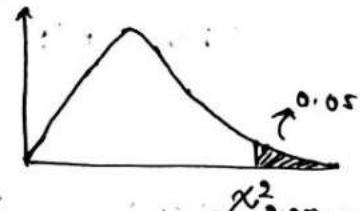
$$\therefore p = e^{-100/\hat{\theta}} = e^{-100/\bar{t}} = \underline{\quad}$$

(iii) Note that  $\frac{2T_i}{\theta} \sim \chi^2_2, i=1(1)10$

$$\Rightarrow 2 \sum_{i=1}^{10} \frac{T_i}{\theta} \sim \chi^2_{20}$$

$$\text{Now, } P\left[2 \sum_{i=1}^{10} \frac{T_i}{\theta} \leq \chi^2_{0.05, 20}\right] = 0.95$$

$$\Rightarrow P\left[\frac{20 \cdot \bar{t}}{\chi^2_{0.05, 20}} \leq \theta < \infty\right] = 0.95$$



Hence, the observed lower confidence limit is

$$\left(\frac{20\bar{t}}{\chi^2_{0.05, 20}}, \infty\right) = (-, \infty)$$

2. Let  $X_1, X_2, \dots, X_{10}$  be ten independent and identically distributed random variables where each  $X_i$  is normally distributed with mean 2.08. Further suppose 9.68% of its value is negative. If  $X_{(1)} < X_{(2)} < \dots < X_{(10)}$  be the ordered arrangement of  $X_i$ 's, then calculate  $P[X_{(3)} > 5.28]$ .

Solution:-

The pdf of  $X_{(n)}$  is

$$f_{X_{(n)}}(x) = \frac{1}{\binom{n}{n-1} \binom{n}{n-n}} \{F(x)\}^{n-1} \cdot f(x) \{1-F(x)\}^{n-n}$$

$$f_{X_{(3)}}(x) = \frac{1}{\binom{10}{2} \binom{7}{7}} \{F(x)\}^2 \{1-F(x)\}^7 f(x)$$

$$\text{and } P[X_{(3)} > 5.28] = \int_{5.28}^{\infty} \frac{1}{\beta(3,8)} \{F(x)\}^2 \{1-F(x)\}^7 f(x) dx$$

$$= \int_{F(5.28)}^1 \frac{1}{\beta(3,8)} z^2 (1-z)^7 dz$$

Now,  $X_i \stackrel{\text{iid}}{\sim} N(2.08, \sigma^2)$

$$\text{and } P[X_i < 0] = \frac{9.68}{100}$$

$$\Rightarrow P\left[\frac{X_i - 2.08}{\sigma} < \frac{0 - 2.08}{\sigma}\right] = 0.0968$$

$$\Rightarrow \Phi\left(-\frac{2.08}{\sigma}\right) = 0.0968 ; X_i \sim N(2.08, \sigma^2)$$

$$\Rightarrow \Phi\left(\frac{2.08}{\sigma}\right) = 0.9132$$

$$\Rightarrow \sigma = \underline{\hspace{2cm}}$$

$$\text{Then } F(5.28) = P[X_i < 5.28] = P\left[\frac{X_i - 2.08}{\sigma} < \frac{5.28 - 2.08}{\sigma}\right] \\ = \Phi\left(\frac{5.28 - 2.08}{\sigma}\right) = \underline{\hspace{2cm}} = x_0, \text{ say.}$$

$$\text{Hence, } P[X_{(3)} > 5.28] = 1 - P[X_{(3)} \leq 5.28]$$

$$= 1 - \int_0^{x_0} \frac{z^2 (1-z)^7}{\beta(3,8)} dz$$

$$= 1 - I_{x_0}(3,8) \quad [\text{Use Pearson table for incomplete beta}]$$

## \* TESTING OF HYPOTHESIS \*

1. Let  $X$  be equal to the thickness of spearmint gum manufactured for bending machine. Assume that the distn. of  $X$  is normal. The target thickness is 7.5 hundred of an age. The following 10 thickness of hundred of an age for pieces of gum that were selected randomly from the production line are

7.65, 7.60, 7.65, 7.70, 7.55, 7.40, 7.40, 7.50, 7.50, 7.55.

- i) At  $\alpha=0.05$  significance level, was the company successful to meet the target thickness?  
 ii) What is the appropriate p-value of your test?  
 iii) Is  $\mu=7.50$  contained in a 95% confidence interval for  $\mu$ ?

Solution:- Let  $X \sim N(\mu, \sigma^2)$

$$\text{From the sample, } \bar{x} = \frac{1}{10} \sum_{i=1}^{10} x_i = 7.55$$

$$\text{and } s^2 = \frac{1}{10} \sum_{i=1}^{10} x_i^2 - \bar{x}^2 = 0.097$$

(i)

To test  $H_0: \mu = 7.5$  vs.  $H_1: \mu \neq 7.5$

Test statistic is  $T = \frac{(\bar{x} - \mu_0) \sqrt{n}}{s} \sim t_{n-1}$ , under  $H_0$ .

$$\text{Here } T = \frac{(\bar{x} - 7.5) \sqrt{10}}{s} \sim t_9, \text{ under } H_0.$$

Critical region:- If the observed value of  $|T| > t_{\alpha/2, 9}$ ; we shall reject  $H_0$  at  $\alpha$  level of significance.

$$\text{i.e. } \left| \left( \frac{\bar{x} - 7.5}{s} \right) \sqrt{10} \right| > t_{0.025, 9} \text{ ; at } \alpha = 0.05$$

From table,  $t_{0.025, 9} = 2.262$ .

$$\text{Observed value, } \frac{(\bar{x} - 7.5) \sqrt{10}}{s} = \frac{(7.55 - 7.5) \sqrt{10}}{0.097} \\ = 1.63 < t_{0.025, 9} = 2.262$$

Hence, there is no reason to reject  $H_0$  at  $\alpha=0.05$  level of significance, i.e., the company was successful to meet the target value.

$$\begin{aligned}
 \text{(ii) } p\text{-value} &= P_{H_0}[|T| \geq |t_{\text{obs}}|] \\
 &= 2P_{H_0}[T \geq |t_{\text{obs}}|], t_{\text{obs}} \text{ is the observed value of } T. \\
 &= 2 \cdot P_{H_0}[T \geq 1.63] \quad [\text{From Biometrika table}] \\
 &= 2 \times 0.1 \\
 &= 0.2.
 \end{aligned}$$

As  $p$ -value is quite large, the observed value is a likely value under  $H_0$ .

$$\begin{aligned}
 \text{(iii) } 95\% \text{ confidence interval for } \mu \text{ is} \\
 (\bar{x} - t_{0.025, 9} \cdot \frac{s}{\sqrt{n}}, \bar{x} + t_{0.025, 9} \cdot \frac{s}{\sqrt{n}}) \\
 = \left( 7.55 - 2.62 \times \frac{0.097}{\sqrt{10}}, \quad \right) = (7.48, 7.62) \text{ when } \mu = 7.5
 \end{aligned}$$

2. The heart weights in groups of 12 females and 15 male cats are given below:

Male: 12.7, 15.6, 9.1, 12.1, 8.3, 11.2, 9.4, 8.0, 14.9, 10.7, 13.6, 9.6, 11.7, 9.3, 7.6

Female: 7.4, 7.3, 7.1, 9.0, 7.6, 9.5, 10.1, 10.2, 10.1, 9.5, 8.7, 7.2

Does the heart of a male cat on an average weight more than that of a female cat. State clearly any assumption you use.

Solutions:- Let  $X$  and  $Y$  denote the heart weights of a male and a female.

We assume that,  $X \sim N(\mu_1, \sigma_1^2)$  independently  
 $Y \sim N(\mu_2, \sigma_2^2)$

To test  $H_0: \mu_1 = \mu_2$  vs.  $H_1: \mu_1 > \mu_2$

We also assume that  $\sigma_1 = \sigma_2 = \sigma$  (unknown)

Let  $x_{11}, x_{12}, \dots, x_{1n}$  be a r.s. from  $N(\mu_1, \sigma^2)$

Let  $x_{21}, x_{22}, \dots, x_{2n}$  be a r.s. from  $N(\mu_2, \sigma^2)$ .

Here  $n_1 = 15, n_2 = 12$

$$\bar{x}_1 = \frac{1}{15} \sum_{i=1}^{15} x_{1i} = \underline{\hspace{2cm}}$$

$$\bar{x}_2 = \frac{1}{12} \sum_{i=1}^{12} x_{2i} = \underline{\hspace{2cm}}$$

$$s_1^2 = \frac{1}{n_1-1} \left\{ \sum x_{1i}^2 - n_1 \bar{x}_1^2 \right\} = \underline{\hspace{2cm}}$$

$$s_2^2 = \frac{1}{n_2-1} \left\{ \sum x_{2i}^2 - n_2 \bar{x}_2^2 \right\} = \underline{\hspace{2cm}}$$

We know,  $s^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}$  is an UE

Test statistic:-

$$T = \frac{\bar{x}_1 - \bar{x}_2}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1+n_2-2}, \text{ under } H_0.$$

Critical Region:-

$$\frac{(\bar{x}_1 - \bar{x}_2)}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} > t_{\alpha; n_1+n_2-2}$$

So, observed value of T is  $\frac{\bar{x}_1 - \bar{x}_2}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} =$

from table,  $t_{0.05, 15+12-2} = t_{0.05, 28} = 1.708.$

Conclusion:-

- 3) In an experiment to investigate the effect of light on root growth in mustard-sidlings, to growth of siblings, where grown in identical condition, except that one was kept in the dark, find the other was exposed in sunlight during the day. After a certain period of time, the root lengths in mm of all the sidlings are measured. The following table gives the data obtained performed appropriate statistical test of significance to assist whether
- i) light effects, root growth, & the variances of the root length in the two population to be equal but unknown.
  - ii) light effects in both growths in root lengths is the same.
- The table makes two rows first for light and second for darkness.

L: 21, 39, 31, 30, 52, 39, 55, 50, 29, 17

D: 22, 16, 20, 14, 32, 28, 36, 41, 17, 22

Solution:-

i) To test :-  $H_0: \mu_1 = \mu_2$  against  $H_1: \mu_1 > \mu_2$ .  
 If  $H_1$  is accepted, then light affects the root growth.  
 (same as previous problem)

ii) To test  $H_0: \sigma_1^2 = \sigma_2^2$  vs.  $H_1: \sigma_1^2 \neq \sigma_2^2$

Test statistic :-  $F = \frac{S_1^2}{S_2^2} \sim F_{n_1-1, n_2-1}$

Critical Region:-

$$\frac{s_1^2}{s_2^2} > F_{\alpha/2; n_1-1, n_2-1} \text{ or } < F_{1-\alpha/2; n_1-1, n_2-1}$$

where,  $n_1 = 10$ ,  $n_2 = 10$

Observed value of F is  $\frac{s_1^2}{s_2^2} = \frac{\dots}{\dots} = \dots$   
 $\alpha = 0.1$ .

From table,  $F_{0.05, 9, 9} = 3.18$ .

$$F_{0.95, 9, 9} = \frac{1}{F_{0.05, 9, 9}} = \frac{1}{3.18}$$

Conclusion:-

A) [CU '2006]

The following table contains observations on the supine systolic and diastolic blood pressures (in mm of Hg) for 15 patients with moderate hypertension, immediately before and two hours after taking a drug, captopril.

- Perform an appropriate test for the hypothesis that the drug is successful in reducing the diastolic blood pressure.
- Obtain a 95% confidence interval for the mean difference in the systolic blood pressure before and after treatment.

| Patient No. | Systolic BP |       | Diastolic BP |       |
|-------------|-------------|-------|--------------|-------|
|             | before      | After | before       | After |
| 1           | 210         | 201   | 130          | 125   |
| 2           | 169         | 165   | 122          | 121   |
| 3           | 187         | 166   | 124          | 121   |
| 4           | 160         | 157   | 104          | 106   |
| 5           | 167         | 147   | 112          | 101   |
| 6           | 176         | 145   | 101          | 85    |
| 7           | 185         | 168   | 121          | 98    |
| 8           | 206         | 180   | 124          | 105   |
| 9           | 173         | 147   | 115          | 103   |
| 10          | 146         | 136   | 102          | 98    |
| 11          | 174         | 151   | 98           | 90    |
| 12          | 201         | 168   | 119          | 98    |
| 13          | 198         | 179   | 106          | 110   |
| 14          | 148         | 129   | 107          | 103   |
| 15          | 154         | 131   | 100          | 82    |

Solution:- (i) Let  $X$  and  $Y$  denote the diastolic BP of a patient before and after taking the drug.

Assume that,  $(X, Y) \sim BN(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$ .

To test whether the drug is successful in reducing the diastolic BP, i.e. to test:

$$H_0: \mu_x = \mu_y \text{ vs. } H_1: \mu_x > \mu_y.$$

Let  $(x_i, y_i), i=1(1)15$ , be 15 paired samples on  $(X, Y)$ .

Define,  $D_i = x_i - y_i, i=1(1)n$ .

$$\Rightarrow D_i \stackrel{iid}{\sim} N(\mu_D, \sigma_D^2), \text{ where } \mu_D = \mu_x - \mu_y.$$

Here to test,  $H_0: \mu_D = 0$  vs.  $H_1: \mu_D > 0$ ,

Test statistic:-  $\frac{\sqrt{n}(\bar{D}-0)}{S_D} \sim t_{n-1}$

Critical region:-  $\frac{\sqrt{n}.\bar{d}}{S_D} > t_{\alpha/2, n-1}$

From the data:- for  $n=15$ ,

| Patient No. | 1 | 2 | 3 | ... | 15 |
|-------------|---|---|---|-----|----|
| $d$ :       |   |   |   |     |    |

$$\bar{d} = \frac{1}{15} \sum d_i = \underline{\hspace{2cm}}$$

$$S_d^2 = \frac{1}{n-1} \sum_{i=1}^n (d_i - \bar{d})^2 = \frac{1}{14} \left\{ \sum d_i^2 - 15\bar{d}^2 \right\} = \underline{\hspace{2cm}}$$

Observed value of the test statistic is  $\frac{\sqrt{n}.\bar{d}}{S_d} = \underline{\hspace{2cm}}$

Tabulated value:  $t_{0.05, 14} = \underline{\hspace{2cm}}$

Conclusion:-

(ii) Let  $U$  and  $V$  denote the systolic BP of a patient before and after taking drug.

Assume that,  $(U, V) \sim BN(\mu_U, \mu_V, \sigma_U^2, \sigma_V^2, \rho)$

Let  $(U_i, V_i), i=1(1)15$  be the 15 given paired sample on  $(U, V)$ .

Define,  $D_i = U_i - V_i \stackrel{iid}{\sim} N(\mu_D, \sigma_D^2), \mu_D = \mu_U - \mu_V$ .

$$\Rightarrow \frac{\sqrt{n}(\bar{D}-\mu_0)}{S_D} \sim t_{n-1}$$

Now, 95% C.I. for  $\mu_D$  is  $\left( \bar{d} - t_{0.05, 14} \cdot \frac{S_d}{\sqrt{15}}, \bar{d} + t_{0.05, 14} \cdot \frac{S_d}{\sqrt{15}} \right)$

Here . . . . .

| Patient No. | 1 | 2 | ... | 15 |
|-------------|---|---|-----|----|
| $d$ :       |   |   |     |    |

$$\bar{d} = \underline{\hspace{2cm}}$$

$$S_d = \underline{\hspace{2cm}}$$

Hence the observed 95% C.I. for  $\mu_D$  is

$$\left( \underline{\hspace{2cm}}, \underline{\hspace{2cm}} \right).$$

6) [CU'2005]

The members of a team of nine men were asked to load and fire naval guns by one of two methods  $M_1$  and  $M_2$ , attempting to get off as many rounds per minute as possible. Two sets of such firing were made per method. The following table gives the outcome of this experiment.

| Participate No. | Number of rounds fired per minute in |          |            |          |
|-----------------|--------------------------------------|----------|------------|----------|
|                 | first set                            |          | second set |          |
|                 | by $M_1$                             | by $M_2$ | by $M_1$   | by $M_2$ |
| 1               | 20.2                                 | 14.2     | 24.1       | 16.2     |
| 2               | 22.0                                 | 14.1     | 23.5       | 16.1     |
| 3               | 23.1                                 | 14.1     | 22.9       | 16.1     |
| 4               | 26.2                                 | 18.0     | 26.9       | 19.1     |
| 5               | 22.6                                 | 14.0     | 24.6       | 18.1     |
| 6               | 22.9                                 | 12.2     | 23.7       | 13.8     |
| 7               | 23.8                                 | 12.5     | 24.9       | 15.4     |
| 8               | 22.9                                 | 13.7     | 25.0       | 16.0     |
| 9               | 21.8                                 | 12.7     | 23.5       | 15.1     |

In the light of this data, are we justified in inferring that method  $M_1$  is significantly better than  $M_2$ ? Obtain a 95% C.I. for the difference in the mean performance by the two methods.

Solution:- Let  $(X_i, Y_i)$  be the paired samples in the 1<sup>st</sup> set by  $M_1$  and  $M_2$ ,  $i=1(1)9$ .

Let  $(X'_i, Y'_i)$  be the paired samples in the 2<sup>nd</sup> set by  $M_1$  and  $M_2$  for  $i=1(1)9$ .

We assume that,  $(X_i, Y_i) \stackrel{\text{iid}}{\sim} BN$ ,  $i=1(1)9$  > independent,

$(X'_i, Y'_i) \stackrel{\text{iid}}{\sim} BN$ ,  $i=1(1)9$

As  $(X_i, Y_i)$ ,  $(X'_i, Y'_i)$  are observed on the  $i^{\text{th}}$  gun man.

Note that,  $X_i \stackrel{\text{iid}}{\sim} \text{Normal}$ ,  $i=1(1)9$ .  
 $X'_i \stackrel{\text{iid}}{\sim} \text{Normal}$

$$\Rightarrow \frac{X_i + X'_i}{2} \stackrel{\text{iid}}{\sim} N(\mu_{M_1}, \sigma_{M_1}^2)$$

$$\text{Similarly, } \frac{Y_i + Y'_i}{2} \stackrel{\text{iid}}{\sim} N(\mu_{M_2}, \sigma_{M_2}^2).$$

But  $X_i + X_{i'}, Y_i + Y_{i'}$  are observed on the same  $i^{\text{th}}$  gun man,  
 $i=1(1)9.$

So,  $\left( \frac{X_i + X_{i'}}{2}, \frac{Y_i + Y_{i'}}{2} \right) \sim_{\text{iid}} \text{BN} \left( (\mu_{M_1}, \mu_{M_2}, \sigma_{M_1}^2, \sigma_{M_2}^2, \rho) \right)$

To test  $H_0: \mu_{M_1} = \mu_{M_2}$  vs.  $H_1: \mu_{M_1} > \mu_{M_2}$ .

Define,  $D_i = \frac{(X_i + X_{i'}) - (Y_i + Y_{i'})}{2} \sim_{\text{iid}} N(\mu_D, \sigma_D^2), i=1(1)9.$   
where,  $\mu_D = \mu_{M_1} - \mu_{M_2}.$

Now, To test,  $H_0: \mu_D = 0$  vs.  $H_1: \mu_D > 0$

(use Paired t-test).

6. (a) The correlation coefficient between head-length and stature for a sample of 36 members of an Indian tribe has been found to be 0.4339. Is it reasonable to assume that in the population of the characters are uncorrelated?

(b) Examine on the basis of the data given in the table below whether the variances of muscle weights of right and left legs of rabbits are equal.

| Sample no. of rabbit | Weight (gms) of anterior muscle of |           |
|----------------------|------------------------------------|-----------|
|                      | Left leg                           | Right leg |
| 1                    | 5.0                                | 4.9       |
| 2                    | 4.8                                | 5.0       |
| 3                    | 4.3                                | 4.3       |
| 4                    | 5.1                                | 5.3       |
| 5                    | 4.1                                | 4.1       |
| 6                    | 4.0                                | 4.0       |
| 7                    | 7.1                                | 6.9       |
| 8                    | 5.9                                | 6.3       |
| 9                    | 5.3                                | 5.2       |
| 10                   | 5.3                                | 5.5       |
| 11                   | 5.3                                | 5.9       |
| 12                   | 5.9                                | 6.8       |
| 13                   | 6.5                                | 6.3       |
| 14                   | 6.3                                | 6.6       |
| 15                   | 6.6                                | 6.3       |
| 16                   | 5.2                                | 6.3       |

(c) Given in the table are the means, the s.d.s and the correlation coefficient of scores  $x, y$  on two halves of a psychological test on 20 days.

(i) Examine whether the scores on the two halves are uncorrelated.  
(ii) Examine whether the scores on the two halves are equally variable.

| Score | Mean | S.d. | Correlation |
|-------|------|------|-------------|
| $x$   | 45.5 | 9.51 | 0.76        |
| $y$   | 50.2 | 5.25 |             |

(a) Let  $(X_i, Y_i), i=1(1)36$ , be the paired values on head-length and stature for 36 Indian tribes. Assume that  $(X_i, Y_i) \sim BN$  with correlation coefficient  $\rho$ .

To test  $H_0: \rho = 0$  vs.  $H_1: \rho \neq 0$

Test statistic:-

$$\frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t_{n-2}, \text{ under } H_0.$$

Critical region:-

$$\left| \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \right| > t_{\alpha/2, n-2}$$

Here  $n=36, r=0.4339$

Observed value of the test statistic is

$$\frac{r\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{0.4339 \times \sqrt{34}}{\sqrt{1-(0.4339)^2}} \\ = 2.808.$$

from Biometrika, Vol-I, for  $\alpha=0.05$

$$t_{0.025, 30} = 2.042$$

$$t_{0.025, 40} = 2.021$$

Here,  $\left| \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \right| = 2.808 > t_{0.025, 34}$

$\Rightarrow H_0$  is rejected at 5% level based on the given information.

$\Rightarrow$  The population characters are not uncorrelated.

(b) Let  $(X_i, Y_i)$  denote the weights of right and left legs of the  $i$ th rabbit,  $i=1(1)6$ .

Assume that  $(X_i, Y_i) \stackrel{iid}{\sim} BN(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$ ,  $i=1(1)n, n=16$ .

To test  $H_0: \sigma_x = \sigma_y$  vs.  $H_1: \sigma_x \neq \sigma_y$

Define,  $U_i = X_i + Y_i$

$V_i = X_i - Y_i$ ,  $i=1(1)16$ .

Note that  $\rho_{UV} = 0$ , under  $H_0$ .

To test,  $H_0: \sigma_x = \sigma_y$

$\Leftrightarrow H_0': \rho_{UV} = 0$

To test  $H_0: \rho_{UV} = 0$ , based on the paired sample  $(U_i; V_i)$ ,  
 $i=1(1)16$ ; compute

$$r_{UV} = \frac{\frac{1}{n} \sum U_i V_i - \bar{U}\bar{V}}{\sqrt{\left\{ \frac{1}{n} \sum U_i^2 - \bar{U}^2 \right\}} \sqrt{\left\{ \frac{1}{n} \sum V_i^2 - \bar{V}^2 \right\}}}$$

| Serial no | 1 | 2 | ... | 16 |
|-----------|---|---|-----|----|
| $U_i$     |   |   |     |    |
| $V_i$     |   |   |     |    |

Test statistic:

$$\frac{r_{UV} \sqrt{n-2}}{\sqrt{1-r_{UV}^2}} \sim t_{n-2}, \text{ under } H_0.$$

(c) (i) here we are to test whether  $\rho = 0$  or not.

(ii) To test  $\sigma_1 = \sigma_2$

If  $\rho = 0$ , then it becomes a testing of equality of variances in uncorrelated test.

If  $\rho \neq 0$ , it becomes a testing of equality of variances in correlated test.

To test  $\mu_1 = \mu_2$

If  $\rho = 0$ , &  $\sigma_1 = \sigma_2$  are accepted, then it becomes a Fisher's t-test.

If  $\rho \neq 0$  then it becomes a paired t-test.

7. (a) In an excavation 10 fossils were discovered, of which 6 could be definitely classified as male and 3 as a female. The sex of the last fossil could not be precisely determined. Are these findings compatible with a sex ratio?

Solution:- Let  $X$  denotes the no. of male fossils out of 9 fossils obtained.  
Assuming probability of getting a male fossil =  $p$ ,

$$X \sim \text{Bin}(9, p).$$

To test  $H_0: p = \frac{1}{2}$  against  $H_1: p \neq \frac{1}{2}$ .

Here the observed value of  $X$  is  $x_0 = 6$ .

$$\text{The p-value} = 2 \min \{ P_{H_0}[X \geq x_0], P_{H_0}[X \leq x_0] \}$$

$$= 2 \times \min \left\{ \sum_{x=6}^9 \binom{9}{x} \left(\frac{1}{2}\right)^9, \sum_{x=0}^6 \binom{9}{x} \left(\frac{1}{2}\right)^9 \right\}$$

$$= 2 \times \frac{\binom{9}{6} + \binom{9}{7} + \binom{9}{8}}{2^9}$$

$$= 0.50$$

Let  $\alpha = 0.05$  be the chosen level of significance.

Clearly, p-value  $> 0.05 = \alpha$ .

$\Rightarrow H_0$  is accepted.

$\Rightarrow$  The sex-ratio is 1:1 is supported by the data.

- (c) The following table gives the result of an experiment to compare the effect of newly discovered medicine on a certain disease with that of the prevailing treatment (control). Do the data confirm the superiority of the new drug?

| Treatment    | Cured | Not cured | Total |
|--------------|-------|-----------|-------|
| Control      | 3     | 5         | 8     |
| New medicine | 5     | 3         | 8     |
| Total        | 8     | 8         | 16    |

Solution:- Let  $X_1$  denotes the no. of cured out of 8 patients treated under control.

Let  $X_2$  " " " " " under new medicine.

Let  $X_1 \sim \text{Bin}(8, p_1)$  > independently ;  $n_1=8, n_2=8$   
 $X_2 \sim \text{Bin}(8, p_2)$

To test whether new medicine is superior.

$\Rightarrow$  To test  $H_0: p_1 = p_2$  vs.  $H_1: p_1 < p_2$ .

From data, the observed value of  $X_1$  and  $X_2$  are  $x_{10}=3, x_{20}=5$

Then  $X_1 + X_2 = X$  has the observed value  $x_0 = 8$ .

In the testing problem,

$$\text{p-value is} = P_{H_0} [X_1 \leq x_{10} \mid X = x_0]$$

$$= \sum_{x_1=0}^{x_{10}} \frac{\binom{n_1}{x_1} \binom{n_2}{x_0 - x_1}}{\binom{n_1+n_2}{x_0}} = \sum_{x_1=0}^3 \frac{\binom{8}{x_1} \binom{8}{8-x_1}}{\binom{16}{8}}$$

$$= \underline{\underline{\dots}}$$

Conclusion:-

- 8) (a) The number of deaths from drowning in a certain river in two consecutive months were 8 and 5. Are these fluctuations due to chance?
- (b) A newspaper in a certain city observed that driving conditions have much improved because the number of fatal automobile accidents in last year was 9 whereas the average number per year over the past several years was 15. Is the statement justified? If further, the number of fatal automobile accidents in the first six months of the current year is given to be 3, could you modify your earlier conclusion?
- (c) A 5 cm specimen of a new type of fibre is found to have 13 defects while the manufacturer claims that there are no more than 150 defects per 100 cm. Do the above data support this claim?

Solution:- (a) Let  $X_1$  and  $X_2$  denote the number of deaths of drowning in a certain river in two consecutive months.

Let  $X_1 \sim P(\lambda_1)$  independently  
 $X_2 \sim P(\lambda_2)$

To test whether the observed values are the values of the same population or not.

$\Leftrightarrow$  To test  $H_0: \lambda_1 = \lambda_2$  vs.  $H_1: \lambda_1 \neq \lambda_2$ .

Define,  $X = X_1 + X_2 \sim P(\lambda)$ , under  $H_0: \lambda_1 = \lambda_2 = \lambda$ .

The observed values of  $X_1$  and  $X$  are  $x_{10} = 8$ ,  $x_0 = 8+5 = 13$ .

$$\begin{aligned} \text{∴ p-value} &= 2 \times \min \left\{ P_{H_0}[X_1 \geq x_{10} | X = x_0], P_{H_0}[X_1 \leq x_{10} | X = x_0] \right\} \\ &= 2 \times \min \left\{ \sum_{x=x_{10}}^{x_0} \binom{x_0}{x} \frac{1}{2}^{\lambda}, \sum_{x=0}^{x_{10}} \binom{x_0}{x} \frac{1}{2}^{\lambda} \right\} \\ &= 2 \times \min \left\{ \sum_{x=8}^{13} \binom{13}{x} \frac{1}{2^{13}}, \sum_{x=0}^8 \binom{13}{x} \frac{1}{2^{13}} \right\} \end{aligned}$$

=

\_\_\_\_\_.

Conclusion:-

(b) Let  $X$  denotes the number of accidents per year.

We assume that  $X \sim P(\lambda)$

To test  $H_0: \lambda = 15$  vs.  $H_1: \lambda < 15$

Observed value of  $X$  is  $x_0 = 9$

$$p\text{-value} = P_{H_0}[X \leq x_0]$$

$$= \sum_{x=0}^9 e^{-15} \cdot \frac{15^x}{x!} = 1 - e^{-15} \sum_{x=0}^{15} \frac{15^x}{x!} = \underline{\quad}$$

Conclusion:- Let  $\alpha = 0.05$ , If p-value  $< \alpha = 0.05$   
then  $H_0$  is rejected.

■ Let  $X'$  denotes the no. of accidents per six month. Then

$$X' \sim Poi(\lambda/2)$$

The observed value of  $X'$  is  $x_0' = 3$ .

$$\therefore p\text{-value} = P_{H_0}[X' \leq x_0'] = \sum_{x=0}^3 e^{-7.5} \frac{(7.5)^x}{x!} = \underline{\quad}$$

Conclusion:-

(c) Let  $X$  denotes the no. of defects in a 5 cm specimen. We assume that  $X \sim P(\lambda)$ .

Manufacturer claims:

There are not more than 150 defects in 100 cm.

$\Rightarrow$  " " . . ; " 7.5 " " 5' cm.

To test  $H_0: \lambda_0 = 7.5$  vs.  $H_1: \lambda > 7.5$

The observed value of  $X$  is  $x_0 = 13$

$$\therefore p\text{-value} = P_{H_0}[X \geq 13]$$

$$= \sum_{x=13}^{\infty} e^{-7.5} \frac{(7.5)^x}{x!}$$

$$= 1 - \sum_{x=0}^{12} e^{-7.5} \cdot \frac{(7.5)^x}{x!}$$

$$= \underline{\quad}$$

9. For 20 pairs of fathers and sons, the regression equation of height of son ( $y$ ) on height of father ( $x$ ), both measured in cm, was found to be

$$y = 9.29 + 0.932x$$

For 20 pairs,  $\bar{x} = 168.17$ ,  $\sum(x_i - \bar{x})^2 = 777.80$  and  $\sum(y_i - \bar{y})^2 = 939.42$ . Test whether the regression coefficient differs significantly from unity. Find the 95% confidence limits to the conditional mean of  $y$  given  $x = 177$ . Also, the prediction limits of the height of a son when the height of the father is known to be 177 cm.

Solution:- Let  $\eta_x = a + bx$  be the regression equation of  $y$  on  $x$ . The least square linear regression equation of  $y$  on  $x$  is

$$\begin{aligned} y &= a + bx \\ \text{i.e. } y &= 9.29 + 0.932x \end{aligned}$$

To test  $H_0: \beta = 1$  vs.  $H_1: \beta \neq 1$

Test statistic:-

$$\frac{(b-1)\sqrt{S_{xx}}}{S_{y,x}} \sim t_{n-2}, \text{ under } H_0.$$

From data,  $b = 0.932$ ;  $n = 20$ .

$$S_{xx} = \sum(x_i - \bar{x})^2 =$$

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

$$= \frac{\sum(y_i - \bar{y})^2 - b^2 \sum(x_i - \bar{x})^2}{n-2}$$

$$= \frac{S_{yy} - b^2 S_{xx}}{n-2} = \underline{\quad}$$

If the observed

$$\left| \frac{(b-1)\sqrt{S_{xx}}}{S_{y,x}} \right| > t_{\alpha/2, n-2}$$

We shall reject  $H_0$  at  $\alpha$ -level.

Now, observed  
and from table,  
Conclusion:-

$$\left| \frac{(b-1)\sqrt{s_{xx}}}{s_{yx}} \right| = \text{_____}$$

for  $\alpha = 0.05$ ,  $t_{0.025, 18} = \text{_____}$

- To find 95% CI, for  $\eta_x = a + bx$ , when  $x = 177$ .

Let  $\hat{Y}_x = a + bx$  be the predictor value when  $x$  is given.

$$\frac{\hat{Y}_x - \eta_x}{s_{yx} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}} \sim t_{n-2}. \quad [ \begin{array}{l} \eta_x = a + bx = \\ \bar{x} = 168.17 \end{array} ]$$

The 95% CI for  $\eta_x$  is  $(\hat{Y}_x - t_{0.025, 18} \cdot s_{yx} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}, \hat{Y}_x + t_{0.025, 18} \cdot s_{yx} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}})$ ,

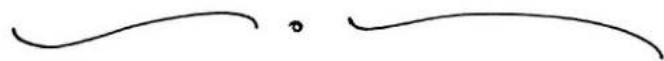
where,  $n = 20$ ,  $x = 177$ ,

- To find 95% prediction limit for  $y$  when  $x = 177$

$$\frac{y - \hat{Y}_x}{s_{yx} \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_{xx}}}} \sim t_{n-2}.$$

The 95% prediction limit for  $y$  when  $x = 177$  is

$$(\hat{Y}_{177} \pm t_{0.025, 18} \cdot s_{yx} \sqrt{1 + \frac{1}{20} + \frac{(177 - \bar{x})^2}{s_{xx}}}).$$



## Problems on Large Sample [only cu problems]

1. In a sample of size 100 from a bivariate popn., the correlation coefficient is found to be 0.25. Test whether it is  
(i) significant (ii) significantly less than 0.5.

Solution:-

Here  $n = 100$

$$r = 0.25$$

(i) To test  $H_0: \rho = 0$  vs  $H_1: \rho > 0$ .

$$\left[ \text{exact test: } \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} \sim t_{n-2} \right]$$

Using Fisher's Z-transformation,

$$\sqrt{n-3} (Z - \xi_0) \sim N(0, 1)$$

$$\text{Here } Z = \frac{1}{2} \log \left( \frac{1+r}{1-r} \right) = \dots$$

$$\xi_{f_0} = \frac{1}{2} \log \left( \frac{1+\rho}{1-\rho} \right) = 0.$$

$\therefore$  Under  $H_0$ ,

$$T = \sqrt{n-3} (Z - 0) \sim N(0, 1).$$

Critical region:- Observed  $T > T_\alpha$ .

(ii) To test  $H_0: \rho = 0.5$  vs  $H_1: \rho < 0.5$

$$\text{Here } \xi_{f_0} = \frac{1}{2} \log \left( \frac{1+0.5}{1-0.5} \right) =$$

Under  $H_0$ ,

$$T = \sqrt{n-3} (Z - \xi_{f_0}) \sim N(0, 1).$$

Critical region:- Observed  $T < -T_\alpha$ .

2. To examine a manufacturer's claim that not more than 1.9% of the products are defectives, then 830 items are put to inspection and 25 found defected can be considered the claim to be justified.

Sol. To  $H_0: p = 0.019$  vs.  $H_1: p > 0.019$

Here  $n = 830$

$$\hat{p} = \frac{25}{830}$$

Using  $\sin^{-1}$  transformation,

under  $H_0$ ,

$$Y = \sqrt{n} \left( \sin^{-1} \sqrt{\hat{p}} - \sin^{-1} \sqrt{0.019} \right) \sim N(0, 1)$$

Critical region: observed  $Y > Z_{\alpha}$ .

3. A proponent of innovative teaching methods wishes to compare the effectiveness of teaching English by the traditional classroom lecture system and by the extensive use of audio-visual aid. To do so, 100 students are selected at random from a class of 250 and assigned to audio-visual instruction, the remaining 150 students are taught English in classroom lecture. At the end of the term, all 250 students are given a test; the no. of students from each group to pass the test is recorded in the following table.

| Medium                   | Pass | Fail |
|--------------------------|------|------|
| Audio-visual instruction | 63   | 37   |
| classroom lecture        | 107  | 43   |

- (a) Find the 95% C.I. for the difference between success rate for the two methods of instructions.  
 (b) Do the data strongly support that a better passing rate is achieved using the classroom lecture than is achieved using A.V. method.  
 (c) Explain whether your inferential procedure crucial depends on the term 'random'.

Sol. Let  $p_1$  and  $p_2$  be the passing rates in Audio-visual and classroom lecture.

$$\text{Here } \hat{p}_1 = \frac{63}{100}, \hat{p}_2 = \frac{107}{150}.$$

$$(a) \text{ Here } P \left[ \sin^2 \left( \sin^{-1} \sqrt{\hat{p}_i} - \frac{Z_{\alpha/2}}{\sqrt{4n_i}} \right) \leq p_i \leq \sin^2 \left( \sin^{-1} \sqrt{\hat{p}_i} + \frac{Z_{\alpha/2}}{\sqrt{4n_i}} \right) \right] \\ = 1 - \alpha, i=1,2.$$

$$\Rightarrow P [L_i \leq p_i \leq U_i] = 1 - \alpha, i=1,2.$$

$$\text{Then } P [L_1 \leq p_1 \leq U_1, L_2 \leq p_2 \leq U_2]$$

$$= P[A_1 \cap A_2] \geq P[A_1] + P[A_2] - 1 = 1 - \alpha/2 + 1 - \alpha/2 - 1 \\ = 1 - \alpha.$$

$$\Rightarrow P[L_1 - U_2 \leq p_1 - p_2 \leq U_1 - L_2] > 1 - \alpha,$$

Here  $L_1 =$

$$L_2 =$$

$$U_1 =$$

$$U_2 =$$

and  $\alpha = 0.05$ .

Hence  $(L_1 - U_2, U_1 - L_2) = ($       )

is an observed C.I for  $(p_1 - p_2)$  with confidence level 0.95.

- (b) To test  $H_0: p_1 = p_2$  against  $H_1: p_1 < p_2$

[ Here  $\sqrt{4n_i} (\sin^{-1} \sqrt{\hat{p}_i} - \sin^{-1} \sqrt{p_i}) \sim N(0, 1)$ ,  $i=1, 2$ , independently  
as the groups are selected randomly.]

Here,  $\sin^{-1} \sqrt{\hat{p}_i} \sim N(\sin^{-1} \sqrt{p_i}, \frac{1}{4n_i})$ ,  $i=1, 2$  independently.

Now,  $\sin^{-1} \sqrt{\hat{p}_1} - \sin^{-1} \sqrt{\hat{p}_2} \sim N(\sin^{-1} \sqrt{\hat{p}_1} - \sin^{-1} \sqrt{\hat{p}_2}, \frac{1}{4n_1} + \frac{1}{4n_2})$

Under  $H_0$ ,  $\gamma = \frac{\sin^{-1} \sqrt{\hat{p}_1} - \sin^{-1} \sqrt{\hat{p}_2}}{\sqrt{\frac{1}{4n_1} + \frac{1}{4n_2}}} \sim N(0, 1) ..$

Critical region:- Observed  $T < -T_\alpha$ .

- (c) As the groups are selected randomly, the data obtained from the two groups then constitute two independent random samples. Therefore the estimates  $\hat{p}_1$  and  $\hat{p}_2$  are independently distributed and accordingly we obtain the test statistic.

4. For 600 beans of particular variety, the frequency distn. of breadth (in mm) has  $\gamma_1 = -0.093$ ,  $\gamma_2 = -0.125$ , where  $\gamma_1$  and  $\gamma_2$  are sample measures of skewness and kurtosis. Examine if the popln. can be supposed to be normal.

Sol. Here to test  $H_0$ : the data is a r.s. from normal popln.. As the measures skewness  $\gamma_1$  and kurtosis  $\gamma_2$  are given, then  $H_0$  reduces to  $H_{01}: \gamma_1 = 0$  and  $H_{02}: \gamma_2 = 0$ . as far as given information is concerned.

If both  $H_{01}$  and  $H_{02}$  are accepted, then we accept  $H_0$ .

Under  $H_0$ ,  $\gamma_1 \sim N(0, \frac{6}{n})$ ,  $\gamma_2 \sim N(0, \frac{24}{n})$ .

Test statistic:  $T_1 = \sqrt{\frac{n}{6}} \gamma_1 \sim N(0, 1)$

$T_2 = \sqrt{\frac{n}{24}} \gamma_2 \sim N(0, 1)$ .

If observed  $|T_1| > T_{\alpha/2}$ , we reject  $H_{01}$ .

If observed  $|T_2| > T_{\alpha/2}$ , we reject  $H_{02}$ .

$$[\text{Prob. of accepting } H_0 = P[|T_1| < T_{\alpha/2}, |T_2| < T_{\alpha/2}] \\ \geq (1-\alpha)(1-\alpha) = 1-2\alpha.$$

$\Rightarrow$  Prob. of rejecting  $H_0 \leq 2\alpha.$  ]

[ Calculation & Conclusion: suggestion: use  $\alpha = 1\%$ . ]

5. The corr. coeff. between stature (cm) and nasal height (cm) of a group of 106 males is 0.672 and that for a female group of 117 females is 0.725. Can you treat the corr. coefficients to differ significantly.

Sol. Let  $\rho_1, \rho_2$  be the correlation coefficient between stature and nasal height of male and female, respectively.

Assuming the popln.s are bivariate.

To test  $H_0: \rho_1 = \rho_2$  vs.  $H_1: \rho_1 \neq \rho_2$ .

$$\text{Here } Z_i = \frac{1}{2} \log \left( \frac{1+\rho_i}{1-\rho_i} \right)$$

$$\epsilon_i = \frac{1}{2} \log \left( \frac{1+\rho_i}{1-\rho_i} \right)$$

$$\text{Here, } Z_i \sim N \left( \epsilon_i + \frac{\rho_i}{2(n-1)}, \frac{1}{n_i-3} \right), i=1, 2, \text{ independently.}$$

Under  $H_0: \rho_1 = \rho_2 = \rho$

$$Z_1 - Z_2 \sim N \left( 0, \frac{1}{n_1-3} + \frac{1}{n_2-3} \right).$$

$$\text{where } E(Z_1 - Z_2) = \frac{1}{2} \left\{ \frac{1}{n_1-1} - \frac{1}{n_2-1} \right\} = \frac{1}{2} \cdot \frac{n_1-n_2}{(n_1-1)(n_2-1)} \approx 0.$$

Test statistic:-

$$T = \frac{(Z_1 - Z_2)}{\sqrt{\frac{1}{n_1-3} + \frac{1}{n_2-3}}} \sim N(0, 1), \text{ under } H_0.$$

6. Consider the following set  $A_1 = \{x; -\infty < x \leq 0\}$

$$A_i = \{x; i-2 \leq x \leq i-1\}, i=2(1)7.$$

$$A_8 = \{x; 6 \leq x < \infty\}$$

A certain Hypothesis  $H_0$  assigns probabilities  $p_{10}$  to this sets

$A_i$  is accordance with

$$p_{10} = \int \frac{1}{2\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-3}{2}\right)^2} dx,$$

the hypothesis  $H_0$  can be tested on the basis of the observed frequencies of the sets  $A_i, i=2(1)8$ , which are respectively 60, 96, 140, 210, 172, 160, 88, 74.

Would  $H_0$  is accepted at the 5% level of significance?

Sol. To test  $H_0: p_{10} = \int_{A_1} \frac{1}{2\sqrt{2\pi}} e^{-\frac{1}{8}(x-3)^2} dx, i=1(1)6.$

$$\Leftrightarrow H_0: x \sim N(3, 2^2).$$

$$H_0: p_{10} = P[X \leq 0] = \Phi\left(-\frac{3}{2}\right)$$

$$p_{20} = P[0 < X \leq 1] = \Phi(-1) - \Phi\left(-\frac{3}{2}\right)$$

$$p_{80} = P[6 < X < \infty] = 1 - \Phi\left(\frac{3}{2}\right)$$

| Class    | Observed freq.<br>( $O_i$ ) | Exp. Freq. $e_i = n p_{10}$ |
|----------|-----------------------------|-----------------------------|
| $A_1$    |                             |                             |
| $A_2$    |                             |                             |
| $\vdots$ |                             |                             |
| $A_8$    |                             |                             |

This is a test of Goodness of fit,  $n = \sum_{i=1}^8 f_i$ .

Test statistic

$$\chi^2 = \sum \frac{O_i^2}{e_i} - n \sim \chi^2_{8-1}, \text{ under } H_0.$$

Critical region:- Observed  $\chi^2 > \chi^2_{\alpha, 7}$ .

7. A survey of drivers was taken to see if they have been in an accident during the previous year, and if so was it the minor and major accident. The results are tabulated by age group

| Age      | Accident type |       |       |
|----------|---------------|-------|-------|
|          | None          | minor | major |
| Under 18 | 67            | 10    | 5     |
| 18 - 25  | 42            | 6     | 5     |
| 26 - 40  | 75            | 8     | 4     |
| 40 - 65  | 56            | 4     | 6     |
| Over 65  | 57            | 15    | 1     |

Use an appropriate test to check if the data suggest that the nature of the accidents depends on age.

Sol. Let A: Nature of Accident

B: Age

To test A and are independent.

The test statistic is:  $\chi^2 = n \left\{ \sum_{i=1}^k \sum_{j=1}^l \frac{f_{ij}^2}{f_{i0} f_{0j}} - 1 \right\} \sim \chi^2_{(n-1)(l-1)}$

Critical region:-

observed  $\chi^2 > \chi^2_{\alpha, (n-1)(l-1)}$ .

8. To compare the results of college 1 and 2 the following data on the final examination results have been obtained.

| College | Excellent | Results |          | Bad |
|---------|-----------|---------|----------|-----|
|         |           | Good    | Mediocre |     |
| 1       | 62        | 173     | 28       | 8   |
| 2       | 51        | 126     | 70       | 17  |

State the null hypothesis and test it to determine if the results vary between the colleges.

Sol.

$$\sum_{i=1}^k \frac{(f_{ii} - n_i p_{ii})^2}{n_i p_{ii}} + \sum_{i=1}^l \frac{(f_{i2} - n_2 p_{i2})^2}{n_2 p_{i2}} + \dots \\ = \sum_{j=1}^l \sum_{i=1}^k \frac{(f_{ij} - n_i p_{ij})^2}{n_i p_{ij}} \sim \chi^2_{(l-1)(k-1)}$$

Under  $H_0$ :  $\sum_{j=1}^l \sum_{i=1}^k \frac{(f_{ij} - n_j \hat{p}_{ij})^2}{n_j \hat{p}_{ij}}$ ;  $\hat{p}_{ij} = \frac{f_{ij}}{n}$ .

$$\chi^2 = n \left\{ \sum_{j=1}^l \sum_{i=1}^k \frac{f_{ij}^2}{n_j f_{ij}} - 1 \right\} \sim \chi^2_{(l-1)(k-1)}$$

Let A denotes the result of a college.

|                         |  | A                       |                         |                         |                        |
|-------------------------|--|-------------------------|-------------------------|-------------------------|------------------------|
|                         |  | A <sub>1</sub>          | A <sub>2</sub>          | A <sub>3</sub>          | A <sub>4</sub>         |
| College I<br>(sample I) |  | f <sub>11</sub><br>= 62 | f <sub>21</sub><br>= 73 | f <sub>31</sub><br>= 28 | f <sub>41</sub><br>= 8 |
| College II<br>sample II |  | f <sub>12</sub>         | f <sub>22</sub>         | f <sub>32</sub>         | f <sub>42</sub>        |
|                         |  | f <sub>10</sub>         | f <sub>20</sub>         | f <sub>30</sub>         | f <sub>40</sub>        |

$$n_1 = \underline{\quad}$$

$$n_2 = \underline{\quad}$$

$$f_{ij}, i=1(1)4, \\ j=1(1)2.$$

To test H<sub>0</sub>: the homogeneity of the result in the colleges.

Test statistic under H<sub>0</sub>,

$$\chi^2 = n \left[ \sum_{j=1}^2 \sum_{i=1}^4 \frac{f_{ij}^2}{n_j f_{i0}} - 1 \right] \sim \chi^2_{(2-1)(4-1)}$$

Critical region:- Observed  $\chi^2 > \chi^2_{\alpha, 3}$ .

9. The no. of occurrences of a word ~~in~~ in 48 essays per thousand 288 essays examined, written by Hamilton and 50 essays written by Maddisson. The following table is obtained.
- | Author   | No. of essays according to rate of occurrence of <del>in</del> words |          |
|----------|----------------------------------------------------------------------|----------|
|          | Upto 45                                                              | above 45 |
| Hamilton | 37                                                                   | 11       |
| Madison  | 48                                                                   | 2        |

Test whether the data occurrence of 'to' vary significantly between the two authors by applying the  $\chi^2$  test, and also by applying the formulae for exact probability.

Sol.  $\rightarrow$  Let A denotes the rate of occurrence of 'to'.  
Let B<sub>1</sub>, B<sub>2</sub> denote the Authors : Hamilton and Madison.

|                |        | A <sub>1</sub> | A <sub>2</sub> |
|----------------|--------|----------------|----------------|
| B <sub>1</sub> | a = 37 | c = 11         | 48             |
| B <sub>2</sub> | b = 48 | d = 2          | 50             |
|                | 85     | 13             | N = 98         |

To test H<sub>0</sub>: A and B are independent.

Applying Yates continuity correction to Pearsonian  $\chi^2$ , we get

$$\chi^2 = \frac{\{(ad-bc) - \frac{N}{2}\}^2}{(a+d)(c+d)(a+c)(b+d)} \sim \chi^2_1.$$

Critical region:- Observed  $\chi^2 > \chi^2_{\alpha, 1}$ .

Exact Prob. Test:- Under  $H_0$ , given marginals, the prob. of obtaining the cell frequency  $\begin{array}{|c|c|} \hline a & c \\ \hline b & d \\ \hline \end{array}$  is

$$p_d = \frac{(a+c)(b+d)}{\binom{N}{a+b}} = \frac{\underline{a+c} \underline{b+d} \underline{a+b} \underline{c+d}}{\underline{N} \underline{a} \underline{b} \underline{c} \underline{d}}$$

$p\text{-value} = p_2 + p_1 + p_0$  [since here  $d=2$ ]

$p\text{-value} < \alpha \Rightarrow \text{reject } H_0$ .

$$p_2 \rightarrow \begin{array}{r|rr} 37 & 11 & 48 \\ \hline 48 & 2 & 50 \\ \hline 85 & 13 & 98 \end{array} \quad p_1 \rightarrow \begin{array}{r|rr} 36 & 12 & 98 \\ \hline 49 & 1 & 50 \\ \hline 85 & 13 & 98 \end{array} \quad p_0 \rightarrow \begin{array}{r|rr} 35 & 13 & 98 \\ \hline 50 & 0 & 50 \\ \hline 85 & 13 & 98 \end{array}$$

10. The variance of the stature in cm for males in 8 different capital throughout the world are given below. Examine whether the variance differs from capital to capital. If not, find the pooled estimates of the variance.

| Capital | Sample size | Variance of the stature (in $\text{cm}^2$ ) |
|---------|-------------|---------------------------------------------|
| 1       | 299         | 38.8306                                     |
| 2       | 77          | 24.8924                                     |
| 3       | 131         | 39.5068                                     |
| 4       | 59          | 24.9685                                     |
| 5       | 124         | 24.3490                                     |
| 6       | 170         | 32.4325                                     |
| 7       | 337         | 39.2262                                     |
| 8       | 139         | 33.4936                                     |

Sol: Let  $\sigma_i$  be the popn. s.d. of the  $i^{\text{th}}$  capital.  
Assuming the popn. distributions are normal,  
 $\ln \sigma_i \sim N(\ln \bar{\sigma}_i, \frac{1}{2n_i})$ ,  $i=1(1)8$ .

$$\Rightarrow \sum_{i=1}^8 \sqrt{2n_i} (\ln \sigma_i - \ln \hat{\sigma})^2 \sim \chi^2_7, \text{ under } H_0: \sigma_1 = \sigma_2 = \dots = \sigma_8 = \bar{\sigma}.$$

$$\text{where, } \ln \hat{\sigma} = \frac{\sum 2n_i \ln \sigma_i}{\sum 2n_i} = \text{pooled estimate}$$

Critical region: Observed  $\chi^2 > \chi^2_{\alpha/2, 7}$ .

11. The correlation between height and weight was found to be 0.7, 0.8, 0.95 for samples of thousands object each for three different ethnic groups of children aged between 1 to 3 years. Test whether there is significant evidence for the dependence of the correlation on ethnicity. If not, find the pool estimate of the common correlation.

Sol. Let  $\rho_i$  be the correlation coefficient between height and weight of the  $i^{\text{th}}$  ethnic group,  $i=1, 2, 3$ .

Assume, the popn. distns are BN.

To test  $H_0: \rho_1 = \rho_2 = \rho_3$  ag.  $H_1: \text{not } H_0$ .

Here,  $\sqrt{n_i - 3} (Z_i - \rho_i) \sim N(0, 1)$ ,  $i=1, 2, 3$ , independently.

$$\text{where, } Z_i = \frac{1}{2} \log \left( \frac{1+n_i}{1-n_i} \right) =$$

$$\rho_i = \frac{1}{2} \log \left( \frac{1+\rho_i}{1-\rho_i} \right) =$$

Under  $H_0: \rho_1 = \rho_2 = \rho_3 = \rho$ ,

$$\sum_{i=1}^3 (n_i - 3)(Z_i - \hat{\rho})^2 \sim \chi^2_{3-1}, \text{ where } \hat{\rho} = \frac{\sum (n_i - 3)Z_i}{\sum (n_i - 3)} = \bar{Z}.$$

Critical region:- Observed  $\chi^2 > \chi^2_{\alpha/2}$

$\square$  If  $H_0: \rho_1 = \rho_2 = \rho_3 = \rho$ ,  
the pooled estimate of the common correlation  $\rho$  is given by

$$\begin{aligned} \hat{\rho} &= \bar{Z} \\ \Rightarrow \frac{1}{2} \log \left( \frac{1+\hat{\rho}}{1-\hat{\rho}} \right) &= \bar{Z} \\ \Rightarrow \hat{\rho} &= \frac{e^{2\bar{Z}} - 1}{e^{2\bar{Z}} + 1} = \underline{\hspace{2cm}}. \end{aligned}$$

12. sample of sizes 75, 125, 100, 150, 130 from 5 poisson popn. give the mean values as  $115/75, 197/125, 141/100, 237/150, 185/130$ , respectively. Do the popn. have the same mean value?

Solution:-

To test  $H_0: \lambda_1 = \lambda_2 = \lambda_3 = \lambda_4 = \lambda_5$

Under  $H_0: \lambda_i = \lambda$ ,  $i=1(1)5$ .

$$4n_i (\sqrt{\lambda_i} - \sqrt{\lambda})^2 \sim \chi^2_{5-1}.$$

$$\text{where } \hat{\lambda} = \frac{\sum 4n_i \sqrt{\lambda_i}}{\sum 4n_i}$$

## Problems on Statistical Inference

1. The following observations ( $y$ ) are drawn from  $N(\theta, \theta^k)$ . Obtain the MLE of  $\theta$  in each case  $k=0, 1$ , and  $2$ . Are they unique?

$20.2, 22.9, 23.3, 20.0, 19.4, 22.0, 22.1, 22.0, 21.9, 21.5, 19.7, 21.5, 20.9.$

Hint:- Let  $y_1, y_2, \dots, y_n$  be a n.s. from  $N(\theta, \theta^k)$ ,  $k=0, 1, 2$ .

$$[k=0] \quad \text{The MLE of } \theta \text{ is } \hat{\theta} = \bar{y} = \underline{\hspace{2cm}}.$$

$$[k=1] \quad N(\theta, \theta) \\ \hat{\theta} = -1 + \frac{\pm \sqrt{\frac{4}{n} \sum y_i^2}}{2} = \underline{\hspace{2cm}}$$

$$[k=2] \quad N(\theta, \theta^2), \theta \neq 0 \\ L(\theta | y_1, \dots, y_n) = \left( \frac{1}{2\theta^2 \pi} \right)^n e^{-\frac{1}{2\theta^2} \sum_{i=1}^n (y_i - \theta)^2}, \theta \neq 0$$

$$\frac{\partial \ln L}{\partial \theta}$$

$$0 = -\frac{n}{\theta} + \frac{\sum y_i^2}{\theta^3} - \frac{\sum y_i}{\theta^2}.$$

$$\Rightarrow \theta^2 + \theta \bar{y} - \frac{1}{n} \sum y_i^2 = 0$$

$$\Rightarrow \hat{\theta} = \frac{-\bar{y} \pm \sqrt{\bar{y}^2 + \frac{4}{n} \sum y_i^2}}{2} = \underline{\hspace{2cm}}, \underline{\hspace{2cm}}$$

[MLE is not unique]

2. In the general experiment on linseed, the observed frequencies for petal and stigma colour are given in the following table along with probabilities of the cell in terms of a parameter  $\theta$ .

| Stigma colour | Petal colour                            |                                        |
|---------------|-----------------------------------------|----------------------------------------|
|               | lilac                                   | deep lilac                             |
| White         | $357 \left( \frac{2+\theta}{4} \right)$ | $33 \left( \frac{1-\theta}{4} \right)$ |
| Purple        | $37 \left( \frac{1-\theta}{4} \right)$  | $94 \left( \frac{\theta}{4} \right)$   |

Estimate  $\theta$  by the method of MLE, [ ( )  $\rightarrow$  Probability ]

Also estimate the variance of this estimate of  $\theta$ .

Sol.

Likelihood function,

$$L(\theta | n_1, n_2, n_3, n_4) = \frac{1^n}{\prod_{i=1}^4 n_i} \left( \frac{2+\theta}{4} \right)^{n_1} \left( \frac{1-\theta}{4} \right)^{n_2} \left( \frac{1-\theta}{4} \right)^{n_3} \left( \frac{\theta}{4} \right)^{n_4}$$

Likelihood equation:

$$\theta = \frac{\partial \ln L}{\partial \theta} = \frac{n_1}{2+\theta} - \frac{n_2+n_3}{1-\theta} + \frac{n_4}{\theta}.$$

$$\Rightarrow n\theta^2 - \{n_1 - 2(n_2+n_3) - n_4\}\theta - 2n_4 = 0, \quad n = \sum_{i=1}^4 n_i$$

$$\Rightarrow \hat{\theta} = \underline{\hspace{2cm}}.$$

For large  $n$ ,  
 $\hat{\theta} \sim N\left(\theta, \frac{1}{I_n(\theta)}\right)$ , where  $I_n(\theta) = E\left(-\frac{\partial^2}{\partial \theta^2} \ln L\right)$

$$\begin{aligned} I_n(\theta) &= E\left\{ \frac{n_1}{(2+\theta)^2} + \frac{n_2+n_3}{(1-\theta)^2} + \frac{n_4}{\theta^2} \right\} \\ &= n \left\{ \frac{\frac{2+\theta}{4}}{(2+\theta)^2} + \frac{\frac{2(1-\theta)}{4}}{(1-\theta)^2} + \frac{\theta/4}{\theta^2} \right\} \\ &= \frac{n}{4} \left\{ \frac{1}{2+\theta} + \frac{2}{1-\theta} + \frac{1}{\theta} \right\} \end{aligned}$$

Asymptotic variance is  $\text{Var}(\hat{\theta}) = \frac{1}{I_n(\theta)}$ .

$$\text{Var}(\hat{\theta}) = \frac{1}{I_n(\theta)} = \frac{n}{4} \left\{ \frac{1}{2+\theta} + \frac{2}{1-\theta} + \frac{1}{\theta} \right\}$$

3) The length of life recorded in hours for 10 electron tubes were

980, 1020, 995, 1015, 990, 1030, 975, 950, 1050, 870.

Assume that life times are distributed in the form:

$$f(t, \theta) = \frac{1}{\theta} e^{-(t-\theta)}, \theta > 0, 0 < t < \infty.$$

Obtain the MLE of  $\theta$  and the estimated standard error of this estimate. Estimate also the probability that an electron tube will survive at least 100 hours. Given estimate of the large sample standard error of this estimated probability. Determine the lower confidence limit. Confidence coefficient = 0.05, to the true prob. of survival for 100 hours or more.

- (i) By using the exact distn of MLE of  $\theta$ ,
- (ii) Assuming some approximate distn for the estimated prob. of survivor.

Sol. The likelihood function is

$$L(\theta | t_1, t_2, \dots, t_n) = \frac{1}{\theta^n} e^{-\sum_{i=1}^n t_i/\theta}, \theta > 0.$$

Here  $n = 10$ ,

$$\text{Likelihood equation is } 0 = \frac{\partial}{\partial \theta} \ln L = -\frac{n}{\theta} + \frac{\sum t_i}{\theta^2}$$

$$\Rightarrow \hat{\theta} = \bar{t} = \text{_____} \text{ is the MLE of } \theta.$$

$$\text{Now, S.E.}(\hat{\theta}) = \sqrt{V(\bar{t})} = \sqrt{\frac{\theta^2}{n}} = \frac{\theta}{\sqrt{n}}.$$

$$\text{and } \text{S.E.}(\hat{\theta}) = \frac{\hat{\theta}}{\sqrt{n}} = \frac{\bar{t}}{\sqrt{n}} = \text{_____}.$$

$$\text{To estimate } p = P[T > 100] = e^{-100/\theta}$$

$$\text{MLE of } p \text{ is } \hat{p} = e^{-100/\hat{\theta}} = \text{_____}.$$

For large sample,

$$\hat{\theta} \sim N\left(\theta, \frac{1}{I_n(\theta)}\right), \text{ where } I_n(\theta) = \frac{n}{\theta^2}.$$
$$\Rightarrow \hat{p} = e^{-100/\theta} \sim N\left(p, \frac{\left\{ \frac{\partial}{\partial \theta} p \right\}^2}{I_n(\theta)}\right)$$

For large  $n$ ,

$$\text{Var}(\hat{p}) = \frac{\left\{ p \cdot \frac{100}{\theta^2} \right\}^2}{\frac{n}{\theta^2}}$$

$$\Rightarrow S.E.(\hat{p}) = \frac{p \cdot \frac{100}{\theta^2}}{\sqrt{\frac{n}{\theta^2}}} = \frac{100p}{\sqrt{n}\theta}.$$

$$\text{Hence, } SE(\hat{p}) \approx \frac{100p}{\hat{\theta}\sqrt{n}} = \text{_____}$$

Confidence limit of  $p$ :-

(i) Exact C.I. :-  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n t_i$

$$\text{Here } y_i = \frac{2t_i}{\theta} \stackrel{iid}{\sim} \chi^2_2, i=1(1)n.$$

$$\Rightarrow 2 \frac{\sum t_i}{\theta} \sim \chi^2_{2n}$$

$$\Rightarrow \frac{2n\hat{\theta}}{\theta} \sim \chi^2_{2n}$$

$$\therefore P\left[\frac{2n\hat{\theta}}{\theta} < \chi^2_{\alpha, 2n}\right] = 1-\alpha$$

$$\Rightarrow P\left[\theta > \frac{2n\hat{\theta}}{\chi^2_{\alpha, 2n}}\right] = 1-\alpha,$$

$$\Rightarrow P\left[e^{-100/\theta} > e^{-\frac{\chi^2_{\alpha, 2n}}{2n\hat{\theta}} \cdot 100}\right] = 1-\alpha,$$

$$\Rightarrow P\left[p > e^{-100 \cdot \frac{\chi^2_{\alpha, 2n}}{2n\hat{\theta}}}\right] = 1-\alpha,$$

Lower confidence limit with confidence coefficient 0.95 is

$$p_L = e^{-\frac{100 \cdot \chi^2_{0.05, 20}}{20 \cdot 2 \cdot \hat{\theta}}} = \text{_____}$$

(ii) Approximate confidence limit :-

$$\hat{p} \sim N\left(p, S.E^2(\hat{p})\right)$$

$$\Rightarrow \frac{\hat{p}-p}{S.E(\hat{p})} \sim N(0,1).$$

Hence for large  $n$ ,

$$P\left[\left| \frac{\hat{p}-p}{S.E(\hat{p})} \right| < T_{\alpha/2}\right] = 1-\alpha.$$

$$\Rightarrow P\left[\hat{p} - T_{\alpha/2} \cdot S.E(\hat{p}) < p < \hat{p} + T_{\alpha/2} \cdot S.E(\hat{p})\right] = 1-\alpha.$$

- 4) At time  $t=0$ , 20 identical components are put on test. The life distn. of each is exponential with mean  $\theta$ , after 24 hours. We were found that 15 of the 20 components are still working. Derive the MLE of  $\theta$ . Also give an estimate of SE. of the estimator.

Hint:- Here  $T \sim \text{Exp}$  with mean  $\theta$   
 Let  $Y =$  the no. of bulbs survived upto 24 hours out of 20 bulbs.  
 clearly,  $Y \sim \text{Bin}(n=20, p)$ ; where  $p = P[T > 24] = e^{-24/\theta}$ ,  
 MLE of  $p$  is  $\hat{p} = \frac{y}{n} = \frac{15}{20} = 0.75$

$$\Rightarrow e^{-24/\hat{\theta}} = \hat{p}$$

$$\Rightarrow \hat{\theta} = \frac{-24}{\ln(0.75)} = \dots$$

- 5) In a life testing experiment, 10 electric lamps are put to test. The lamps are burnt at a stretch for 20 hours. 2 of them survive time to termination (in hours) of life, for the remaining lamps, are observed and given below

9.8, 15.8, 17.2, 11.2, 13.8, 18.9, 14.8, 19.6

find an estimate of the mean life  $\theta$  assuming the life distn. to be exponential. If all the bulbs survive, what could have been your estimate of  $\theta$ .

Sol.  $T \sim \text{Exp. with mean } \theta$ .

$$p = P[T > 20] = e^{-20/\theta}$$

Let  $x_1, \dots, x_8$  be the lifetime of 8 lamps.

$$\begin{aligned} \text{Likelihood function is } L(\theta) &= \prod_{i=1}^8 \left\{ \frac{1}{\theta} \cdot e^{-x_i/\theta} \right\} \cdot p^2 \\ &= \frac{1}{\theta^8} \cdot e^{-\sum_{i=1}^8 x_i/\theta} \cdot e^{-40/\theta} \\ &= \frac{1}{\theta^8} \cdot e^{-\frac{8(\bar{x}+5)}{\theta}} \end{aligned}$$

$$\text{Likelihood equation is: } 0 = \frac{\partial}{\partial \theta} \ln L(\theta) = -\frac{8}{\theta} + \frac{8(\bar{x}+5)}{\theta^2}$$

$$\Rightarrow \hat{\theta} = \bar{x} + 5 = \dots$$

If all 10 bulbs survived, the likelihood equation is

$$L(\theta) = \left\{ e^{-20/\theta} \right\}^{10} = e^{-200/\theta}$$

$L(\theta)$  is max. iff  $200/\theta$  is minimum.

iff  $\theta$  is max.

Hence, MLE does not exist.

6) The following in the freq. distn. of 154 obs'n drawn at random from a multinomial popn. with 6 classes

| Class no                         | 1  | 2  | 3 | 4 | 5  | 6  |
|----------------------------------|----|----|---|---|----|----|
| No. of observation in the sample | 79 | 32 | 5 | 6 | 17 | 15 |

denoting the probabilities of 6 classes by  $\pi_1, \pi_2, \dots, \pi_6$ , respectively, find the MLE of  $\pi_1 - 2\pi_2 + 6\pi_5$ .

Sol. → Let  $f_i, i=1(1)6$  be the frequency of the  $i^{\text{th}}$  class in a n.s. of size  $n=154$ .

Likelihood function:-

$$L(\pi_i | f_i) = \frac{1}{\prod_{i=1}^n f_i!} \cdot \prod_{i=1}^6 \pi_i^{f_i}, \text{ where } \sum_{i=1}^6 \pi_i = 1.$$

and  $0 < \pi_i < 1, i=1(1)6$ .

To maximize  $\ln L$  subject to  $\sum_{i=1}^6 \pi_i = 1$ .

$$\text{let, } F = \ln L + \lambda \left( \sum_{i=1}^6 \pi_i - 1 \right)$$

$$= \text{constant} + \sum_{i=1}^6 f_i \cdot \ln \pi_i + \lambda \left( \sum_{i=1}^6 \pi_i - 1 \right)$$

$$\text{solve: } 0 = \frac{\partial F}{\partial \pi_i} = \frac{f_i}{\pi_i} + \lambda$$

$$\Rightarrow \pi_i = -\frac{f_i}{\lambda}$$

$$\text{and } 1 = \sum \pi_i = -\frac{\sum f_i}{\lambda} = -\frac{n}{\lambda}.$$

$$\Rightarrow -\frac{1}{\lambda} = \frac{1}{n}$$

$$\Rightarrow \hat{\pi}_i = \frac{f_i}{n}, \text{ are the MLE's.}$$

$$\Rightarrow \hat{\pi} = \left( \frac{f_1}{n}, \frac{f_2}{n}, \dots, \frac{f_6}{n} \right).$$

The MLE of  $(\pi_1 - 2\pi_2 + 6\pi_5)$

$$= (1, -2, 0, 0, 6, 0) \begin{pmatrix} \pi_1 \\ \pi_2 \\ \pi_3 \\ \pi_4 \\ \pi_5 \\ \pi_6 \end{pmatrix}$$

$$= \hat{\pi}' \pi$$

$$\text{is } \hat{\pi}' \hat{\pi} = \hat{\pi}_1 - 2\hat{\pi}_2 + 6\hat{\pi}_5.$$

→ The IQ's of 10 teenagers are 98, 114, 105, 101, 123, 117, 106, 92, 110, 108, whereas those of 8 teenagers belonging to other ethnic group are 122, 105, 95, 126, 114, 108.

Assuming that the data can be looked upon as independent s.s. from Normal popn. with mean  $\mu_1$  and  $\mu_2$  and common variance  $\sigma^2$ . Estimate the parameters by the method of MLE.

Sol. Let  $X_1, X_2, \dots, X_{10} \stackrel{iid}{\sim} N(\mu_1, \sigma^2)$  independent.

$Y_1, Y_2, \dots, Y_8 \stackrel{iid}{\sim} N(\mu_2, \sigma^2)$

Likelihood Equation:-

$$L(\mu_1, \mu_2, \sigma | x_i, y_j) = \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^{16} \cdot \exp \left[ - \frac{\sum (x_i - \mu_1)^2 + \sum (y_j - \mu_2)^2}{2\sigma^2} \right]$$

Likelihood Equation:-

$$\textcircled{1} \rightarrow 0 = \frac{\partial \ln L}{\partial \mu_1} = - \frac{2 \sum (x_i - \mu_1) (-1)}{2\sigma^2}$$

$$\Rightarrow \mu_1 = \bar{x}.$$

$$\textcircled{2} \rightarrow 0 = \frac{\partial \ln L}{\partial \mu_2} \Rightarrow \mu_2 = \bar{y}.$$

$$\textcircled{3} \rightarrow 0 = \frac{\partial \ln L}{\partial \sigma} = -\frac{16}{\sigma} + \frac{\sum (x_i - \mu_1)^2 + \sum (y_j - \mu_2)^2}{\sigma^3}$$

$$\Rightarrow \sigma^2 = \frac{\sum (x_i - \hat{\mu}_1)^2 + \sum (y_j - \hat{\mu}_2)^2}{16} = \dots$$

8) Scores obtained from 10 students in two parallel forms of a test of a test are recorded below:

Assuming that scores in the two sets arise from a bivariate normal popn. with the common mean  $\mu$ , common s.d.  $\sigma$  and the corr. coeff.  $\rho$ , find the MLE of  $\mu, \sigma, \rho$ .

1st form                    2nd form

|    |    |
|----|----|
| 91 | 47 |
| 52 | 54 |
| 49 | 50 |
| 44 | 48 |
| 32 | 40 |
| 62 | 65 |
| 40 | 38 |
| 45 | 42 |
| 39 | 45 |
| 48 | 45 |

Sol. Let  $(X_i, Y_i), i=1(1)10$ , be a r.s. from  $\text{BN}(\mu, \mu, \sigma^2, \sigma^2, \rho)$   
 Let  $U_i = X_i + Y_i \sim_{\text{iid}} N(2\mu, 2\sigma^2(1+\rho)) = N(\mu_U, \sigma_U^2)$   
 $V_i = X_i - Y_i \sim_{\text{iid}} N(0, 2\sigma^2(1-\rho)), i=1(1)n$ , indep.  
 $\equiv N(0, \sigma_V^2)$

MLE :- (i)  $\begin{cases} \hat{\mu}_U = \bar{U} \\ 2\hat{\mu} = \bar{U} = \bar{X} + \bar{Y} \\ \Rightarrow \hat{\mu} = \frac{\bar{X} + \bar{Y}}{2}. \end{cases}$

(ii)  $\hat{\sigma}_U^2 = s_U^2 = \frac{1}{n} \sum (U_i - \bar{U})^2 = \frac{1}{n} \sum (X_i - \bar{X} + Y_i - \bar{Y})^2$   
 $= \frac{1}{n} \sum_i [(X_i - \bar{X})^2 + (Y_i - \bar{Y})^2 + 2(X_i - \bar{X})(Y_i - \bar{Y})]$   
 $\Rightarrow 2\hat{\sigma}^2(1+\rho) = s_x^2 + s_y^2 + 2s_{xy}.$

(iii)  $\hat{\sigma}_V^2 = s_V^2 = \frac{1}{n} \sum (V_i - \bar{V})^2 = \frac{1}{n} \sum_i [(X_i - \bar{X})^2 + (Y_i - \bar{Y})^2 - 2(X_i - \bar{X})(Y_i - \bar{Y})]$   
 $\Rightarrow 2\hat{\sigma}^2(1-\rho) = s_x^2 + s_y^2 - 2s_{xy}$

(ii) + (iii) gives  $4\hat{\sigma}^2 = 2(s_x^2 + s_y^2)$   
 $\Rightarrow \hat{\sigma}^2 = \frac{s_x^2 + s_y^2}{2} = \underline{\hspace{2cm}}$ .

(ii) - (iii) gives

$$4\hat{\sigma}^2\rho = 4s_{xy}$$

$$\Rightarrow \hat{\rho} = \frac{2s_{xy}}{s_x^2 + s_y^2} = \underline{\hspace{2cm}}$$

9. Suppose heights of 6 pairs of identical adult male Bengali films, there observed that

5'6"; 5'5"; 6'1"; 5'11"; 5'7"; 5'8"; 6'2"; 6'; 5'9"; 5'3";  
 5'5"; 5'3";

Suppose it is known that average height of adult Bengali is 5'5" and S.D. is 3'3". Find the MLE of  $\rho$ .

Assuming that pairs of height follow a BN distn., with unknown corr. coeff.  $\rho$ .

Sol. Assume that,  $(X, Y) \sim \text{BN}(\mu, \mu, \sigma, \sigma, \rho)$

$$\mu = 5'5"$$

$$\sigma = 3'3"$$

and  $\rho$  is unknown.

$$\Rightarrow (U, V) = \left( \frac{X-\mu}{\sigma}, \frac{Y-\mu}{\sigma} \right) \sim \text{BN}(0, 0, 1, 1, \rho)$$

Likelihood function is :-

$$L(\rho | (u_i, v_i)) = \left\{ \frac{1}{2\pi\sqrt{1-\rho^2}} \right\}^n \cdot e^{-\frac{1}{2(1-\rho^2)} \{ \sum u_i^2 + \sum v_i^2 - 2\rho \sum u_i v_i \}}$$

Likelihood equation:-

$$\begin{aligned} 0 = \frac{\partial}{\partial \rho} \ln L &= - \frac{n(-2\rho)}{2(1-\rho^2)} - \frac{1}{2(1-\rho^2)} \{ -2 \sum u_i v_i \} \\ &\quad + \frac{\sum u_i^2 + \sum v_i^2 - 2\rho \sum u_i v_i}{2(1-\rho^2)^2} (-2\rho) \end{aligned}$$

$$\Rightarrow \frac{\rho}{(1-\rho^2)} \left\{ n + \frac{1}{\rho} \sum u_i v_i - \frac{1}{1-\rho^2} (\sum u_i^2 + \sum v_i^2 - 2\rho \sum u_i v_i) \right\} = 0$$

$$\Rightarrow np^3 - \rho^2 \sum u_i v_i + \rho (\sum u_i^2 + \sum v_i^2 - n) - \sum u_i v_i = 0$$

$$\Rightarrow \rho^3 + a\rho^2 + b\rho + c = 0, \text{ say.}$$

Use some numerical method of solution.

10. Estimate the value of  $\theta$  by the ML method and find out 95% C.I. on the basis of this sample

3.58, 2.86, 3.16, 2.46, 6.33, 8.31, 15.17, 9.59

drawn from the Cauchy population -

$$DF(x) = \frac{1}{\pi \{ 1 + (x - \theta)^2 \}}, \theta \in \mathbb{R}, -\infty < x < \infty.$$

Sol.  $L(\theta | x_1, \dots, x_n) = \frac{1}{\pi^n \prod_{i=1}^n \{ 1 + (x_i - \theta)^2 \}}, \theta \in \mathbb{R}.$

Likelihood equation:-

$$0 = \frac{\partial}{\partial \theta} \ln L = \sum_{i=1}^n \frac{2(x_i - \theta)}{1 + (x_i - \theta)^2} = f(\theta), \text{ say.}$$

By NR method,

$$\theta_{n+1} = \theta_n - \left\{ \frac{f(\theta)}{f'(\theta)} \right\}_{\theta=\theta_n}$$

$$= \theta_n + \left\{ \frac{\frac{\partial \ln L}{\partial \theta}}{-\frac{\partial^2 \ln L}{\partial \theta^2}} \right\}_{\theta=\theta_n}$$

$$\approx \theta_n + \left\{ \frac{\frac{\partial \ln L}{\partial \theta}}{\ln(\theta)} \right\}_{\theta=\theta_n}$$

replacing  $-\frac{\partial^2 \ln L}{\partial \theta^2}$  by its expectation  $\ln(\theta)$ .

Here  $\ln(\theta) = \frac{n}{2}.$

$$\theta_{n+1} = \theta_n + \frac{4}{n} \sum_{i=1}^n \frac{(x_i - \theta_n)}{1 + (x_i - \theta_n)^2}$$

Trial root  $\hat{\theta}_1 = \bar{x} =$  the sample median = \_\_\_\_\_.

$$\hat{\theta}_2 = \text{_____}$$

$$\hat{\theta}_3 = \text{_____}$$

etc.  $\hat{\theta} = \text{_____}$ .

For large sample,

$$\hat{\theta} \sim N(\theta, \frac{1}{In(\theta)})$$

$$\Rightarrow \hat{\theta} \sim N(\theta, \frac{2}{n})$$

$$\Rightarrow 1 - \alpha = P\left[ \left| \frac{\hat{\theta} - \theta}{\sqrt{\frac{2}{n}}} \right| < \chi_{\alpha/2} \right]$$

$$= P\left[ \theta - \sqrt{\frac{2}{n}} \chi_{\alpha/2} < \hat{\theta} < \theta + \sqrt{\frac{2}{n}} \chi_{\alpha/2} \right]$$

11. To test a null hypothesis  $H_0: p = 0.6$  of a popl'n.  $\sim B(m=5, p)$  vs.  
 $H_1: p > 0.6$ . Suggest an UMP test with exact size  $\alpha$ . The choice of  $\alpha$   
is left to you. If the observed no. of success is 4, construct a  
randomized test of exact size 0.1 and conclude.

Sol.

To test  $H_0: p = 0.6$  vs.  $H_1: p > 0.6$

Let  $x$  be an observation from  $B(m=5, p)$

MP test of its size of testing  $H_0: p = p_0$  vs.  $H_1: p = p_1, p_1 > p_0 = 0.6$

is

$$\phi(x) = \begin{cases} 1 & \text{if } \frac{f(x, p_1)}{f(x, p_0)} > c \\ 0 & \text{ow} \end{cases}$$

$$= \begin{cases} 1 & , x > k \\ 0 & , \text{ow} \end{cases}$$

where,  $P_{H_0}[\phi(x)] = \alpha$  (say)

and the test  $\phi(x)$  is independent of  $p_1 (> p_0)$ .

Hence, the test  $\phi(x)$  is UMP of  $H_0: p = 0.6$  vs.  $H_1: p > 0.6$  of its size.

Let  $k=4$ ,

then  $\phi(x) = \begin{cases} 1 & , x > 4 \\ 0 & , \text{ow} \end{cases}$

is a UMP test of size =  $E_{H_0}[\phi(x)] = P_{H_0}[X > 4]$

$$= \binom{5}{4} (0.6)^4 (0.4)^1$$

$$= 0.07776$$

$$\text{If } \phi(x) = \begin{cases} 1, & x > 3 \\ 0, & \text{otherwise} \end{cases}$$

$$\text{then size} = P_{H_0}[x > 3]$$

$$= \binom{5}{4} (0.6)^4 (0.4)^1 + (0.6)^5$$

$$= 0.2592 + 0.07776 > 0.1$$

To get the exact size 0.1, it is required to randomize at  $x=4$ .

$$\text{Let } \phi(x) = \begin{cases} 1, & x > 4 \\ \gamma, & x = 4 \\ 0, & \text{otherwise} \end{cases}$$

where  $\gamma$  is such that

$$0.1 = E_{H_0}[\phi(x)]$$

$$= 1 \cdot P_{H_0}[x > 4] + \gamma \cdot P_{H_0}[x = 4]$$

$$= 0.07776 + \gamma (0.2592)$$

$$\Rightarrow \gamma = 0.086$$

$$\phi(x) = \begin{cases} 1, & x > 4 \\ 0.086, & x = 4 \\ 0, & x < 4 \end{cases}$$

is the UMP test at level  $\alpha = 0.1$ .

Here  $x=4$  is the observed value.

Hence, we reject  $H_0$  with probability  $\gamma = 0.086$  and accept  $H_0$  with prob.  $1-\gamma$ .

Draw a 3-digit random number, then the prob. of  $R = \{\text{the selected no.} \leq 0.85\}$  is  $P(R) = \frac{86}{1000}$ .

Let the selected no. be 126.

Then we accept  $H_0: p = 0.6$  at level 0.1.

12. Suppose the coating time denoted by  $x$  in minutes for a bus is uniformly distributed over  $U[0, \theta]$ . To test the hypothesis  $H_0: \theta = 10$  vs.  $H_1: \theta > 10$  on the basis of a sample of size 6 to decision rules are proposed. Reject  $H_0$  if (i)  $\max\{x_1, \dots, x_6\} > c_1$

$$(ii) (\text{no. of } \{x_1, \dots, x_6\} > 8) > c_2$$

Find the values of  $c_1$  and  $c_2$  taking level of significance is 0.05. Also, draw power curve for both these procedure and comment on the relative performance.

Sol. Let  $X_1, \dots, X_6 \sim \text{iid } R(\theta, \theta)$

To test  $H_0: \theta = 10$  vs.  $H_1: \theta \neq 10$

$$(a) \phi_1(x) = \begin{cases} 1 & \text{if } x_{(6)} > c_1 \\ 0 & \text{ow} \end{cases}$$

$$(b) \phi_2(x) = \begin{cases} 1 & \text{if } Y > c_2 \\ 0 & \text{ow} \end{cases}$$

where  $Y = \text{the no. of } x_i's \text{ which are greater than } 8$ ,  
 $\therefore Y \sim \text{Bin}(6, p)$ .

$$\therefore p = \int_{\frac{8}{\theta}}^{\frac{10}{\theta}} \frac{1}{\theta} d\theta = \left(1 - \frac{8}{\theta}\right) = P[X_i > 8]$$

Power function:-

$$(a) \beta_1(\theta) = P_\theta [X_{(6)} > c_1] \\ = 1 - P_\theta [X_{(6)} \leq c_1] \\ = 1 - \left(\frac{c_1}{\theta}\right)^6.$$

$$(b) \beta_2(\theta) = P_\theta [Y > c_2]$$

$$P = \sum_{y=c_2+1}^{6} \binom{6}{y} \left(1 - \frac{8}{\theta}\right)^y \left(\frac{8}{\theta}\right)^{6-y}$$

Here  $\alpha = 0.05$ ,

$$\therefore 0.05 = \beta_1(\theta) \Big|_{\theta=10} = 1 - \left(\frac{c_1}{10}\right)^6$$

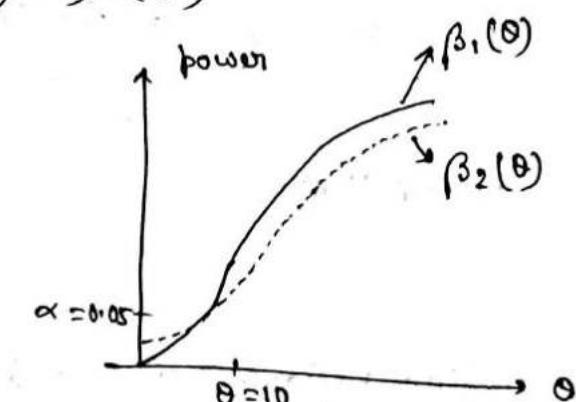
$$\Rightarrow c_1 = 10 \cdot (0.95)^{1/6}.$$

$$\text{and } 0.05 = \beta_2(\theta) \Big|_{\theta=10}$$

$$= \sum_{y=c_2+1}^{6} \binom{6}{y} \left(\frac{1}{5}\right)^y \left(\frac{4}{5}\right)^{6-y}$$

Comment:- The power of the test

(a) which is based on exact  
 $x_i's$  is greater than that  
of the test (b) which is  
based on Bernoulli trial.



## - : Non-Parametric Inference:-

13. The following are 15 measurements of the octane rating of a certain kind of Gasoline  
 $98.5, 95.2, 97.3, 96.0, 96.8, 100.3, 97.4; 95.3, 93.2, 99.1, 96.1,$   
 $97.6, 98.2, 98.5, 94.9.$

Use an exact non-parametric test and also an approximate test to examine whether or not the average octane rating of the given kind of Gasoline is 98.5.

Sol.  $X$ : the measurement of the octane of gasoline.  
Let  $x_1, x_2, \dots, x_n, n=15$  be given random samples.  
To test  $H_0: \mu_{Y/2} = 98.5$ , vs.  $H_1: \mu_{Y/2} \neq 98.5$ .

(a) Exact Non-Parametric Test:- count the no. of +ve signs among

$$x_1 - \mu_0, \dots, x_n - \mu_0.$$

We have two +ve sign and 12 -ve sign and one observation  $x_{14}$  is  $= 98.5$ . We ignore the value  $x_{14}$  from the sample.

Now, we have a r.s. of size  $n=14$ .

Let  $Y$  = the no. of +ve signs among  $(x_i - \mu_0)$ 's,  $i=1(1)14$ .

Under  $H_0$ ,  $Y \sim \text{Bin}(n=14, p=\frac{1}{2})$ .

The observed value of  $Y$  is  $y_0 = 2$ .

$$\begin{aligned}\text{The p-value} &= 2 \cdot \min \{ P_{H_0}[Y \geq y_0], P_{H_0}[Y \leq y_0] \} \\ &= 2 \cdot P_{H_0}[Y \leq y_0] \\ &= 2 \cdot \sum_{y=0}^2 \binom{14}{y} \cdot \frac{1}{2^{14}} = \underline{\underline{\quad}}.\end{aligned}$$

If p-value  $< 0.05$ , we reject  $H_0$ , at 5% level of significance.

(b) Assume that the distn. of  $X$  is  $N(\mu, \sigma^2)$ , whether the actual distn. is normal or not.

Here,  $\mu_{Y/2} = \mu$ .

To test  $H_0: \mu = 98.5$  vs.  $H_1: \mu \neq 98.5$

$$\text{t-test: } t = \frac{\bar{x} - 98.5}{\frac{s}{\sqrt{n}}} \sim t_{n-1}, \text{ under } H_0.$$

$n=15.$

Computation:-

Critical region:- Observed  $|t| > t_{\alpha/2, 14}$ .

The test procedure described is an approximate test procedure, as the actual distn. of  $x$  may not be normal.

14. The following table shows Hamilton depression scale factors measurements in 9 patients suffering from depression, taken before (X) and after (Y) a visit to therapist:

|   |       |       |       |      |      |      |      |      |      |
|---|-------|-------|-------|------|------|------|------|------|------|
| X | 1.83  | 0.50  | 1.62  | 2.48 | 1.68 | 1.88 | 1.55 | 3.06 | 1.3  |
| Y | 1.878 | 0.647 | 0.598 | 2.05 | 1.06 | 1.29 | 1.06 | 8.19 | 1.29 |

Perform a suitable (a) parametric (b) non parametric test to judge whether the therapy can be considered to be effective.

Sol. (a) Parametric:-

Assume that  $(X, Y) \sim BN(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$

To judge whether the therapy is considered to be effective, i.e.

$H_0: \mu_x = \mu_y$  alternative  $H_1: \mu_x > \mu_y$ .

This is a paired-t testing problem.

Define  $d_i = x_i - y_i \stackrel{iid}{\sim} N(\mu_d, \sigma_d^2)$ ,

where  $\mu_d = \mu_x - \mu_y$ .

Test statistic  $\frac{\sqrt{n}(\bar{d} - 0)}{s_d} \sim t_{n-1}$ , under  $H_0$ .

Computation:-

$$\begin{array}{c|ccccccccc}
& 1 & 2 & \dots & 9 \\
\hline d_i & \text{---} & \text{---}
\end{array}$$

$$\bar{d} = \bar{d}_F = \frac{1}{n} \sum d_i$$

$$s_d^2 = \frac{1}{n-1} \left\{ \sum d_i^2 - n\bar{d}^2 \right\}$$

$$= \text{_____}$$

Critical regions:- Observed  $t > t_{\alpha/2}$ .

(b) Non-Parametric:-

Here  $d_i = x_i - y_i, i=1(1)9$ .

And we are interested in testing  $H_0: \text{median}(d) = 0$  vs.  $H_1: \text{median}(d) > 0$

Test:-  $H_0: \gamma_{1/2}(d) = 0$  vs.  $H_1: \gamma_{1/2}(d) > 0$

We shall use sign test based on the values

$d_1, d_2, \dots, d_9$  as a r.s. from an absolutely continuous univariate distn.

Let  $Z = \text{the no. of +ve } d_i's$

$\sim \text{Bin}(n=9, p=1/2)$ , under  $H_0$ .

The observed value of  $Z$  is  $Z_0 = \text{_____}$ .

p-value =  $P_{H_0}[Z \geq Z_0] = \text{_____}$ .

15. Conduct the parametric and non-parametric test for differences of location for the data given below:

loop I : 9 7 12 16 14  
 loop II : 5 2 8 4 20

Sol. Let  $x_1, \dots, x_m, m=5$  and  $y_1, \dots, y_n, n=5$  be the given data from group I, group II, respectively.

Parametric test:- (Fisher's t-test)

Assumption:  $x_1, \dots, x_m \stackrel{iid}{\sim} N(\mu_x, \sigma_x^2)$  > independently distributed  
 $y_1, \dots, y_n \stackrel{iid}{\sim} N(\mu_y, \sigma_y^2)$

Assume that  $\sigma_x = \sigma_y$

To test:  $H_0: \mu_x = \mu_y$  vs.  $H_1: \mu_x \neq \mu_y$

Fisher's t-statistic:

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{1}{m} + \frac{1}{n}}} \sim t_{m+n-2}$$

Non-parametric test:- (Median test)

Assume that the distns. of the two groups are absolutely continuous and independent.

To test  $H_0: \text{Eq}_{1/2}(x) = \text{Eq}_{1/2}(y)$  vs.  $H_1: \text{Eq}_{1/2}(x) \neq \text{Eq}_{1/2}(y)$

Let  $\tilde{z}$  be the median in the combined set.

Define,  $V = \text{the no. of } x_i's \text{ which are } \leq \tilde{z}$

Under  $H_0$ ,  $P[V=v] = \frac{\binom{m}{v} \binom{n}{p-v}}{\binom{m+n}{p}}$ ,  $v=0(1)n$   
 $m+n=10=2p$   
 $\therefore p=5$ .

Combine the data in an increasing order, we have

2, 4, 5, 7, 9, 12, 14, 16, 20.

$$\tilde{z} = \frac{8+9}{2} = 8.5$$

The observed  $V = 1 = v$

$$\begin{aligned} \text{∴ p-value} &= 2 \min \{ P_{H_0}[V \geq v_0], P_{H_0}[V \leq v_0] \} \\ &= 2 P_{H_0}[V \leq v_0 = 1] \\ &= 2 \cdot \frac{\binom{5}{v} \binom{n}{5-v}}{\binom{10}{5}} \end{aligned}$$

## ANOVA

1. The following measurements refer to the no. of hours in which 9 patients are free from pain after taking placebo, a new drug and aspirin:

| Group    | Observations |          |         |
|----------|--------------|----------|---------|
|          | Placebo      | New drug | Aspirin |
| Placebo  | 0.0          | 1.0      |         |
| New drug | 2.8          | 3.5      | 2.8     |
| Aspirin  | 3.1          | 2.7      | 3.8     |

- (i) Test whether the new drug is more effective than the others  
(ii) Test if the effect of the new drug is the same as the average effect of the other two.

Sol. The factor considered here is drug (A) with levels Placebo ( $A_1$ ), New drug ( $A_2$ ) and Aspirin ( $A_3$ ). The data given is one-way classified data. Let  $y_{ij}$  denotes the  $j^{\text{th}}$  observations in the  $i^{\text{th}}$  level or group,  $i=1, 2, 3$ ,  $j=1(n_i)$ ;  $n_1=2$ ,  $n_2=4$ ,  $n_3=3$ .

Model [Fixed Effects] :-

$$y_{ij} = \mu + \alpha_i + \epsilon_{ij} \text{ with } \sum n_i \alpha_i = 0.$$

where,  $\epsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma_e^2)$ ,  $i=1(1)^3$ ,  $j=1(n_i)$ .

Here  $\mu + \alpha_i$  is mean-effect (fixed) of the  $i^{\text{th}}$  level.

To test  $H_0: \alpha_1 = \alpha_2 = \alpha_3 = 0$ .

computation:-  $G_1 = \sum_{i=1}^3 \sum_{j=1}^{n_i} y_{ij} = \underline{\hspace{10cm}}$

$$C.F. = \frac{G_1^2}{n} = \underline{\hspace{10cm}} \text{ with } n=9.$$

$$SS(\text{Total}) = \sum \sum y_{ij}^2 - CF = \underline{\hspace{10cm}}.$$

$$SS(\text{between}) = \sum_{i=1}^3 \frac{T_{i0}^2}{n_i} - CF = \underline{\hspace{10cm}}$$

$$= \frac{T_{10}^2}{n_1} + \frac{T_{20}^2}{n_2} + \frac{T_{30}^2}{n_3} - CF = \underline{\hspace{10cm}}.$$

$$\text{where } T_{i0} = \sum_{j=1}^{n_i} y_{ij}$$

ANOVA Table:-

| Source of Variation      | d.f.    | SS                                 | MS    | F                     |                                              |
|--------------------------|---------|------------------------------------|-------|-----------------------|----------------------------------------------|
| Between groups.          | $3-1=2$ | $SS(\text{between}) =$             | $MSB$ | $F = \frac{MSB}{MSE}$ | $F_{0.05, 2, 6} = \underline{\hspace{10cm}}$ |
| Within groups<br>[Error] | $9-3=6$ | $SS_E = \underline{\hspace{10cm}}$ | $MSE$ |                       | $= \underline{\hspace{10cm}}$                |
| Total =                  | $9-1=8$ | $SS(\text{total}) =$               |       |                       |                                              |

If observed  $F > F_{0.05}; 2, G$ , reject  $H_0$ .

(i) To test  $H_{01}: \mu_2 = \mu_1$ , vs.  $H_{01}': \mu_2 > \mu_1$   
 and  $H_{02}: \mu_2 = \mu_3$  vs.  $H_{02}': \mu_2 > \mu_3$ .  
 To test  $H_{01}$ , the test statistic

$$\frac{\bar{y}_{20} - \bar{y}_{10}}{\sqrt{MSE\left(\frac{1}{n_2} + \frac{1}{n_1}\right)}} \sim t_6, \text{ under } H_{01}.$$

Reject  $H_{01}$ , if observed  $t > t_{\alpha/2, 6}$ .

To test  $H_{02}$ ,  
 the test statistic is  $t_1 = \frac{\bar{y}_{20} - \bar{y}_{30}}{\sqrt{MSE\left(\frac{1}{n_2} + \frac{1}{n_3}\right)}}$

(ii) To test  $H_{03}: \mu_2 = \frac{\mu_1 + \mu_3}{2}$ ,

$$\Leftrightarrow H_{03}: \mu_1 - 2\mu_2 + \mu_3 = 0$$

$$\text{vs. } H_{03}': \mu_1 - 2\mu_2 + \mu_3 \neq 0.$$

Test statistic ( $t$ ) =  $\frac{\bar{y}_{10} - 2\bar{y}_{20} + \bar{y}_{30}}{\sqrt{MSE\left(\frac{1}{n_1} + \frac{4}{n_2} + \frac{1}{n_3}\right)}} \sim t_6, \text{ under } H_{03}.$

Critical region: observed  $|t| > t_{\alpha/2, 6}$ .

2. 4 randomly selected typists work on each of 4 given type writers of different make and speeds were recorded wise of typist and typewriter.

| Type-writer | Typist |        |        |        |
|-------------|--------|--------|--------|--------|
| 1           | 30, 32 | 40, 30 | 25, 40 | 20, 30 |
| 2           | 40, 35 | 32, 25 | 20, 30 | 31, 21 |
| 3           | 20, 31 | 25, 25 | 26, 30 | 25, 21 |
| 4           | 27, 28 | 31, 24 | 30,    | 19, 25 |

Analyse the data and give estimate of the parameters.

Sol. Type writer [factors A] :  $A_1, A_2, A_3, A_4$

Typists [factor B] :  $B_1, B_2, B_3, B_4$

The effects of different levels of A are fixed. As 4 typists are selected randomly, their effects are random, the data given is a two-way classification data with 2 obs.n. per cell. Let  $y_{ijk}$  be the  $k^{\text{th}}$  obs. corresponding to the  $i^{\text{th}}$  level of A and  $j^{\text{th}}$  level of B.

$$k=1(1)m=2, i=1(1)p=4; j=1(1)q=4.$$

Model:- [Two way classified data with mixed effects]

$$y_{ijk} = \mu + \alpha_i + \beta_j + c_{ij} + e_{ijk}$$

(fixed) (random)

$$\text{where } \sum_{i=1}^p \alpha_i = 0, \sum_{i=1}^p c_{ij} = 0, j=1(1)p$$

and  $\{c_{ij}\}$ ,  $\{e_{ijk}\}$  and  $\{e_{ijk}\}$  are jointly normal,  
 $e_{ijk} \sim_{\text{iid}} N(0, \sigma_e^2)$ .

$$\text{Define, } \sigma_A^2 = \frac{1}{p-1} \sum_{i=1}^p \alpha_i^2$$

$$\sigma_B^2 = \text{Var}(b_j)$$

$$\sigma_{AB}^2 = \frac{1}{p-1} \sum_{i=1}^p \text{Var}(c_{ij})$$

To test  $H_0: \sigma_A^2 = 0$  and to estimate the parameters  
 $\sigma_e^2, \sigma_B^2, \sigma_{AB}^2$ .

Computation:-

$$G_1 = \sum_i \sum_j \sum_k y_{ijk} = \text{_____}$$

$$\text{C.F. } G_1^2/n = \text{_____}$$

$$n = p \times m = 32.$$

$$SS(\text{Total}) = \sum_i \sum_j \sum_k y_{ijk}^2 - \text{C.F.}$$

$$\text{Defined, } T_{100} = \sum_j \sum_k y_{ijk}$$

$$T_{0j0} = \sum_i \sum_k y_{ijk}$$

$$T_{ij0} = \sum_k y_{ijk}.$$

| A\B            | B <sub>1</sub>           | B <sub>2</sub>           | B <sub>3</sub>           | B <sub>4</sub>           | Total                    |
|----------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| A <sub>1</sub> | T <sub>110</sub> = _____ | T <sub>120</sub> = _____ | T <sub>130</sub> = _____ | T <sub>140</sub> = _____ | T <sub>100</sub> = _____ |
| A <sub>2</sub> | T <sub>210</sub> = _____ | T <sub>220</sub> = _____ |                          |                          | T <sub>200</sub> = _____ |
| A <sub>3</sub> |                          |                          |                          |                          |                          |
| A <sub>4</sub> |                          |                          |                          |                          |                          |
|                | T <sub>010</sub> = _____ | T <sub>020</sub> = _____ |                          |                          | G <sub>1</sub> = _____   |

$$SSA = \sum_{i=1}^p \frac{T_{i00}^2}{qm} - CF = \underline{\quad}$$

$$SSB = \sum_j \frac{T_{j00}^2}{pm} - CF = \underline{\quad}$$

$$SS(AB) = \sum_i \sum_j \frac{T_{ij0}^2}{m} - \sum_i \frac{T_{i00}^2}{qm} - \sum_j \frac{T_{j00}^2}{pm} + CF.$$

$$= \left( \sum_i \sum_j \frac{T_{ij0}^2}{m} - CF \right) - SSA - SSB.$$

ANOVA Table

| Source of Variation | d.f.             | SS     | MS                          | E(MS)                                        |
|---------------------|------------------|--------|-----------------------------|----------------------------------------------|
| Due to A (fixed)    | $p-1 = 3$        | SSA    | $MSA = \frac{SSA}{3}$       | $\sigma_e^2 + m\sigma_{AB}^2 + qm\sigma_A^2$ |
| Due to B (random)   | $q-1 = 3$        | SSB    | $MSB = \frac{SSB}{3}$       | $\sigma_e^2 + pm\sigma_B^2$                  |
| Due to AXB          | $(p-1)(q-1) = 9$ | SS(AB) | $MS(AB) = \frac{SS(AB)}{9}$ | $\sigma_e^2 + m\sigma_{AB}^2$                |
| Error               | $pq(m-1) = 16$   | SSE    | $MSE = \frac{SSE}{16}$      | $\sigma_e^2 = \underline{\quad}$             |
| Total               | $pqm-1 = 31$     | SST    |                             |                                              |

$$F_A = \frac{MSA}{MS(AB)}$$

If observed  $F > F_{\alpha}; 3, 9$ , reject  $H_A: \sigma_A^2 = 0$ .

[ If  $H_A$  is rejected, give  $\hat{\alpha}_i = \bar{y}_{i00} - \bar{y}_{000} = \underline{\quad}$  ]

The estimates of the parameters  $\hat{\sigma}_e^2 = MSE = \underline{\quad}$

$$\hat{\sigma}_{AB}^2 = \frac{MS(AB) - MSE}{m} = \underline{\quad}$$

$$\hat{\sigma}_B^2 = \frac{MSB - MSE}{pm} = \underline{\quad}$$

g. In an experiment on yield of sugar beet (tons/acre) there were two levels of irrigation treatment and three of fertilizer treatment and each combination of treatments was carried out in 5 replicates. The following values of sum of squares (SS) were obtained : SS (Irrigation) = 120.0 ; SS (Fertilizer) = 221.7 ; SS (Interaction) = 35.0 ; SS (Error) = 108.0;

Assuming that the irrigation & fertilizer effects are random, estimate the components of variance. Also, test whether the variance of the irrigation component is zero.

Sol. Here two effects are : Irrigation (A) : A<sub>1</sub> A<sub>2</sub>  
Fertilizer (B) : B<sub>1</sub> B<sub>2</sub> B<sub>3</sub>

Both the effects are random.

It is two way classified data with m=5 obsn. per cell.

$$p=2, q=3, m=5.$$

Model:-  $y_{ijk} = \mu + a_i + b_j + c_{ij} + e_{ijk}$ , where

$$a_i \sim \text{iid } N(0, \sigma_A^2)$$

$$b_j \sim \text{iid } N(0, \sigma_B^2)$$

$$c_{ij} \sim \text{iid } N(0, \sigma_{AB}^2) \quad \text{are independent.}$$

$$e_{ijk} \sim \text{iid } N(0, \sigma_e^2)$$

$$\text{Note that, } \sigma^2 = V(y_{ijk}) = \sigma_A^2 + \sigma_B^2 + \sigma_{AB}^2 + \sigma_e^2.$$

ANOVA MODEL:-

| Source of variation | Df           | SS    | MS     | E(MS)                                          |
|---------------------|--------------|-------|--------|------------------------------------------------|
| due to A            | 2-1=1        | 120   | 120    | $\sigma_e^2 + m\sigma_{AB}^2 + 2m\sigma_A^2$   |
| due to B            | 3-1=2        | 221.7 | 110.85 | $\sigma_e^2 + m\sigma_{AB}^2 + p m \sigma_B^2$ |
| due to AB           | 2            | 35    | 17.5   | $\sigma_e^2 + m\sigma_{AB}^2$                  |
| Error               | 29 = pq(m-1) | 108   | 4.5    | $\sigma_e^2 = 4.5$                             |
| Total               | 29 = pqm-1   |       |        |                                                |

$$\text{Here } \hat{\sigma}_e^2 = \text{MSE} = 4.5$$

$$\hat{\sigma}_{AB}^2 = \frac{\text{MS}(AB) - \text{MSE}}{m} = \frac{17.5 - 4.5}{5} = 2.6$$

$$\hat{\sigma}_A^2 = \frac{\text{MS}A - \text{MS}(AB)}{qm} = \frac{120 - 17.5}{18} = \underline{\quad}$$

$$\hat{\sigma}_B^2 = \frac{\text{MS}B - \text{MS}(AB)}{pm} = \frac{110.85 - 17.5}{10} = \underline{\quad}$$

Reject H<sub>0</sub>:  $\sigma_A^2 = 0$  if observed  $F = \frac{\text{MS}A}{\text{MS}(AB)} > F_{\alpha; 1, 2}$ .

## DESIGN

1. In the experiment described below 4 materials were tested in each of 4 runs on a machine with 4 different position. The letters A to D refer to the 4 materials. The layout of the expt. is given below. Here the figures denote the loss in weight in a run of standard length.

| Run | Position in machine |        |        |        |
|-----|---------------------|--------|--------|--------|
|     | 4                   | 2      | 1      | 3      |
| 2   | A(251)              | B(241) | D(227) | C(229) |
| 3   | D(234)              | C(237) | A(274) | B(226) |
| 1   | C(235)              | D(236) | B(218) | A(268) |
| 4   | B(195)              | A(270) | C(230) | D(225) |

(a) Analyse the data and comment.

(b) If the variation due to the different position of the machine is ignored, will you modify your conclusion?

Sol. Factor A (Row): Factor runs are the four levels A  
 Factor B (Column): Four positions of the machine are the four levels of B.  
 Treatment: A, B, C, D.

(a) Experiment is conducted according to the LSD with four treatments A, B, C, D.

Let  $y_{ijk}$  be the obsn. of the  $k$ th treatment in the  $(i, j)$ th cell.

Model:-  $y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_k + \epsilon_{ijk}$

with  $\sum \alpha_i = \sum \beta_j = \sum \gamma_k = 0$ ,

$\epsilon_{ijk} \sim \text{iid}, N(0, \sigma^2_\epsilon)$

To test:  $H_0: \gamma_1 = \gamma_2 = \gamma_3 = \gamma_4 = 0$ ,

Computation:-

$$G_i = \sum_j \sum_k y_{ijk} = \text{_____}$$

$$\text{C.F.} = G_i^2 / n, \quad n=4.$$

$T_{i00} \rightarrow$  row totals

$T_{0j0} \rightarrow$  column totals

$T_{00k} \rightarrow$  treatment totals

$$SS(\text{row}) = \sum_{i=1}^m \frac{T_{i00}^2}{m} - \text{C.F.}$$

$$SS(\text{column}) = \sum_{j=1}^n \frac{T_{0j0}^2}{m} - \text{C.F.}$$

$$SS(\text{treatment}) = \sum_k \frac{T_{00k}^2}{m} - \text{C.F.}$$

### ANOVA Table

| Source of Variation | d.f.       | SS        | MS                       |
|---------------------|------------|-----------|--------------------------|
| Row                 | $m-1=3$    | SSR       | $F = \frac{MS(tr)}{MSE}$ |
| Column              | = 3        | SSE       |                          |
| Treatment           | = 3        | SS(tr)    | MS(tr)                   |
| Error               | 6          | —         | MSE                      |
| Total               | $m^2-1=15$ | SS(Total) |                          |

If observed  $F > F_{\alpha/2, 3, 6}$ , we reject  $H_0$ .

If  $H_0$  is rejected, then, the different treatment has different effect in general.

- (b) If the variation due to different position of the machine is ignored that is columns are ignored, then corresponding 4 rows as 4 blocks, the design of experiment reduced to RBD. Let  $y_{ik}$  be the observation on the  $k^{\text{th}}$  treatment in the  $i^{\text{th}}$  block (row).

Model:-  $y_{ik} = \mu + \alpha_i + \tau_k + \epsilon_{ik}$   
where,  $\epsilon_{ik} \sim i.i.d N(0, \sigma_e^2)$ .

$$\sum_i \alpha_i = 0 = \sum_k \tau_k$$

To test  $H_0: \tau_k = 0, k = 1, 2, 3, 4$ .

Computation:-

$$SS(\text{Block}) = SS(\text{Row})$$

SS(treatment) in RBD is same as the SS(tr) in LSD.

But SS(column) is added to SSE of LSD to get the SSE\* in RBD.

### ANOVA table

| Source of Variation | d.f.       | SS                                |
|---------------------|------------|-----------------------------------|
| Row                 | $m-1=3$    | SSR                               |
| Treatment           | 3          | SS(tr)                            |
| Error               | $3+6=9$    | $SSE^* = SSE + SE(\text{column})$ |
| Total               | $m^2-1=15$ | SS(Total)                         |

$$F^* = \frac{MS(tr)}{MSE^*}$$

If observed  $F^* > F_{0.05, 3, 9}$ , reject  $H_0$ .

2) In order to compare the hardness of alloy, three furnaces (F) and three levels of moulds (M) were tried. The layout as well as the hardness (in suitable unit) are shown below:

|                | Rep-I          |                |                | Rep-II         |                |                | Rep-III        |                |                |
|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|----------------|
|                | M <sub>3</sub> | M <sub>2</sub> | M <sub>1</sub> | M <sub>2</sub> | M <sub>1</sub> | M <sub>3</sub> | M <sub>1</sub> | M <sub>2</sub> | M <sub>3</sub> |
| F <sub>1</sub> | 156            | 118            | 140            | 104            | 89             | 117            | 103            | 126            | 149            |
|                | M <sub>3</sub> | M <sub>2</sub> | M <sub>1</sub> | M <sub>3</sub> | M <sub>1</sub> | M <sub>2</sub> | M <sub>3</sub> | M <sub>1</sub> | M <sub>2</sub> |
| F <sub>2</sub> | 130            | 174            | 157            | 112            | 89             | 81             | 144            | 124            | 129            |
|                | M <sub>1</sub> | M <sub>3</sub> | M <sub>2</sub> | M <sub>1</sub> | M <sub>2</sub> | M <sub>3</sub> | M <sub>3</sub> | M <sub>1</sub> | M <sub>2</sub> |
| F <sub>3</sub> | 114            | 161            | 141            | 103            | 132            | 133            | 100            | 91             | 97             |
|                | M <sub>1</sub> | M <sub>2</sub> | M <sub>3</sub> | M <sub>1</sub> | M <sub>2</sub> | M <sub>3</sub> | M <sub>1</sub> | M <sub>2</sub> | M <sub>3</sub> |

Analyse the data.

Sol. Split Plot Design:-

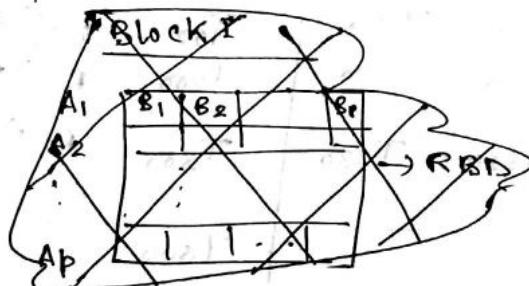
factors A: (A<sub>1</sub>, A<sub>2</sub>, ..., A<sub>p</sub>)

↳ they exclude use of small plots  
i.e. they are applied on large plots.

→ main effects are known to be different / effect of A<sub>i</sub>'s are to be tested with less precision.

Another factor:- B<sub>1</sub>, B<sub>2</sub>, ..., B<sub>q</sub>.

→ they are applicable  
small plots and we want test B and AxB more accurately than A.



F denoted by F<sub>1</sub>, F<sub>2</sub>, F<sub>3</sub>, M denoted by M<sub>1</sub>, M<sub>2</sub>, M<sub>3</sub>,

Here p=3, levels of F arranged is a R.B.D. using n=3 blocks or replicable and the factor M at q=3 levels, are applied to the plots of a block after subdividing each plot into q=3 subplots. This is a split plot design.

Model:-  $y_{ijk} = \{ \mu + b_i + T_j + e_{ij} \} + \gamma_k + \delta_{jk} + \epsilon_{ijk}$

model for whole plot  
treatment F by RBD

model for subplot treatment  
in whole plot.

$$i=1(1)n=3, j=1(1)p=3, k=1(1)q=3.$$

$$\sum T_j = \sum \gamma_k = \sum \delta_{jk} = \sum \epsilon_{ijk} = 0$$

| I              |                |                |
|----------------|----------------|----------------|
| F <sub>3</sub> | M <sub>3</sub> | M <sub>2</sub> |
| F <sub>1</sub> | M <sub>1</sub> | M <sub>3</sub> |
| F <sub>2</sub> |                | M <sub>2</sub> |

cohen,  $b_i$ ,  $e_{ij}$ ,  $e_{ijk}$  are independently normal ( $0, \sigma_b^2; \sigma_e^2; \sigma_{e'}^2$ ).

$$\text{To test } H_0: \gamma_j = 0$$

$$H_0: \gamma_K = 0$$

$$H_0: \delta_{jk} = 0$$

Computation:-

$$G_i = \sum_j \sum_k y_{ijk} = \dots$$

$$C.F. = G_i^2 / npq$$

$$SS(\text{total}) = \sum_j \sum_k y_{ijk} - C.F. = \dots$$

Whole Plot Analysis:-

$$SS(\text{Block}) = \sum_{i=1}^n \frac{T_{i00}^2}{pq} - C.F.$$

$$SS(F) = \sum \frac{T_{0j0}^2}{pq} - C.F.$$

$$SSE_I = \left\{ \sum_i \sum_j \frac{T_{ij0}^2}{pq} - C.F. \right\} - SSE_B - SSE_F$$

Whole Plot analysis Table:-

| Replicate \ F | F <sub>1</sub>  | F <sub>2</sub> | F <sub>3</sub> | Total     |
|---------------|-----------------|----------------|----------------|-----------|
| I             | $T_{110} = 461$ | $T_{120}$      | $T_{130}$      | $T_{100}$ |
| II            | $T_{210}$       | $T_{220}$      | $T_{230}$      | $T_{200}$ |
| III           | $T_{310}$       | $T_{320}$      | $T_{330}$      | $T_{300}$ |
|               |                 |                |                | $T_{000}$ |

Split plot Analysis:-

$$SS(M) = \sum \frac{T_{00k}^2}{np} - C.F.$$

$$SS(F \times M) = \left( \sum \frac{T_{0jk}^2}{n} - C.F. \right) - SSE - SSM$$

Table for SS(F × M)

| F \ M          | M <sub>1</sub>   | M <sub>2</sub>   | M <sub>3</sub>   |                  |
|----------------|------------------|------------------|------------------|------------------|
| F <sub>1</sub> | T <sub>011</sub> | T <sub>012</sub> | T <sub>013</sub> | T <sub>010</sub> |
| F <sub>2</sub> | T <sub>021</sub> | T <sub>022</sub> | T <sub>023</sub> | T <sub>020</sub> |
| F <sub>3</sub> | T <sub>031</sub> | T <sub>032</sub> | T <sub>033</sub> | T <sub>030</sub> |
|                | T <sub>001</sub> | T <sub>002</sub> | T <sub>003</sub> | T <sub>000</sub> |

$$F_1(M_1 + M_2 + M_3) = T_{011} \\ \text{under R-I, II, III}$$

$$F_1(M_2 + M_2 + M_3) = T_{012} \\ \text{under R-I, II, III}$$

ANOVA Table

| <u>Source of Variation</u> | <u>df</u>           | <u>SS</u> | <u>MS</u>        | <u>F</u>                                |
|----------------------------|---------------------|-----------|------------------|-----------------------------------------|
| Replicates                 | $M = 2$             |           | $MS(F)$          | $F_1 = \frac{MSF}{MSE_I}$               |
| whole plot treat (F)       | $b-1=2$             |           | $MS(E_I)$        |                                         |
| Error (I)                  | $(n-1)(p-1)=9$      |           | $MS(M)$          | $F_2 = \frac{MSM}{MSE_{II}}$            |
| Subplot treat (M)          | $q-1=2$             |           | $MS(F \times M)$ | $F_3 = \frac{MS(F \times M)}{MSE_{II}}$ |
| Interaction (F × M)        | $(p-1)(q-1)=4$      |           | $MS(E_{II})$     |                                         |
| Error (II)                 | $p(q-1)(n-1)=12$    |           |                  |                                         |
| Total                      | $df = bpq - 1 = 26$ |           |                  |                                         |

$$F_1 \leftrightarrow F \alpha; 2, 4$$

$$F_2 \leftrightarrow F \alpha; 2, 12$$

$$F_3 \leftrightarrow F \alpha; 4, 12$$

SAMPLE SURVEY

- ①. (a) Draw a random sample of size 7 from an exp. popn with mean 2.345.
- (b) Draw a n.s. of size 5 from a Cauchy popn. with median 0 and scale 2.
- (c) Draw a n.s. of size 6 from a univariate normal popn. with mean 17.95 and s.d. 6.28.
- (d) Draw a n.s. of size 5 from the distn.  
 $P[X=0] = \frac{1}{5}$ ,  $P[X=1] = \frac{2}{5}$ ,  $P[X=2] = \frac{2}{5}$ .

Solution: (a) Let  $X \sim \text{Exp}(\theta = 2.345)$

$$\left[ f_X(x) = \frac{1}{\theta} e^{-x/\theta} \right]$$

$$F_X(x) = \int_{-\infty}^x f_X(t) dt = 1 - e^{-x/\theta}$$

D.F. of  $X$  is

$$F(x) = 1 - e^{-x/\theta}, x > 0$$

By probability integral transformation,

$$U = F(X) \sim R(0,1)$$

If  $U$  is an observed sample from  $R(0,1)$ , then  $U = F(x)$

$$\Rightarrow u = 1 - e^{-x/\theta}$$

$$\Rightarrow x = -\ln(1-u)$$

$= -2.345 \ln(1-u)$  is an observed sample  
from Exp. with mean  $\theta = 2.345$ ,

We take seven 3-digit random number from Fisher-Yates table.

Page No: 125, row=5, column=4

260, 573, 375, 204, 056, 930, 001.

We place decimal points before the selected nos.

Then 0.260, 0.573, 0.375, 0.204, 0.056, 0.930, 0.001 are 7 n.s.  
from  $R(0,1)$ .

| Serial No. | $U_i$ | $x_i = -2.345 \ln(1-u_i)$ |
|------------|-------|---------------------------|
| 1          |       |                           |
| 2          |       |                           |
| 3          |       |                           |
| 4          |       |                           |
| 5          |       |                           |
| 6          |       |                           |
| 7          |       |                           |

$$(b) \quad f_X(x) = \frac{1}{\pi \sqrt{\sigma^2 + (x-\mu)^2}}$$

$$\Rightarrow F_X(x) = \int_{-\infty}^x f_X(t) dt$$

$$= \frac{1}{2} + \frac{1}{\pi} \tan^{-1}\left(\frac{x-\mu}{\sigma}\right). \text{ Here } \mu=0, \sigma=2.$$

$$F(x) = \frac{1}{2} + \frac{1}{\pi} \tan^{-1}\left(\frac{x}{2}\right)$$

where  $X \sim R(0,1)$

$$\text{Here } U=F(x)=\frac{1}{2} + \frac{1}{\pi} \tan^{-1}\left(\frac{x}{2}\right).$$

$$\Rightarrow x = 2 \tan\left\{\pi\left(U - \frac{1}{2}\right)\right\}$$

$$= 2 \tan\pi\left[U - \frac{\pi}{2}\right]$$

$$= -2 \tan\pi U.$$

$$(c) \quad X \sim N(17.95, (6.23)^2)$$

$$\Rightarrow Z = \frac{X-17.95}{6.23} \sim N(0,1)$$

$$\text{Hence, } U = \Phi\left(\frac{X-17.95}{6.23}\right) \sim R(0,1)$$

$$\Rightarrow U = \Phi\left(\frac{X-17.95}{6.23}\right)$$

$$\Rightarrow \Phi^{-1}(U) = \frac{X-17.95}{6.23} \rightarrow \text{from Biometrika.}$$

(d) Draw a n.r. of size 5 from the distn

$$P[X=0] = \frac{1}{5}, \quad P[X=1] = \frac{2}{5}, \quad P[X=2] = \frac{2}{5}.$$

Let us define the n.v.  $Y \sim U(0,1)$ .

$$\text{So, } P[0 < Y < .2] = .2 = P[X=0]$$

$$P[.2 \leq Y < .6] = .4 = P[X=1].$$

$$P[.6 \leq Y < 1] = .4 = P[X=2]$$

so, we take 5 3-digit random no. If the no. selected is in between (000-199) then we consider is equivalent to choosing  $X=0$ , if the no. is in between (200-599), we equivalently choose  $X=1$ , and finally, if the selected no. is in between (600-999), then we choose  $X=2$ .

- Q1. For the cost function a)  $C = c_0 + \sum c_n n_h$   
 b)  $C = c_0 + \sum c_n \sqrt{n_h}$ .

where  $c_0$  and  $c_n$  are known constants. Find the optimum values of  $n_h$ 's by minimizing  $\text{Var}(\bar{Y}_{\text{st}})$  for fixed total cost with total sample size being 60, given the following data:

| Stratum | 1   | 2   | 3   | 4   |
|---------|-----|-----|-----|-----|
| $N_h$   | 30  | 40  | 60  | 70  |
| $S_h$   | 1.5 | 2.0 | 3.5 | 4.0 |
| $C_h$   | 1   | 2   | 3   | 4   |

Also, compute the optimum allocation due to Neyman & compare the efficiency of the optimum allocations with that of proportional allocation to estimate the population mean  $\bar{Y}$ .

Solution:- (a) Define,  $F = \frac{1}{N^2} \sum_{n=1}^L N_h^2 \left( \frac{1}{n_h} - \frac{1}{N_h} \right) S_h^2 + \lambda \left( C_0 + \sum_{n=1}^L C_n n_h - C \right)$

Now,  $\frac{\partial F}{\partial n_h} = 0$

$$\Rightarrow \frac{1}{N^2} N_h^2 S_h^2 \left( -\frac{1}{n_h^2} \right) + \lambda \cdot c_n = 0$$

$$\Rightarrow n_h = \lambda_1 \cdot \frac{N_h \cdot S_h}{\sqrt{C_n}}$$

Given,  $\sum_h n_h = 60$

$$\Rightarrow \lambda_1 = \frac{60}{\sum_h \frac{N_h S_h}{\sqrt{C_n}}} = \dots$$

| Stratum | $N_h$ | $S_h$ | $C_h$ | $\frac{N_h S_h}{\sqrt{C_n}}$ | $n_h = \lambda_1 \cdot \frac{N_h S_h}{\sqrt{C_n}}$ |
|---------|-------|-------|-------|------------------------------|----------------------------------------------------|
| 1       |       |       |       |                              |                                                    |
| 2       |       |       |       |                              |                                                    |
| 3       |       |       |       |                              |                                                    |
| 4       |       |       |       |                              |                                                    |

(b) Define  $F = \frac{1}{N} \sum_{n=1}^L N_h^2 \left( \frac{1}{n_h} - \frac{1}{N_h} \right) \cdot S_h^2 + \lambda \left( C_0 + \sum_{n=1}^L C_n n_h^{1/2} - C \right)$

Now,  $0 = \frac{\partial F}{\partial n_h}$

$$\Rightarrow \frac{1}{N^2} \cdot N_h^2 \cdot S_h^2 \left( -\frac{1}{n_h^2} \right) + \lambda \cdot \frac{c_n}{2 \sqrt{n_h}} = 0$$

$$\Rightarrow n_h = \lambda_1^* \left( \frac{N_h S_h}{\sqrt{C_n}} \right)^{4/3}$$

Given  $\sum_h n_h = 60 \Rightarrow \lambda_1^* = \frac{60}{\sum_h \left( \frac{N_h S_h}{\sqrt{C_n}} \right)^{4/3}} = \dots$

| Stratum | $N_h$ | $S_h$ | $C_h$ | $\frac{N_h S_h}{\sqrt{C_n}}$ | $n_h = \lambda_1^* \left( \frac{N_h S_h}{\sqrt{C_n}} \right)^{4/3}$ |
|---------|-------|-------|-------|------------------------------|---------------------------------------------------------------------|
| 1       |       |       |       |                              |                                                                     |
| 2       |       |       |       |                              |                                                                     |
| 3       |       |       |       |                              |                                                                     |
| 4       |       |       |       |                              |                                                                     |

Neyman's optimum and proportional allocations:-

Here  $n_h \propto \frac{N_n S_n}{\sqrt{C_n}}$

① If  $C_n = \text{constant}$ , then  $n_h = \lambda_2 N_n S_n$  and  $60 = \sum n_h$ .

$$\Rightarrow \lambda_2 = \frac{60}{\sum N_n S_n}.$$

$$\therefore n_h = \left( \frac{60}{\sum N_n S_n} \right) \cdot N_n S_n.$$

→ optimum allocation.

② Here  $C_n = \text{constant}$

$S_n = \text{constant}$

$$n_h = \lambda_3 \cdot N_n$$

$$\text{and } \lambda_3 = \frac{60}{\sum N_n} = \frac{60}{N}.$$

$$\therefore n_h = \left( \frac{60}{N} \right) N_n \rightarrow \text{proportional allocation.}$$

③ Ques:- Using the information given below, find the relative standard error that one could expect if a sample of 1% of the villages group of district be selected by SRSWOR for estimating the total population  $\gamma$  of the region.

| District serial No. | No. of Villages ( $N_n$ ) | Average pop.ln for village ( $\bar{Y}_n$ ) | Standard deviation ( $S_n$ ) |
|---------------------|---------------------------|--------------------------------------------|------------------------------|
| 1                   | 1953                      | 487                                        | 564                          |
| 2                   | 1864                      | 829                                        | 931                          |
| 3                   | 1381                      | 822                                        | 996                          |
| 4                   | 1184                      | 1083                                       | 1167                         |
| 5                   | 531                       | 1956                                       | 1990                         |
| 6                   | 1391                      | 664                                        | 625                          |
| 7                   | 1996                      | 956                                        | 779                          |
| 8                   | 1951                      | 392                                        | 856                          |
| 9                   | 3369                      | 389                                        | 591                          |

Solution:-

$$\text{Relative standard Error} = RSE(\hat{\gamma}_{st})$$

$$= \frac{S.E(\hat{\gamma}_{st})}{\hat{\gamma}_{st}} = \frac{N \cdot SE(\hat{\gamma}_{st})}{N \cdot \hat{\gamma}_{st}}$$

$$= \frac{\left\{ \frac{1}{N^2} \sum N_n^2 S_n^2 \left( \frac{1}{n_h} - \frac{1}{N_n} \right) \right\}^{1/2}}{\frac{\sum N_n \cdot \bar{Y}_n}{N}}$$

$$\left[ \text{Here } \frac{n_h}{N_h} = \frac{1}{100} \right]$$

$$= \sqrt{\frac{2 N_h \left( N_h - \frac{N_h}{100} \right) \cdot \frac{s_h^2}{N_h}}{100}}$$

$$2 N_h \bar{Y}_h$$

$$= \sqrt{\frac{992 N_h s_h^2}{\sum N_h \bar{Y}_h}} = \dots$$

| District | $N_h$ | $\bar{Y}_h$ | $s_h$ | $N_h \bar{Y}_h$ | $N_h s_h^2$ |
|----------|-------|-------------|-------|-----------------|-------------|
| 1        |       |             |       |                 |             |
| 2        |       |             |       |                 |             |
| :        |       |             |       |                 |             |
| 9        |       |             |       |                 |             |

- ③ All the farms of a country have been stratified according to size and the following data have been obtained.

| Form Size | No. of Farm ( $N_h$ ) | Mean ( $\bar{Y}_h$ ) | S.D. ( $s_h$ ) |
|-----------|-----------------------|----------------------|----------------|
| 0-40      | 394                   | 5.4                  | 8.3            |
| 41-80     | 464                   | 16.3                 | 13.3           |
| 81-120    | 391                   | 24.3                 | 15.1           |
| 121-160   | 334                   | 34.5                 | 19.8           |
| 161-200   | 169                   | 42.1                 | 24.5           |
| 201-240   | 113                   | 50.1                 | 26.0           |
| 241-280   | 148                   | 63.8                 | 35.2           |

If the first 3 stratum combine into a 1 stratum and last 4 into another. Find the relative loss in efficiency as compared to 7 strata situation in estimating the popn. mean by stratified Random sampling (a) Proportional allocation (b) optimum allocation that the total no. of farms are selected 100.

Sol. In the new stratified popn

| Farm Size | No. of firms             | $\bar{Y}'_h$ | $s_h'^2$ |
|-----------|--------------------------|--------------|----------|
| 0-120     | $N_1' = \frac{3}{7} h_n$ |              |          |
| 121-280   | $N_2' = \frac{4}{7} h$   |              |          |

$$\text{where } \bar{Y}_1' = \frac{\sum_{h=1}^3 Y_h \bar{Y}_h}{\sum_{h=1}^3 N_h}, \quad \bar{Y}_2' = \frac{\sum_{h=4}^7 N_h \bar{Y}_h}{\sum_{h=4}^7 N_h}$$

$$\text{and } S_1'^2 = \frac{1}{\sum_{h=1}^3 N_h - 1} \left\{ \sum_{h=1}^3 (N_h - 1) S_h^2 + \sum_{h=1}^3 N_h (\bar{Y}_h - \bar{Y}_1')^2 \right\}$$

$$\text{and } S_2'^2 = \frac{1}{\sum_{h=4}^7 N_h - 1} \left\{ \sum_{h=4}^7 (N_h - 1) S_h^2 + \sum_{h=4}^7 N_h (\bar{Y}_h - \bar{Y}_2')^2 \right\}$$

(a)  $V_{\text{prob}} = \frac{1-f}{n} \sum_{h=1}^7 \frac{N_h}{N} \cdot S_h^2 \quad (\text{Given stratification})$

and  $V_{\text{prob}'} = \frac{1-f}{n} \sum_{h=1}^2 \frac{N_h'}{N} S_h'^2 \quad (\text{New stratification})$

Loss in eff =  $\frac{V_{\text{prob}} - V_{\text{prob}'}}{V_{\text{prob}}} = \dots$

(b)  $V_{\text{opt}} = \frac{\left( \sum_{h=1}^7 \frac{N_h}{N} \cdot S_h^2 \right)^2}{\sum_{h=1}^7 \frac{N_h}{N} \cdot S_h^2} - \frac{\sum_{h=1}^7 \frac{N_h}{N} \cdot S_h^2}{N}$

$$V_{\text{opt}'} = \frac{\sum_{h=1}^2 \frac{N_h'}{N} \cdot S_h'^2}{n} - \frac{\sum_{h=1}^2 \frac{N_h'}{N} \cdot S_h'^2}{N}$$

Loss in eff =  $\frac{V_{\text{opt}'} - V_{\text{opt}}}{V_{\text{opt}}} = \dots$

(4) (Ratio - Regression Estimator)

An experimenter makes an farmer's eye-estimate of the weight of peaches on each tree in a orchard of 200 trees. He finds a total weight of 11600 lbs and weight from a SRS of 10 trees which yield the following result:

Serial No. of tree: 1 2 3 4 5 6 7 8 9 10

Actual weight: 61 42 50 58 67 45 39 57 71 53

Estimated is: 59 47 52 60 67 48 44 58 76 58

Compute the ratio and regression estimates of the total actual weight of Peaches of all the 200 trees in the orchard and compare the precisions of 2 estimates.

Solution:-  $Y = \text{total weight}$ ,  $X = \text{an eye estimate}$

To estimate the "the weight of Peaches ( $y$ )" of all the 200 trees.  
Here auxiliary information "an eye estimate of weight of Peaches ( $x$ )" is given.

$$X = 11600 \text{ lbs.}$$

Data:  $(x_i, y_i), i=1(1)n$ .

(a) Ratio estimation: Ratio estimate  $\hat{Y}_R = \frac{\bar{y}}{\bar{x}} \cdot X = \frac{\bar{y}}{\bar{x}} \cdot X$   
where  $\bar{y} = \frac{1}{n} \sum y_i = \underline{\quad}$ ,  $\bar{x} = \frac{1}{n} \sum x_i = \underline{\quad}$ .

$$\Rightarrow \hat{Y}_R = \underline{\quad}$$

$$\text{MSE}(\hat{Y}_R) \cong N^2 \left( \frac{1}{n} - \frac{1}{N} \right) (\delta_y^2 + \delta_x^2 - 2\hat{R}\delta_{xy})$$

$$\text{where } \delta_x^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2, \delta_y^2 = \frac{1}{n-1} \sum (y_i - \bar{y})^2,$$

$$\delta_{xy} = \frac{1}{n-1} \sum (x_i - \bar{x})(y_i - \bar{y}), \hat{R} = \frac{\bar{y}}{\bar{x}} = \underline{\quad}$$

(b) Regression Estimation:

$$\text{Regression estimate}, \hat{Y}_{LR} = N \cdot \hat{Y}_{LR} = N \left\{ \bar{y} + b(\bar{x} - \bar{\bar{x}}) \right\}$$

$$= N\bar{y} + b(\bar{x} - N\bar{x})$$

$$b = \frac{\delta_{xy}}{\delta_x^2} = \underline{\quad},$$

$$\text{MSE}(\hat{Y}_{LR}) \cong N^2 \left( \frac{1-f}{n} \right) \delta_e^2; \delta_e^2 = \delta_y^2 (1-f^2) = \delta_y^2 \cdot x$$

$$\text{where } \delta_y^2 \cdot x = \frac{1}{n-2} \sum_i (y_i - \bar{y} - b(x_i - \bar{x}))^2$$

$$= \frac{1}{n-2} \left\{ \sum (y_i - \bar{y})^2 - b^2 \sum (x_i - \bar{x})^2 \right\}$$

$$= \frac{n-1}{n-2} \left\{ \delta_y^2 - b^2 \delta_x^2 \right\}$$

$$= \delta_e^2$$

$$\text{Comparison: } \text{Eff (regression/ratio)} = \frac{\text{MSE}(\hat{Y}_R)}{\text{MSE}(\hat{Y}_{LR})} = \underline{\quad}$$

⑤ Two-stage-cluster: In an experimental investigation, 100 fields each consisting of 16 plots of equal size, were sown with wheat. Out of 100 fields, 10 fields are selected by SRSWOR and out of each field 80 selected plots are selected by WOR to observe the yield. From the given observation following values are estimated.

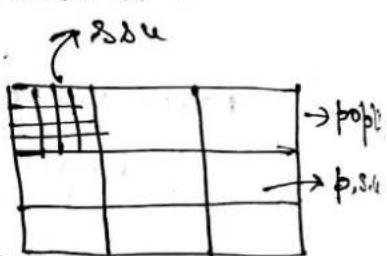
Sample mean (in kg) for the selected fields are:

4.290, 4.255, 3.795, 4.220, 4.070, 3.636, 4.550, 4.285, 4.375, 3.790. Sample average with variance = .018215 (kg)<sup>2</sup>.

Estimate the total yield of wheat in the experimental station along with its standard error.

Compare the efficiency of this estimate with one that would have been obtained by selecting an SRSWOR of 40 plots out of 1600 plots in the station.

Sol.: [First we select  $n$  p.s.u from  $N$  p.s.u. From selected in p.s.u., with in each p.s.u. have  $m$  units, we select  $m$  units.]



$N = 100$  fields (p.s.u.), each consisting of  $M = 16$  plots (s.s.u.). Out of  $N$  p.s.u.,  $n = 10$  fields and out of each selected (p.s.u.) fields  $m = 4$  plots (s.s.u.) are selected under SRSWOR.

Let  $y_{ij}$  be the value obtained for the  $j$ <sup>th</sup> selected s.s.u. in the  $i$ <sup>th</sup> selected p.s.u.,  $i = 1(1)10 = n$ ;  $j = 1(1)4 = m$ .

An unbiased estimate of  $\gamma$  (total yield) is

$$\hat{\gamma} = N \cdot \bar{y} = N \left( \frac{1}{n} \sum_{i=1}^n \bar{y}_i \right), \bar{y}_i = \frac{1}{m} \sum_{j=1}^m y_{ij}.$$

$$V(\hat{\gamma}) = (MN)^2 \left\{ \frac{1-f_1}{n} \cdot S_b^2 + \frac{1-f_2}{mn} \cdot S_w^2 \right\}$$

$$\text{and } V(\bar{y}_i) = (MN)^2 \left\{ \frac{1-f_1}{n} \cdot S_b^2 + f_1 \frac{(1-f_2)}{mn} \cdot S_w^2 \right\}$$

where  $S_b^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{y}_i - \bar{\bar{y}})^2 = \underline{\hspace{10cm}}$

and  $S_w^2 = \frac{1}{n(m-1)} \sum_{i=1}^n \sum_{j=1}^m (y_{ij} - \bar{y}_i)^2 = \underline{\hspace{10cm}}$

$$\hat{V}(\hat{\gamma}) = \underline{\hspace{10cm}},$$

In case of SRSWOR, with sample of size  
 $mn = 40 = n'$ . From  $1MN = 1600 = N'$  plots.

$$V_{SRS}(\hat{Y}) = N'^2 \cdot \frac{N' - n'}{N'n'} \cdot s_y^2,$$

$$s_y^2 = \frac{\sum_{i=1}^n \sum_{j=1}^m (y_{ij} - \bar{y})^2}{nm-1}.$$

$$= \frac{\sum_i \sum_j (y_{ij} - \bar{y}_i)^2 + m \sum_{i=1}^n (\bar{y}_i - \bar{y})^2}{nm-1}.$$

$$= \frac{n(m-1) \cdot s_w^2 + m(n-1) \cdot s_b^2}{mn-1}.$$

$$= \text{_____}.$$

.....