

**MUTIVARIATE ANALYSIS**  
**(DMSTT24)**  
**(MSC - STATISTICS)**



**ACHARYA NAGARJUNA UNIVERSITY**

**CENTRE FOR DISTANCE EDUCATION**

**NAGARJUNA NAGAR,**

**GUNTUR**

**ANDHRA PRADESH**

## UNIT-I

### **Concept of Multivariate Analysis:**

In botanical and zoological survey the interest of scientists is to investigate the new species of animals or trees and to identify their classes and families. A doctor wants to diagnose the disease of a newly arrived patient. The agricultural scientists investigate the suitability of a piece of land for the cultivation of a particular crop. The selection committee for appointing a group of individuals for a particular job applies different test procedures to select the candidates. Selection committee in educational institute selects students for admission into a course. In all the above cases the investigation is not done on the basis of a single criterion. For example, a doctor only examining the temperature of a patient cannot diagnose it as typhoid. He needs to observe some other symptoms for his proper diagnostic decision. The suitability of a piece of land for a particular crop is judge on test scales of fertility of land, the amount of potash in the land, the soil type, etc. The test scores on educational qualification, age, health condition, behavior, etc., of candidates help the selection committee in selecting a candidate.

The above discussion indicates that a decision regarding an object or individual depends on the simultaneous study of several characteristics observed from the objects. Since the characteristics are observed from an objects, these are inter-related. For a reliable and valid conclusion about any population parameter vector, it needs the study of inter-related variables observed from a sample of objects selected randomly from the population under study. Since the variables cannot be splitted off from each other as these are dependent among themselves.

From the above discussion it is clear that one needs to analyze all the variables observed from n sample objects simultaneously. Multivariate analysis is a statistical technique for simultaneous analysis of two or more variables observed from one or more sample objects. In this analysis, the inter-relationships of the variables are studied along with the

study of mean, variance and some other characteristics related to univariate analysis. However, the main objective of the analysis is to estimate the extent or amount of relationship among the variables. For example, one may need to observe the magnitude and direction of influences of some socio-economic variables which mobilize the couples to adopt family planning method, or to study the causes of preferring a particular industrial item by consumers, or to identify the causes of a disease which affects the community or to identify the class and family of newly observed species. All these analyses depend on multivariate data observed on same occasion from sample individuals.

From the objective of the analysis it is clear that the multivariate data are of two types, viz., dependent and independent. Accordingly, one needs to study the extent of relationship of dependent and independent sets. Also it needs the analysis of the structure of inter-relationships of the variables altogether. Thus multivariate analysis can be classified into two type, viz., (a) Inter-dependence analysis, and (b) Dependent analysis.

The main components of inter-dependence analysis are (i) Principal Component Analysis, (ii) Factor Analysis, (iii) Cluster analysis, where as the dependent analysis deals with (i) discriminant Analysis, (ii) Canonical Correlation Analysis, (iii) Multivariate Analysis of Variance, and (iv) Multivariate Regression Analysis.

**Random Vector:**

Definition: A random vector is a vector whose elements are random variables i.e., if  $X_1, X_2, \dots, X_n$  are random variables then the n-tuple

$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix} \text{ is called a random vector.}$$

**Distribution function of  $\underline{\mathbf{X}}$ :**

The Joint distribution function of  $\underline{\mathbf{X}}$  is defined by

$$\begin{aligned} F(\underline{\mathbf{X}}) &= F(x_1, x_2, \dots, x_n) \\ &= P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_n \leq x_n) \end{aligned}$$

Which is a monotonic non-decreasing in each of its arguments with basic consideration

$$F(-\infty, -\infty, \dots, -\infty) = 0$$

$$F(\infty, \infty, \dots, \infty) = 1$$

**Joint density function of  $\underline{\mathbf{X}}$ :**

The Joint probability density function of  $\underline{\mathbf{x}}$  is defined by

$$\begin{aligned} f(\underline{\mathbf{x}}) &= f(X_1, X_2, \dots, X_n) = \frac{\partial^n F(X_1, X_2, \dots, X_n)}{\partial x_1 \partial x_2 \dots \partial x_n} \\ &= \lim_{\Delta x_1 \rightarrow 0, \dots, \Delta x_n \rightarrow 0} P(x_1 < X_1 < x_1 + \Delta x_1, \dots, x_n < X_n < x_n + \Delta x_n) \end{aligned}$$

Conversely, the distribution function can be expressed in terms of the density function as follows

$$\begin{aligned} F(\underline{\mathbf{x}}) &= \int_{-\infty}^{\underline{\mathbf{x}}} f(\underline{\mathbf{x}}) \partial \underline{\mathbf{x}} \\ &= \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} \dots \int_{-\infty}^{x_n} f(X_1, X_2, \dots, X_n) \partial X_1 \partial X_2 \dots \partial X_n \end{aligned}$$

**Conditional density function:**

**Definition:** The conditional density function of  $\underline{\mathbf{X}}$  given that the random vector  $\underline{\mathbf{Y}}$  has a specified value  $\underline{\mathbf{y}}$  is defined by

$$\begin{aligned} f(\underline{\mathbf{x}}/\underline{\mathbf{y}}) &= \lim_{\Delta \underline{\mathbf{y}} \rightarrow 0} P(\underline{\mathbf{x}}/\underline{\mathbf{y}} \leq \underline{\mathbf{Y}} \leq \underline{\mathbf{y}} + \Delta \underline{\mathbf{y}}) \\ &= f(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) / f(y_1, y_2, \dots, y_n) \\ &= f(\underline{\mathbf{x}}, \underline{\mathbf{y}}) / f(\underline{\mathbf{y}}) \end{aligned}$$

where  $f(\underline{\mathbf{y}})$  is the marginal density of  $\underline{\mathbf{y}}$  and is obtained by the integration.

$$f(\underline{\mathbf{y}}) = \int_{\underline{\mathbf{x}}=-\infty}^{\infty} f(\underline{\mathbf{x}}, \underline{\mathbf{y}}) \partial \underline{\mathbf{x}}$$

**Mean vector or Expectation of  $\underline{\mathbf{X}}$ :**

The mean vector of  $\underline{\mathbf{X}}$  is denoted as  $\underline{\boldsymbol{\mu}}$  and is defined as  $E(\underline{\mathbf{X}})$  and

$$E(\underline{\mathbf{X}}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix} = \underline{\boldsymbol{\mu}} \text{ (say)}$$

**The Variance-Covariance matrix of  $\underline{\mathbf{X}}$ :**

The Variance-Covariance matrix of  $\underline{\mathbf{X}}$  is denoted by  $\boldsymbol{\Sigma}$  and is defined as

$$\boldsymbol{\Sigma} = E[(\underline{\mathbf{X}} - E(\underline{\mathbf{X}}))(\underline{\mathbf{X}} - E(\underline{\mathbf{X}}))']$$

$$= E[(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})']$$

$$= E \left[ \begin{bmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \\ \vdots \\ X_n - \mu_n \end{bmatrix} [X_1 - \mu_1, X_2 - \mu_2, \dots, X_n - \mu_n] \right]$$

$$= E \begin{bmatrix} [X_1 - \mu_1]^2 & [(X_1 - \mu_1)(X_2 - \mu_2)] & \dots & [(X_1 - \mu_1)(X_n - \mu_n)] \\ [(X_2 - \mu_2)(X_1 - \mu_1)] & [X_2 - \mu_2]^2 & \dots & [(X_2 - \mu_2)(X_n - \mu_n)] \\ \vdots & \vdots & \vdots & \vdots \\ [(X_n - \mu_n)(X_1 - \mu_1)] & [(X_n - \mu_n)(X_2 - \mu_2)] & \dots & [X_n - \mu_n]^2 \end{bmatrix}$$

$$= \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1n} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_{nn} \end{bmatrix}_{n \times n} \quad \text{where } \sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)] = \sigma_{ji}$$

We can see  $\boldsymbol{\Sigma}$  is a symmetric matrix

**Un-Correlated, Orthogonal and Independent Random Vectors:**

Definiton: Two random vectors  $\underline{\mathbf{X}}$  and  $\underline{\mathbf{Y}}$  are called **Uncorrelated** if

$$E(\underline{\mathbf{X}}\underline{\mathbf{Y}}') = E(\underline{\mathbf{X}})E(\underline{\mathbf{Y}}') \text{ (or) } \text{cov}(\underline{\mathbf{X}}, \underline{\mathbf{Y}}) = 0$$

**Orthogonal** if  $E(\underline{\mathbf{X}}\underline{\mathbf{Y}}') = 0$

**Independent** if  $f(\underline{\mathbf{X}}, \underline{\mathbf{Y}}) = f(\underline{\mathbf{X}})f(\underline{\mathbf{Y}})$

Note:

1. If  $\underline{\mathbf{X}}$  and  $\underline{\mathbf{Y}}$  are uncorrelated and either  $E(\underline{\mathbf{X}}) = 0$  (or)  $E(\underline{\mathbf{Y}}) = 0$  or then  $\underline{\mathbf{X}}$  and  $\underline{\mathbf{Y}}$  are orthogonal
2. If  $\underline{\mathbf{X}}$  and  $\underline{\mathbf{Y}}$  are independent then

$$\begin{aligned} E(\underline{\mathbf{X}}\underline{\mathbf{Y}}') &= \int \int \underline{\mathbf{x}}\underline{\mathbf{y}}' f(\underline{\mathbf{x}}, \underline{\mathbf{y}}) d\underline{\mathbf{x}} d\underline{\mathbf{y}} \\ &= \int \int \underline{\mathbf{x}}\underline{\mathbf{y}}' f(\underline{\mathbf{x}}) f(\underline{\mathbf{y}}) d\underline{\mathbf{x}} d\underline{\mathbf{y}} \\ &= \left[ \int \underline{\mathbf{x}} f(\underline{\mathbf{x}}) d\underline{\mathbf{x}} \right] \left[ \int \underline{\mathbf{y}}' f(\underline{\mathbf{y}}) d\underline{\mathbf{y}} \right] = E(\underline{\mathbf{X}})E(\underline{\mathbf{Y}}') \end{aligned}$$

i.e., and are independent implies uncorrelated

3. The converse of '2' need not be true

## MULTIVARIATE NORMAL DISTRIBUTION

MULTIVARIATE NORMAL DENSITY:

Suppose X is a scalar normal variate with mean  $\mu$  and variance  $\sigma^2$  then the p.d.f of X can be written as

$$f(x; \mu, \sigma) = k e^{-\frac{1}{2}(x-\mu)(\sigma^2)^{-1}(x-\mu)}, \sigma^2 > 0, -\infty < \mu < \infty \quad \rightarrow (1)$$

Where,  $k = \frac{1}{\sigma\sqrt{2\pi}}$

Now suppose  $\underline{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix}$  is a p-variate random vector and

Its mean vector is given by

$$E(\underline{X}) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix} = \underline{\mu} \quad \rightarrow (2)$$

and its variance –covariance matrix is given by

$$\begin{aligned} V(\underline{X}) &= E[(\underline{X} - E(\underline{X}))(\underline{X} - E(\underline{X}))'] \\ &= E[(\underline{X} - \underline{\mu})(\underline{X} - \underline{\mu})'] \\ V(\underline{X}) &= \begin{bmatrix} E[X_1 - \mu_1]^2 & E[(X_1 - \mu_1)(X_2 - \mu_2)] & \dots & E[(X_1 - \mu_1)(X_p - \mu_p)] \\ E[(X_2 - \mu_2)(X_1 - \mu_1)] & E[X_2 - \mu_2]^2 & \dots & E[(X_2 - \mu_2)(X_p - \mu_p)] \\ \vdots & \vdots & \vdots & \vdots \\ E[(X_p - \mu_p)(X_1 - \mu_1)] & E[(X_p - \mu_p)(X_2 - \mu_2)] & \dots & E[X_p - \mu_p]^2 \end{bmatrix} \\ &= \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{bmatrix} = \underline{\Sigma} \text{ (say)} \quad \rightarrow (3) \end{aligned}$$

Where,

$$\sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)] = \sigma_{ji}$$

clearly,  $\underline{\Sigma}$  is symmetric & positive definite matrix.

Now the multivariate normal density of  $\underline{X}$  can be obtained by replacing the positive quantity  $(x - \mu)(\sigma^2)^{-1}(x - \mu)$  by the quadratic form

$$(\underline{x} - \underline{\mu})' \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu}) \quad \rightarrow (4)$$

and is given by

$$f(\underline{\mathbf{x}}; \underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}}) = ke^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \underline{\boldsymbol{\Sigma}}^{-1}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})} \rightarrow (5)$$

Where ( $k > 0$ ) is chosen so that the integral over the entire p-dimensional Euclidean space of  $X_1, X_2, \dots, X_p$  is unity. we observe that

$$f(\underline{\mathbf{x}}; \underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}}) \geq 0 \quad (\because k \text{ is chosen as positive})$$

since  $\underline{\boldsymbol{\Sigma}}$  is positive definite

$$\begin{aligned} & (\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \underline{\boldsymbol{\Sigma}}^{-1}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}}) > 0 \\ \Rightarrow & -\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \underline{\boldsymbol{\Sigma}}^{-1}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}}) < 0 \\ \Rightarrow & e^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \underline{\boldsymbol{\Sigma}}^{-1}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})} < e^0 = 1 \end{aligned}$$

i.e.  $0 \leq f(\underline{\mathbf{x}}; \underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}}) \leq k$  i.e.  $f(\underline{\mathbf{x}})$  is bounded.

Now we should find  $k(>0)$  such that

$$\int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(x) = k \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \underline{\boldsymbol{\Sigma}}^{-1}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})} d\underline{\mathbf{x}} = 1 \rightarrow (6)$$

$$k^{-1} = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \underline{\boldsymbol{\Sigma}}^{-1}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})} d\underline{\mathbf{x}} \rightarrow (7)$$

since  $\underline{\boldsymbol{\Sigma}}^{-1}$  is positive definite  $\exists$  a non singular matrix A such that

$$\underline{\boldsymbol{\Sigma}}^{-1} = A'A \rightarrow (8)$$

then (7) can be written as

$$k^{-1} = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' A'A(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})} d\underline{\mathbf{x}} \rightarrow (9)$$



If we use the linear transformation from  $\underline{\mathbf{X}}$  to a new random vector  $\underline{\mathbf{Y}}$  such that

$$\underline{\mathbf{Y}} = A(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}}) \quad \rightarrow (10)$$

then (9) becomes

$$k^{-1} = J \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{-\frac{1}{2} \underline{\mathbf{y}}' \underline{\mathbf{y}}} d\underline{\mathbf{y}} \quad \rightarrow (11)$$

where J is the Jacobian obtained when  $\underline{\mathbf{X}}$  is transformed into  $\underline{\mathbf{Y}}$  and is given by

$$J^{-1} = \text{mod} \begin{vmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_p} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_p} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y_p}{\partial x_1} & \frac{\partial y_p}{\partial x_2} & \dots & \frac{\partial y_p}{\partial x_p} \end{vmatrix}$$

$$= \text{mod} \begin{vmatrix} a_{11} & a_{12} & \dots & a_{1p} \\ a_{21} & a_{22} & \dots & a_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pp} \end{vmatrix} = \text{mod} |A|$$

Where  $|A|$  is determinant of A.

$\therefore$  Equation (11) becomes

$$k^{-1} = \frac{1}{\text{mod} |A|} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} e^{-\frac{1}{2} \underline{\mathbf{y}}' \underline{\mathbf{y}}} d\underline{\mathbf{y}}$$

$$\begin{aligned}
&= \frac{1}{\text{mod}|A|} \prod_{i=1}^p \left( \int_{-\infty}^{\infty} e^{-\frac{1}{2}y_i^2} dy_i \right) \\
&= \frac{1}{\text{mod}|A|} \prod_{i=1}^p \sqrt{2\pi} \left( \because \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}y_i^2} dy_i = 1 \right) \\
&= \frac{(2\pi)^{p/2}}{\text{mod}|\Sigma^{-1}|^{1/2}} \quad (\because |\Sigma^{-1}| = |A'A| = |A|^2) \\
\text{i.e.,} \quad k &= \frac{1}{|\Sigma|^{1/2} (2\pi)^{p/2}} \\
&\quad (\because \Sigma^{-1} \text{ is positive definite to } \text{mod}|\Sigma^{-1}|^{1/2} = |\Sigma^{-1}|^{1/2} = \frac{1}{|\Sigma|^{1/2}})
\end{aligned}$$

substituting k in (5) we get the p.d.f of the random normal vector  $\underline{\mathbf{X}}$  and is given by

$$f(\underline{\mathbf{x}}; \underline{\boldsymbol{\mu}}, \Sigma) = \frac{1}{|\Sigma|^{1/2} (2\pi)^{p/2}} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \Sigma^{-1} (\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})} \quad \rightarrow (12)$$

thus (12) is the p.d.f of a multivariate normal vector  $\underline{\mathbf{X}}$  whose mean vector

and variance-covariance matrix are respectively given by  $\underline{\boldsymbol{\mu}}$  and  $\Sigma$

and is denoted as  $n(\underline{\mathbf{x}}; \underline{\boldsymbol{\mu}}, \Sigma)$  and its distribution is denoted as  $N_p(\underline{\boldsymbol{\mu}}, \Sigma)$ .

**NOTE 1:**

From (10), we may see that

$$\begin{aligned}
E(\underline{\mathbf{Y}}) &= AE(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}}) \\
&= A(E(\underline{\mathbf{X}}) - \underline{\boldsymbol{\mu}}) \\
&= A(\underline{\boldsymbol{\mu}} - \underline{\boldsymbol{\mu}}) \\
&= \mathbf{0}
\end{aligned}$$

i.e.  $\underline{\mathbf{Y}}$  has zero mean vector.

The variance-covariance matrix  $\underline{\mathbf{Y}}$  of is given by

$$\begin{aligned}
V(\underline{\mathbf{Y}}) &= E[(\underline{\mathbf{y}} - E(\underline{\mathbf{y}}))(\underline{\mathbf{y}} - E(\underline{\mathbf{y}}))'] \\
&= E[\underline{\mathbf{y}}\underline{\mathbf{y}}'] \quad (\because E(\underline{\mathbf{y}}) = \underline{\mathbf{0}}) \\
&= E[A(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})'A'] \\
&= AV(\underline{\mathbf{X}})A' \\
&= A\Sigma A'
\end{aligned}$$

but from (8),

$$\begin{aligned}
\Sigma &= (A'A)^{-1} \\
&= A^{-1}(A')^{-1} \quad (\because A \text{ is a non-singular}) \\
\therefore V(\underline{\mathbf{Y}}) &= AA^{-1}(A')^{-1}A' \\
&= I_k I_k \\
&= I_k
\end{aligned}$$

Thus if  $\underline{\mathbf{X}}$  is  $N_p(\underline{\boldsymbol{\mu}}, \Sigma)$ , then the random vector  $\underline{\mathbf{Y}}$  defined as

$$\underline{\mathbf{Y}} = A(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}}) \quad (\text{where } A \text{ is defined as in (8)}) \text{ follows } N_p(\underline{\mathbf{0}}, I_k).$$

In other words, the individual element of  $\underline{\mathbf{Y}}$  are standard normal variates and mutually independent i.e.  $Y_i \sim N(0,1)$  with  $\text{cov}(Y_i, Y_j) = 0$ .

#### NOTE 2.

In the practical situations 'A' can be computed as follows. since  $\Sigma$  is a symmetric p.d. matrix we may write  $\Omega'\Sigma\Omega = \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$ , when  $\Omega$  is the normalized latent vector matrix and  $\Lambda$  is the latent root matrix and since  $\Sigma$  is p.d. all  $\lambda_1, \lambda_2, \dots, \lambda_p$  are positive. Therefore  $\Lambda$  can be written  $\Lambda = (\Lambda^{1/2})'(\Lambda^{1/2})$

$$\text{Where, } \Lambda^{1/2} = \text{diag}(\sqrt{\lambda_1}, \sqrt{\lambda_2}, \dots, \sqrt{\lambda_p})$$

then

$$\begin{aligned}
\Omega'\Sigma\Omega &= (\Lambda^{1/2})'(\Lambda^{1/2}) \\
\Rightarrow \Sigma &= (\Omega^{-1})'(\Lambda^{1/2})'\Lambda^{1/2}\Omega^{-1} = A^{-1}(A^{-1})'
\end{aligned}$$

where

$$\begin{aligned}
A^{-1} &= (\Omega^{-1})'(\Lambda^{1/2})' \\
\Rightarrow A &= \Lambda^{-1/2}\Omega' \quad (\because (\Lambda^{1/2})' = \Lambda^{1/2})
\end{aligned}$$

$$\text{thus, } \underline{\mathbf{Y}} = \Lambda^{1/2}\Omega'(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})$$

$$\text{where } \Lambda^{1/2} = \text{diag}\left(\frac{1}{\sqrt{\lambda_1}}, \dots, \frac{1}{\sqrt{\lambda_p}}\right)$$

The transformation from  $\underline{\mathbf{X}}$  to  $\underline{\mathbf{Y}}$  follows  $N_p(\mathbf{0}, \mathbf{I}_k)$ . This transformation is called “whitening”.

**ALTERNATIVE METHOD OF OBTAINING THE P.D.F. OF MULTIVARIATE NORMAL VARIATE:-**

$$\text{Suppose } \underline{\mathbf{Y}} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_p \end{pmatrix}$$

is a multivariate normal vector with mean vector  $\mathbf{0}$  and variance-covariance matrix  $I_p$ .

i.e.  $E(\underline{\mathbf{Y}}) = \mathbf{0}$

i.e.  $E(y_i) = 0 \forall i$

→ (a)

and

$$\begin{aligned} V(\underline{\mathbf{Y}}) &= \begin{bmatrix} E(y_1^2) & E(y_1 y_2) & \dots & E(y_1 y_p) \\ E(y_2 y_1) & E(y_2^2) & \dots & E(y_2 y_p) \\ \vdots & \vdots & \dots & \vdots \\ E(y_p y_1) & E(y_p y_2) & \dots & E(y_p^2) \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \dots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} = I_p \end{aligned}$$

i.e.  $E(Y_i^2) = 1$  &  $E(Y_i Y_j) = 0 \quad \forall i \neq j$

→ (b)

Thus  $Y_i$ 's are i.i.d. with '0' mean and unit variance. now consider the p.d.f. of  $\underline{\mathbf{Y}}$ , which is given as

$$\begin{aligned} f(\underline{\mathbf{Y}}) &= f(y_1, y_2, \dots, y_p) \\ &= f(y_1) f(y_2) \dots f(y_p) \\ &= \prod_{i=1}^p f(y_i) \quad (\text{since } y_i \text{'s are i.i.d.}) \end{aligned}$$

$$\begin{aligned}
&= \prod_{i=1}^p \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y_i^2} \quad (\text{since } Y_i \sim N(0,1), f(y_i) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}y_i^2}) \\
&= \frac{1}{(2\pi)^{p/2}} e^{-\frac{1}{2}\sum_{j=1}^n y_j^2} = \frac{1}{(2\pi)^{p/2}} e^{-\frac{1}{2}\underline{\mathbf{y}}'\underline{\mathbf{y}}} \quad \rightarrow (1)
\end{aligned}$$

let us make use of the following linear transformation of  $\underline{\mathbf{Y}}$  into  $\underline{\mathbf{X}}$   
 $\underline{\mathbf{X}} = \mathbf{A}\underline{\mathbf{Y}} + \mathbf{b}$  where  $\mathbf{A}$  is non-singular → (2)

Now the p.d.f of  $\underline{\mathbf{X}}$  becomes

$$\begin{aligned}
f(\underline{\mathbf{x}}) &= f(x_1, x_2, \dots, x_p) \\
&= \frac{1}{(2\pi)^{p/2}} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \mathbf{b})'(\mathbf{A}^{-1})'\mathbf{A}^{-1}(\underline{\mathbf{x}} - \mathbf{b})} J(\underline{\mathbf{x}})
\end{aligned}$$

where  $J(\underline{\mathbf{x}})$  is the jacobian and is given by

$$\begin{aligned}
J(\underline{\mathbf{x}}) &= \left| \frac{d\underline{\mathbf{y}}}{d\underline{\mathbf{x}}} \right| \\
&= \begin{vmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_p} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_p} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial y_p}{\partial x_1} & \frac{\partial y_p}{\partial x_2} & \dots & \frac{\partial y_p}{\partial x_p} \end{vmatrix} \\
&= |\mathbf{A}^{-1}| \\
&= \frac{1}{|\mathbf{A}|} \quad (\because |\mathbf{A}^{-1}| = \frac{1}{|\mathbf{A}|}) \\
\therefore f(\underline{\mathbf{x}}) &= \frac{1}{|\mathbf{A}|(2\pi)^{p/2}} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \mathbf{b})'(\mathbf{A}\mathbf{A}')^{-1}(\underline{\mathbf{x}} - \mathbf{b})}
\end{aligned}$$

Which is the p.d.f. of  $\underline{\mathbf{X}}$  in terms of  $\mathbf{A}$  &  $\mathbf{b}$  and now let us interpret  $\mathbf{b}$  and  $\mathbf{A}$   
From (2),

$$\begin{aligned} E(\underline{\mathbf{X}}) &= AE(\underline{\mathbf{Y}}) + \underline{\mathbf{b}} = \underline{\mathbf{b}} \quad (\because E(\underline{\mathbf{Y}}) = \underline{\mathbf{0}}) \\ \text{i.e } \underline{\mathbf{b}} &= E(\underline{\mathbf{X}}) = \underline{\boldsymbol{\mu}} \text{ say} \end{aligned} \quad \rightarrow (4)$$

$$\begin{aligned} V(\underline{\mathbf{X}}) &= AV(\underline{\mathbf{Y}})'A' = AI_pA' = AA' \quad (\text{from } \mathbf{(b)}) \\ &= \boldsymbol{\Sigma} \text{ say} \end{aligned} \quad \rightarrow (5)$$

$\therefore \underline{\mathbf{b}}$  is the mean vector of  $\underline{\mathbf{X}}$  and  $AA'$  is the variance-covariance matrix of  $\underline{\mathbf{X}}$

using (4) & (5) in (3) we get,

$$\begin{aligned} f(\underline{\mathbf{x}}) &= \frac{1}{(2\pi)^{p/2} \text{mod}|\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{\mathbf{x}}-\underline{\boldsymbol{\mu}})'\boldsymbol{\Sigma}^{-1}(\underline{\mathbf{x}}-\underline{\boldsymbol{\mu}})} \quad \rightarrow (6) \\ & \quad (\because |AA'| = |\boldsymbol{\Sigma}| \Rightarrow |A|^2 = |\boldsymbol{\Sigma}| \Rightarrow |A| = |\boldsymbol{\Sigma}|^{1/2}) \end{aligned}$$

since  $|\boldsymbol{\Sigma}| = |A|^2 > 0$ ,  $\boldsymbol{\Sigma}$  is positive definite and therefore  $\text{mod}|\boldsymbol{\Sigma}|^{1/2} = |\boldsymbol{\Sigma}|^{1/2}$  and (6) can be written as

$$f(\underline{\mathbf{x}}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{\mathbf{x}}-\underline{\boldsymbol{\mu}})'\boldsymbol{\Sigma}^{-1}(\underline{\mathbf{x}}-\underline{\boldsymbol{\mu}})} \quad \rightarrow (7)$$

Eq (7) is the p.d.f. of the multivariate normal variate  $\underline{\mathbf{X}}$  where mean is  $\underline{\boldsymbol{\mu}}$  and variance-covariance matrix is  $\boldsymbol{\Sigma}$  and is denoted by  $n(\underline{\mathbf{x}}/\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$ .

The distribution function of  $\underline{\mathbf{X}}$  is denoted as  $N_p(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$ .

**Definition:** If  $\underline{\mathbf{X}}$  is a random vector with mean  $\underline{\boldsymbol{\mu}}$  and the variance-covariance matrix,  $\boldsymbol{\Sigma}$  and its p.d.f. is given by (7) then  $\underline{\mathbf{X}}$  is said to follow p-variate normal distribution and is denoted as  $\underline{\mathbf{X}} \sim N_p(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$ .

## PROPERTIES OF MULTIVARIATE NORMAL DISTRIBUTION

### THEOREM:

Let  $\underline{\mathbf{X}}$  (with p components) be distributed according to  $N_p(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$  then,  $\underline{\mathbf{Y}} = C\underline{\mathbf{X}}$  is distributed according to  $N(C\underline{\boldsymbol{\mu}}, C\uboldsymbol{\Sigma}C')$  for C non-singular.

### PROOF:

Since  $\underline{\mathbf{X}} \sim N_p(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$  & its p.d.f. is given as

$$f(\underline{\mathbf{x}}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})} \rightarrow (1)$$

Now, consider the linear transformation

$$\begin{aligned} \underline{\mathbf{Y}} &= C\underline{\mathbf{X}} \text{ where } C \text{ is non-singular} \\ \Rightarrow \underline{\mathbf{X}} &= C^{-1}\underline{\mathbf{Y}} \end{aligned} \rightarrow (2)$$

Now the p.d.f. 1 becomes in terms of  $\underline{\mathbf{Y}}$  as

$$g(\underline{\mathbf{y}}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(C^{-1}\underline{\mathbf{y}} - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (C^{-1}\underline{\mathbf{y}} - \underline{\boldsymbol{\mu}})} J(\underline{\mathbf{y}}) \rightarrow (3)$$

where  $J(\underline{\mathbf{y}})$  is the Jacobian and is given by

$$\begin{aligned} J(\underline{\mathbf{y}}) &= \text{mod} \left| \frac{\partial \underline{\mathbf{x}}}{\partial \underline{\mathbf{y}}} \right| = \text{mod} |C^{-1}| \\ &= \frac{1}{\text{mod} |C|} \\ &= \frac{1}{\sqrt{|C|^2}} \\ &= \sqrt{\frac{|\boldsymbol{\Sigma}|}{|C| |\boldsymbol{\Sigma}| |C'|}} \\ &= \frac{|\boldsymbol{\Sigma}|^{1/2}}{|C \boldsymbol{\Sigma} C'|^{1/2}} \end{aligned} \rightarrow (4)$$

Using (4) & (3) ,we get

$$\begin{aligned}
g(\underline{\tilde{y}}) &= \frac{1}{(2\pi)^{p/2} |\mathbf{C}\Sigma\mathbf{C}'|^{1/2}} e^{-\frac{1}{2}(\mathbf{C}^{-1}\underline{\tilde{y}}-\underline{\tilde{\mu}})'\Sigma^{-1}(\mathbf{C}^{-1}\underline{\tilde{y}}-\underline{\tilde{\mu}})} \\
&= \frac{1}{(2\pi)^{p/2} |\mathbf{C}\Sigma\mathbf{C}'|^{1/2}} e^{-\frac{1}{2}[\mathbf{C}^{-1}(\underline{\tilde{y}}-\mathbf{C}\underline{\tilde{\mu}})]'\Sigma^{-1}[\mathbf{C}^{-1}(\underline{\tilde{y}}-\mathbf{C}\underline{\tilde{\mu}})]} \\
&= \frac{1}{(2\pi)^{p/2} |\mathbf{C}\Sigma\mathbf{C}'|^{1/2}} e^{-\frac{1}{2}(\underline{\tilde{y}}-\mathbf{C}\underline{\tilde{\mu}})'(\mathbf{C}\Sigma\mathbf{C}')^{-1}(\underline{\tilde{y}}-\mathbf{C}\underline{\tilde{\mu}})} \\
&= n(\underline{\tilde{y}} / \mathbf{C}\underline{\tilde{\mu}}, \mathbf{C}\Sigma\mathbf{C}') \quad \rightarrow (5)
\end{aligned}$$

But

$$E(\underline{\mathbf{Y}}) = \mathbf{C}E(\underline{\mathbf{X}}) = \mathbf{C}\underline{\mu} \quad \rightarrow (6)$$

$$\& V(\underline{\mathbf{Y}}) = \mathbf{C}V(\underline{\mathbf{X}})\mathbf{C}' = \mathbf{C}\Sigma\mathbf{C}' \quad \rightarrow (7)$$

Now, if we write the multivariate normal p.d.f. of  $\underline{\mathbf{Y}}$  with mean  $\underline{\mu}$  and the variance-covariance matrix  $\mathbf{C}\Sigma\mathbf{C}'$  that will becomes as (5) and therefore

$$\mathbf{C}\underline{\mathbf{X}} \sim N(\mathbf{C}\underline{\mu}, \mathbf{C}\Sigma\mathbf{C}').$$

Hence the proof.

### **THEOREM:**

If a multivariate normal vector is divided into two subvectors and one sub -vector is uncorrelated with other sub-vector ,then those two sub-vectors of variables are independent and each sub-vector is also a multivariate normal vector.

**(OR)**

$$\text{Let } \underline{\mathbf{X}}_{p \times 1} \sim N_p(\underline{\mu}, \Sigma) \& \underline{\mathbf{X}} = \begin{pmatrix} \underline{\mathbf{X}}_1 \\ \underline{\mathbf{X}}_2 \end{pmatrix}$$

$$\text{Where } \underline{\mathbf{X}}_1 \text{ is } q \times 1 \text{ and } \underline{\mathbf{X}}_2 \text{ is } (p-q) \times 1 \text{ and } \underline{\mu} = \begin{pmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{pmatrix}$$



$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

where ,  $\Sigma_{11}$  is variance-covariance matrix of  $\underline{X}_1$   
 $\Sigma_{22}$  is variance-covariance matrix of  $\underline{X}_2$   
and  $\Sigma_{12}$  is covariance matrix of  $\underline{X}_1$  &  $\underline{X}_2$   
 $\Sigma_{21}$  is covariance matrix of  $\underline{X}_2$  &  $\underline{X}_1$ .

Now if  $\Sigma_{12} = \Sigma_{21}' = \mathbf{0}_{q \times p-q}$

then ,  $\underline{X}_2$  &  $\underline{X}_1$  are independent and  $\underline{X}_1 \sim N_q(\underline{\mu}_1, \Sigma_{11})$   
 $\underline{X}_2 \sim N_{p-q}(\underline{\mu}_2, \Sigma_{22})$

**PROOF:**

We are given  $\Sigma_{12} = \mathbf{0}_{q \times p-q} = \Sigma_{21}'$

i.e. the covariance matrix of  $\underline{X}_{p \times 1}$  is given by

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22} \end{pmatrix}$$

In order to show that, the random vectors  $\underline{X}_1$  &  $\underline{X}_2$  are independently normally distributed, we have to show that

$$n(\underline{x}/\underline{\mu}, \Sigma) = n(\underline{x}_1/\underline{\mu}_1, \Sigma_{11})n(\underline{x}_2/\underline{\mu}_2, \Sigma_{22})$$

we have,

$$n(\underline{x}/\underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})'\Sigma^{-1}(\underline{x}-\underline{\mu})} \rightarrow (1)$$

consider the Q.F in (1),

$$\text{i.e., } Q = (\underline{x}-\underline{\mu})'\Sigma^{-1}(\underline{x}-\underline{\mu})$$

$$\begin{aligned} &= \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}' \begin{bmatrix} \Sigma_{11} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22} \end{bmatrix}^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \\ &= \begin{bmatrix} (x_1 - \mu_1)' & (x_2 - \mu_2) \end{bmatrix}' \begin{bmatrix} \Sigma_{11}^{-1} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22}^{-1} \end{bmatrix} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix} \end{aligned}$$

( $\because \Sigma_{11}$  is the variance-covariance matrix of  $\underline{X}_1$  and hence positive definite)

$$\begin{aligned}
&= \begin{bmatrix} (\underline{x}_1 - \underline{\mu}_1)' \underline{\Sigma}_{11}^{-1} & (\underline{x}_2 - \underline{\mu}_2)' \underline{\Sigma}_{22}^{-1} \end{bmatrix} \begin{bmatrix} \underline{x}_1 - \underline{\mu}_1 \\ \underline{x}_2 - \underline{\mu}_2 \end{bmatrix} \\
&= (\underline{x}_1 - \underline{\mu}_1)' \underline{\Sigma}_{11}^{-1} (\underline{x}_1 - \underline{\mu}_1) + (\underline{x}_2 - \underline{\mu}_2)' \underline{\Sigma}_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2) \\
&= Q_1 + Q_2 \quad \rightarrow (2)
\end{aligned}$$

also we have,

$$|\underline{\Sigma}| = \begin{vmatrix} \underline{\Sigma}_{11} & 0 \\ 0 & \underline{\Sigma}_{22} \end{vmatrix} = |\underline{\Sigma}_{11}| |\underline{\Sigma}_{22}| \quad \rightarrow (3)$$

Using (2) & (3) in (1) we get ,

$$n(\underline{x}/\underline{\mu}, \underline{\Sigma}) = \frac{1}{(2\pi)^{q/2} |\underline{\Sigma}_{11}|^{1/2}} e^{-\frac{1}{2} Q_1} \frac{1}{(2\pi)^{(p-q)/2} |\underline{\Sigma}_{22}|^{1/2}} e^{-\frac{1}{2} Q_2}$$

where  $Q_1$  &  $Q_2$  are as given in (2),

$$\therefore n(\underline{x}/\underline{\mu}, \underline{\Sigma}) = n(\underline{x}_1/\underline{\mu}_1, \underline{\Sigma}_{11}) \cdot n(\underline{x}_2/\underline{\mu}_2, \underline{\Sigma}_{22})$$

Thus, the joint p.d.f. of the normal variates  $X_1, X_2, \dots, X_p$  is the product of the marginal p.d.f. of  $X_1, X_2, \dots, X_q$  and the marginal p.d.f. of  $X_{q+1}, \dots, X_p$ . Thus, the two sets of normal variates are independent.

### THEOREM:

If  $\underline{X}_1$  &  $\underline{X}_2$  are independent and are distributed as  $N_q(\underline{\mu}_1, \underline{\Sigma}_{11})$  &  $N_{p-q}(\underline{\mu}_2, \underline{\Sigma}_{22})$

respectively then,  $\begin{pmatrix} \underline{X}_1 \\ \underline{X}_2 \end{pmatrix} \sim N_p \left( \begin{bmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{bmatrix}, \begin{bmatrix} \underline{\Sigma}_{11} & 0 \\ 0 & \underline{\Sigma}_{22} \end{bmatrix} \right)$ .

### PROOF:-

we have given ,

$$\underline{X}_1 \sim N_q(\underline{\mu}_1, \underline{\Sigma}_{11})$$

$$\underline{X}_2 \sim N_{p-q}(\underline{\mu}_2, \underline{\Sigma}_{22})$$

and  $\underline{X}_1$  &  $\underline{X}_2$  are independent i.e.  $\underline{X}_1, \underline{X}_2$  are uncorrelated.

$$\text{i.e. cov}(\underline{X}_1, \underline{X}_2) = \underline{\Sigma}_{12} = \underline{\Sigma}_{21} = 0.$$

We have to find out the joint p.d.f. of  $f(\underline{\mathbf{x}})$  of  $\underline{\mathbf{X}} = \begin{pmatrix} \underline{\mathbf{X}}_1 \\ \underline{\mathbf{X}}_2 \end{pmatrix}$

we have,

$$g(\underline{\mathbf{x}}) = f(\underline{\mathbf{x}}_1)f(\underline{\mathbf{x}}_2) \quad (\because \underline{\mathbf{X}}_1 \& \underline{\mathbf{X}}_2 \text{ are independent})$$

$$\begin{aligned} &= n(\underline{\mathbf{x}}_1/\underline{\boldsymbol{\mu}}_1, \boldsymbol{\Sigma}_{11})n(\underline{\mathbf{x}}_2/\underline{\boldsymbol{\mu}}_2, \boldsymbol{\Sigma}_{22}) \\ &= \frac{1}{(2\pi)^{q/2} |\boldsymbol{\Sigma}_{11}|^{1/2}} e^{-\frac{1}{2}Q_1} \frac{1}{(2\pi)^{(p-q)/2} |\boldsymbol{\Sigma}_{22}|^{1/2}} e^{-\frac{1}{2}Q_2}, \text{ (where } Q_i = (\underline{\mathbf{x}}_i - \underline{\boldsymbol{\mu}}_i)' \boldsymbol{\Sigma}_{ii}^{-1} (\underline{\mathbf{x}}_i - \underline{\boldsymbol{\mu}}_i), i=1,2) \\ &= \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2}(Q_1+Q_2)} \left( \because |\boldsymbol{\Sigma}| = \begin{vmatrix} \boldsymbol{\Sigma}_{11} & 0 \\ 0 & \boldsymbol{\Sigma}_{22} \end{vmatrix} = |\boldsymbol{\Sigma}_{11}| |\boldsymbol{\Sigma}_{22}| \right) \end{aligned} \quad \rightarrow (1)$$

$$\text{Where } Q_1 + Q_2 = (\underline{\mathbf{x}}_1 - \underline{\boldsymbol{\mu}}_1)' \boldsymbol{\Sigma}_{11}^{-1} (\underline{\mathbf{x}}_1 - \underline{\boldsymbol{\mu}}_1) + (\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2)' \boldsymbol{\Sigma}_{22}^{-1} (\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2) \quad \rightarrow (2)$$

Let us consider

$$Q = (\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\underline{\mathbf{x}} - \underline{\boldsymbol{\mu}}) \quad \rightarrow (3)$$

where ,  $\underline{\mathbf{x}} = \begin{pmatrix} \underline{\mathbf{x}}_1 \\ \underline{\mathbf{x}}_2 \end{pmatrix}, \underline{\boldsymbol{\mu}} = \begin{pmatrix} \underline{\boldsymbol{\mu}}_1 \\ \underline{\boldsymbol{\mu}}_2 \end{pmatrix}$  is  $E(\underline{\mathbf{X}})$  and the variance-covariance matrix  $\underline{\mathbf{X}}$  is

$$\begin{aligned} \boldsymbol{\Sigma} &= \begin{pmatrix} V(\underline{\mathbf{X}}_1) & \text{cov}(\underline{\mathbf{X}}_1, \underline{\mathbf{X}}_2) \\ \text{cov}(\underline{\mathbf{X}}_2, \underline{\mathbf{X}}_1) & V(\underline{\mathbf{X}}_2) \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \end{aligned}$$

But ,since  $\boldsymbol{\Sigma}_{12} = \boldsymbol{\Sigma}'_{21} = \mathbf{0}_{q \times p-q}$

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & 0 \\ 0 & \boldsymbol{\Sigma}_{22} \end{pmatrix} \quad \rightarrow (4)$$

$$\begin{aligned} \text{Now, } Q &= \begin{pmatrix} \underline{\mathbf{x}}_1 - \underline{\boldsymbol{\mu}}_1 \\ \underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2 \end{pmatrix}' \begin{pmatrix} \boldsymbol{\Sigma}_{11}^{-1} & 0 \\ 0 & \boldsymbol{\Sigma}_{22}^{-1} \end{pmatrix} \begin{pmatrix} \underline{\mathbf{x}}_1 - \underline{\boldsymbol{\mu}}_1 \\ \underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2 \end{pmatrix} \\ &= (\underline{\mathbf{x}}_1 - \underline{\boldsymbol{\mu}}_1)' \boldsymbol{\Sigma}_{11}^{-1} (\underline{\mathbf{x}}_1 - \underline{\boldsymbol{\mu}}_1) + (\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2)' \boldsymbol{\Sigma}_{22}^{-1} (\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2) \end{aligned} \quad \rightarrow (5)$$

from (2) & (5) ,  $Q_1 + Q_2 = Q$

$$\therefore g(\underline{\mathbf{X}}) = \frac{1}{(2\pi)^{p/2} |\underline{\Sigma}|^{1/2}} e^{-\frac{1}{2}Q}$$

Where, Q is given by (3) but  $g(\underline{\mathbf{X}})$  is nothing but  $n(\underline{\mathbf{x}}/\underline{\boldsymbol{\mu}}, \underline{\Sigma})$ .

Thus  $\underline{\mathbf{X}} = \begin{pmatrix} \underline{\mathbf{X}}_1 \\ \underline{\mathbf{X}}_2 \end{pmatrix} \sim N_p(\underline{\boldsymbol{\mu}}, \underline{\Sigma})$  where,  $\underline{\Sigma}$  is as given by (4).

**THEOREM:**

If  $X_1, X_2, \dots, X_p$  have a joint normal distribution, a necessary & sufficient condition for one subset of some random variables and the subset consisting of the remaining random variables be independent is that each covariance of a variable from one set and a variable from the other set be '0'.

**PROOF:-**

**Necessary condition:-**

With out loss of generality let us assume that the first q variables form the first subset and the remaining p-q variables form the second subset.

In order to prove the necessary condition, we have given that the variables of  $X_1, X_2, \dots, X_q$  are independently distributed with the variables  $X_{q+1}, X_{q+2}, \dots, X_p$  and we have to prove

$$\text{cov}(X_i, X_j) = E[(X_i - E(X_i))(X_j - E(X_j))] = 0$$

where,  $1 \leq i \leq q$  &  $q+1 \leq j \leq p$

we have

$$\begin{aligned} \text{cov}(X_i, X_j) &= \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (x_i - E(X_i))(x_j - E(X_j))f(x_1, x_2, \dots, x_p) dx_1 dx_2 \dots dx_p \\ &= \left( \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (x_i - E(X_i))f_1(x_1 \dots x_q) dx_1 \dots dx_q \right) \\ &\quad \left( \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} (x_j - E(X_j))f_2(x_{q+1} \dots x_p) dx_{q+1} \dots dx_p \right) \\ &\quad \left( \because f(x_1, \dots, x_p) = f_1(x_1, \dots, x_q)f_2(x_{q+1}, \dots, x_p) \right) \\ &= E[X_i - E(X_i)]E[X_j - E(X_j)] \\ &= [E(X_i) - E(X_i)][E(X_j) - E(X_j)] \\ &= 0.0 \\ &= 0 \end{aligned}$$

Thus if one set of variables is independent of the remaining variables then, the set of variables are uncorrelated with the other set of variables.

**SUFFICIENT CONDITION:**

Here we have given

$$\underline{\mathbf{X}} = \begin{pmatrix} \underline{\mathbf{X}}_1 \\ \underline{\mathbf{X}}_2 \end{pmatrix} \& \underline{\mathbf{X}} : \mathbf{N}(\underline{\boldsymbol{\mu}}, \Sigma) \quad \& \text{cov}(X_i, X_j) = 0$$

where ,  $X_i$  is from  $\underline{\mathbf{X}}_1$

$X_j$  is from  $\underline{\mathbf{X}}_2$

i.e.  $\text{cov}(\underline{\mathbf{X}}_1, \underline{\mathbf{X}}_2) = \Sigma_{12} = 0_{q \times p-q}$  and we have to prove  $\underline{\mathbf{X}}_1$  &  $\underline{\mathbf{X}}_2$  are independently distributed.

The proof is already available.

**NOTE:**

To prove the necessary condition of the above theorem we need not assume  $X_1, \dots, X_p$  are normally distributed.

**THEOREM:**

If  $\underline{\mathbf{X}} \sim \mathbf{N}_p(\underline{\boldsymbol{\mu}}, \Sigma)$  and if a set of components of  $\underline{\mathbf{X}}$  is uncorrelated with the set of other components, the marginal distribution of the set is multivariate normal with means, variances and co-variances obtained by taking the proper components of  $\underline{\boldsymbol{\mu}}$  and  $\Sigma$  respectively.

**PROOF:**

Without loss of generality let us assume that the set consists of first ‘q’ components of  $\underline{\mathbf{X}}$  is uncorrelated with other components.

i.e. if  $\underline{\mathbf{X}} = \begin{pmatrix} \underline{\mathbf{X}}_1 \\ \underline{\mathbf{X}}_2 \end{pmatrix}$  where  $\underline{\mathbf{X}}_1$  is  $q \times 1$  &  $\underline{\mathbf{X}}_2$  is  $(p - q) \times 1$

$\text{cov}(\underline{\mathbf{X}}_1, \underline{\mathbf{X}}_2) = 0 = \Sigma_{12} = \Sigma'_{21}$

i.e.  $\underline{\mathbf{X}}_1$  &  $\underline{\mathbf{X}}_2$  are independent.

What all we have to prove is the marginal distribution of  $\underline{\mathbf{X}}_1$  is  $N_q(\underline{\boldsymbol{\mu}}_1, \Sigma_{11})$ .

This is already proved above.

**THEOREM:**

If  $\underline{\mathbf{X}}$  has multivariate normal distribution, then any subset of the components of  $\underline{\mathbf{X}}$  have a (multivariate) normal distribution.

(OR)

If  $\underline{X}$  is distributed as  $N_p(\underline{\mu}, \underline{\Sigma})$ , the marginal distribution of any (sub) set of components of  $\underline{X}$  is multivariate normal with means, variances and co-variances obtained by taking the proper components of  $\underline{\mu}$  and  $\underline{\Sigma}$  respectively.

**PROOF:**

$$\text{Let } \underline{X} = \begin{pmatrix} \underline{X}_1 \\ \underline{X}_2 \end{pmatrix}, \underline{\mu} = \begin{pmatrix} \underline{\mu}_1 \\ \underline{\mu}_2 \end{pmatrix}, \underline{\Sigma} = \begin{pmatrix} \underline{\Sigma}_{11} & \underline{\Sigma}_{12} \\ \underline{\Sigma}_{21} & \underline{\Sigma}_{22} \end{pmatrix}$$

where

$$\underline{\mu}_1 = E(\underline{X}_1), \underline{\mu}_2 = E(\underline{X}_2)$$

$$\underline{\Sigma}_{11} = V(\underline{X}_1), \underline{\Sigma}_{22} = V(\underline{X}_2)$$

$$\underline{\Sigma}_{12} = \text{cov}(\underline{X}_1, \underline{X}_2') = \underline{\Sigma}_{21}' = [\text{cov}(\underline{X}_2, \underline{X}_1)']$$

Now we shall make a non singular linear transformation to sub vectors

$$\underline{Y}_1 = \underline{X}_1 + \underline{M}\underline{X}_2$$

$$\underline{Y}_2 = \underline{X}_2$$

→ (1)

choosing  $\underline{M}$  so that the components of  $\underline{Y}_1$  are uncorrelated with the components of  $\underline{Y}_2 = \underline{X}_2$ . The matrix 'M' must satisfy the equation.

$$\begin{aligned} \text{cov}(\underline{Y}_1, \underline{Y}_2) &= \mathbf{0}_{q \times p-q} = E[(\underline{Y}_1 - E(\underline{Y}_1))(\underline{Y}_2 - E(\underline{Y}_2))'] \\ &= E\left[\{(\underline{X}_1 - E(\underline{X}_1)) + \underline{M}(\underline{X}_2 - E(\underline{X}_2))\}'\{\underline{X}_2 - E(\underline{X}_2)\}'\right] \\ &= E\left\{\{\underline{X}_1 - E(\underline{X}_1)\}'\{\underline{X}_2 - E(\underline{X}_2)\}'\right\} + \underline{M}E\left\{\{(\underline{X}_2 - E(\underline{X}_2))\}'\{\underline{X}_2 - E(\underline{X}_2)\}'\right\} \\ &= \text{cov}(\underline{X}_1, \underline{X}_2') + \underline{M} \text{cov}(\underline{X}_2, \underline{X}_2') \\ &= \underline{\Sigma}_{12} + \underline{M}\underline{\Sigma}_{22} \end{aligned}$$

$$\text{Thus, } \underline{M} = -\underline{\Sigma}_{12}\underline{\Sigma}_{22}^{-1} \rightarrow (2)$$

and the vector  $\underline{Y}_1$  becomes

$$\underline{Y}_1 = \underline{X}_1 - \underline{\Sigma}_{12}\underline{\Sigma}_{22}^{-1}\underline{X}_2 \rightarrow (3)$$

and the vector

$$\begin{aligned}
\tilde{\mathbf{Y}} &= \begin{pmatrix} \tilde{\mathbf{Y}}_1 \\ \tilde{\mathbf{Y}}_2 \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{X}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\tilde{\mathbf{X}}_2 \\ \tilde{\mathbf{X}}_2 \end{pmatrix} \\
&= \begin{pmatrix} I_q & -\Sigma_{12}\Sigma_{22}^{-1} \\ \mathbf{0}_{p-q \times q} & I_{p-q} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{X}}_1 \\ \tilde{\mathbf{X}}_2 \end{pmatrix} = \begin{pmatrix} I_q & -\Sigma_{12}\Sigma_{22}^{-1} \\ \mathbf{0} & I_{p-q} \end{pmatrix} \tilde{\mathbf{X}} \\
&= \mathbf{C}\tilde{\mathbf{X}}
\end{aligned}$$

Since  $\mathbf{C}$  is a non singular matrix,  $\tilde{\mathbf{Y}}$  is non-singular transformation of  $\tilde{\mathbf{X}}$  and therefore has a normal distribution with mean vector given by

$$\begin{aligned}
\tilde{\boldsymbol{\mu}} &= \begin{pmatrix} I_q & -\Sigma_{12}\Sigma_{22}^{-1} \\ \mathbf{0} & I_{p-q} \end{pmatrix} \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \\
&= \begin{pmatrix} \boldsymbol{\mu}_1 - \Sigma_{12}\Sigma_{22}^{-1}\boldsymbol{\mu}_2 \\ \boldsymbol{\mu}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{pmatrix} \quad (\text{say})
\end{aligned}$$

and the variance –covariance matrix is given by

$$\begin{aligned}
\Omega = \mathbf{V}(\tilde{\mathbf{Y}}) &= \begin{pmatrix} \mathbf{V}(\mathbf{y}_1) & \text{cov}(\mathbf{y}_1, \mathbf{y}_2) \\ \text{cov}(\mathbf{y}_2, \mathbf{y}_1) & \mathbf{V}(\mathbf{y}_2) \end{pmatrix} \\
&= \begin{pmatrix} \Sigma_{11} + \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{22}\Sigma_{22}^{-1}\Sigma_{21} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} & \Sigma_{12} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{22} \\ \Sigma_{21} - \Sigma_{21}\Sigma_{22}^{-1}\Sigma_{21} & \Sigma_{22} \end{pmatrix} \\
&= \begin{pmatrix} \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22} \end{pmatrix}
\end{aligned}$$

which implies that  $\tilde{\mathbf{Y}}_1$  &  $\tilde{\mathbf{Y}}_2$  are uncorrelated and further  $\tilde{\mathbf{Y}} = \begin{pmatrix} \tilde{\mathbf{Y}}_1 \\ \tilde{\mathbf{Y}}_2 \end{pmatrix}$  has multivariate normal distribution.

Therefore  $\tilde{\mathbf{Y}}_1$  &  $\tilde{\mathbf{Y}}_2$  are independent.

$\therefore \tilde{\mathbf{Y}}_2 = \tilde{\mathbf{X}}_2$  has the marginal distribution  $N_{p-q}(\boldsymbol{\mu}_2, \Sigma_{22})$ .

### CONDITIONAL DISTRIBUTION:-

In the above we have seen  $\tilde{\mathbf{Y}}_1 = \tilde{\mathbf{X}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\tilde{\mathbf{X}}_2$  &  $\tilde{\mathbf{Y}}_2 = \tilde{\mathbf{X}}_2$  are uncorrelated and therefore they are independently distributed.

Write down,

$$\tilde{\mathbf{Y}} = \begin{pmatrix} \tilde{\mathbf{Y}}_1 \\ \tilde{\mathbf{Y}}_2 \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{X}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\tilde{\mathbf{X}}_2 \\ \tilde{\mathbf{X}}_2 \end{pmatrix}$$

$$\begin{aligned}
&= \begin{pmatrix} I_q & -\Sigma_{12}\Sigma_{22}^{-1} \\ 0 & I_{p-q} \end{pmatrix} \begin{pmatrix} \underline{\mathbf{X}}_1 \\ \underline{\mathbf{X}}_2 \end{pmatrix} \\
&= A\underline{\mathbf{X}} \quad \rightarrow (1)
\end{aligned}$$

The joint p.d.f. of  $\underline{\mathbf{Y}}$  is  $g(\underline{\mathbf{Y}}) = g(\underline{\mathbf{Y}}_1)g(\underline{\mathbf{Y}}_2)$  ( $\because \underline{\mathbf{Y}}_1$  &  $\underline{\mathbf{Y}}_2$  are independent).

Also we know that (from the above theorem),

$$\underline{\mathbf{Y}}_1 \sim N_q(\underline{\boldsymbol{\mu}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\underline{\boldsymbol{\mu}}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21})$$

$$\underline{\mathbf{Y}}_2 \sim N_{p-q}(\underline{\boldsymbol{\mu}}_2, \Sigma_{22})$$

$$\begin{aligned}
\therefore g(\underline{\mathbf{Y}}) &= n\left(\underline{\mathbf{y}}_{111}/\underline{\boldsymbol{\mu}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\underline{\boldsymbol{\mu}}_2, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\right).n\left(\underline{\mathbf{y}}_2/\underline{\boldsymbol{\mu}}_2, \Sigma_{22}\right) \\
&= \frac{1}{(2\pi)^{q/2}|\Sigma_{11.2}|} \exp\left\{-\frac{1}{2}(\underline{\mathbf{Y}}_1 - \underline{\boldsymbol{\mu}}_1 + \Sigma_{12}\Sigma_{22}^{-1}\underline{\boldsymbol{\mu}}_2)' \Sigma_{11.2}^{-1}(\underline{\mathbf{Y}}_1 - \underline{\boldsymbol{\mu}}_1 + \Sigma_{12}\Sigma_{22}^{-1}\underline{\boldsymbol{\mu}}_2)\right\} \\
&= \frac{1}{(2\pi)^{(p-q)/2}|\Sigma_{22}|} \exp\left\{-\frac{1}{2}(\underline{\mathbf{Y}}_2 - \underline{\boldsymbol{\mu}}_2)' \Sigma_{22}^{-1}(\underline{\mathbf{Y}}_2 - \underline{\boldsymbol{\mu}}_2)\right\}
\end{aligned}$$

where,  $\Sigma_{11.2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}$ .

If we make use of the linear transformation (non singular) as given in(1).The density function of  $\underline{\mathbf{X}}$  is given by

$$f(\underline{\mathbf{x}}) = g(\underline{\mathbf{y}}(\underline{\mathbf{x}})).\mathbf{J}(\underline{\mathbf{x}})$$

Where,  $J(\underline{\mathbf{X}})$  is the jacobian and is given by

$$J(\underline{\mathbf{X}}) = \left| \frac{\partial \underline{\mathbf{y}}}{\partial \underline{\mathbf{x}}} \right| = |A| = |I_q| \cdot |I_{p-q}| = 1$$

$$\therefore f(\underline{\mathbf{X}}) = f(\underline{\mathbf{x}}_1, \underline{\mathbf{x}}_2)$$

$$\begin{aligned}
&= \frac{1}{(2\pi)^{q/2}|\Sigma_{11.2}|} e^{-\frac{1}{2}(\underline{\mathbf{x}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\underline{\boldsymbol{\mu}}_2)' \Sigma_{11.2}^{-1}(\underline{\mathbf{x}}_1 - \Sigma_{12}\Sigma_{22}^{-1}\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_1 + \Sigma_{12}\Sigma_{22}^{-1}\underline{\boldsymbol{\mu}}_2)} \\
&\quad \frac{1}{(2\pi)^{(p-q)/2}|\Sigma_{22}|} e^{-\frac{1}{2}(\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2)' \Sigma_{22}^{-1}(\underline{\mathbf{x}}_2 - \underline{\boldsymbol{\mu}}_2)} \quad \rightarrow (2)
\end{aligned}$$

Now, By the definition conditional density of  $\underline{\mathbf{X}}$ , given that  $\underline{\mathbf{X}}_2 = \underline{\mathbf{x}}_2$  is that

$$f(\underline{\mathbf{x}}_1/\underline{\mathbf{x}}_2) = \frac{f(\underline{\mathbf{x}}_1, \underline{\mathbf{x}}_2)}{f(\underline{\mathbf{x}}_2)} \quad \rightarrow (3)$$



where  $f(\underline{x}_1, \underline{x}_2)$  is given by (2) and  $f(\underline{x}_2)$  is the marginal density of  $\underline{X}_2$  at the point  $\underline{x}_2$  where is nothing but  $n(\underline{x}_2/\underline{\mu}_2, \underline{\Sigma}_{22})$ .

$$\begin{aligned} \text{i.e. } f(\underline{x}_2) &= n(\underline{x}_2/\underline{\mu}_2, \underline{\Sigma}_{22}) \\ &= \frac{1}{(2\pi)^{(p-q)/2} |\underline{\Sigma}_{22}|} \exp\left\{-\frac{1}{2}(\underline{x}_2 - \underline{\mu}_2)' \underline{\Sigma}_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2)\right\} \end{aligned}$$

Using (2) & (4) in (3) we get,

$$f(\underline{x}_1/\underline{x}_2) = \frac{1}{(2\pi)^{q/2} |\underline{\Sigma}_{11.2}|} e^{-\frac{1}{2}[(\underline{x}_1 - \underline{\mu}_1) - \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2)]' \underline{\Sigma}_{11.2}^{-1} [(\underline{x}_1 - \underline{\mu}_1) - \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2)]} \rightarrow (5)$$

which is the conditional p.d.f. of  $\underline{X}_1$  given that  $\underline{X}_2 = \underline{x}_2$ .

From (5), it is clear that the density  $f(\underline{x}_1/\underline{x}_2)$  is clearly a q-variate normal density with mean,

$$\begin{aligned} E(\underline{X}_1/\underline{X}_2 = \underline{x}_2) &= \underline{\mu}_1 + \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2) \\ &= \nu(\underline{x}_2) \text{ , say} \end{aligned} \rightarrow (6)$$

and the variances matrix,

$$\text{var}(\underline{x}_1/\underline{X}_2 = \underline{x}_2) = \underline{\Sigma}_{11.2} = \underline{\Sigma}_{11} - \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\Sigma}_{21} \rightarrow (7)$$

From (6)&(7) we may observe that the conditional mean of  $\underline{x}_1$  is simply a linear function of  $\underline{x}_2$  and the conditional co-variance of  $\underline{x}_1$  does not depend on  $\underline{x}_2$  at all.

**The above result may be put in the following theorem:-**

Let the components of  $\underline{X}$  be divided in to two groups composing the sub vectors  $\underline{X}_1$  &  $\underline{X}_2$ . Suppose the mean  $\underline{\mu}$  is similarly divided into  $\underline{\mu}_1$  &  $\underline{\mu}_2$  and suppose the co-variance matrix  $\underline{\Sigma}$  of  $\underline{X}$  is divided into  $\underline{\Sigma}_{11}, \underline{\Sigma}_{12} = \underline{\Sigma}_{21}, \underline{\Sigma}_{22}$  the co-variance matrices of  $\underline{X}_1$  of  $\underline{X}_1$  &  $\underline{X}_2$ , and of  $\underline{X}_2$  respectively. Then if the distribution of  $\underline{X}$  is normal, the conditional distribution of  $\underline{X}_1$  is given  $\underline{X}_2 = \underline{x}_2$  is normal with mean  $\underline{\mu}_1 + \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} (\underline{x}_2 - \underline{\mu}_2)$  and co-variance matrix  $\underline{\Sigma}_{11} - \underline{\Sigma}_{12} \underline{\Sigma}_{22}^{-1} \underline{\Sigma}_{21}$ .

**NOTE:-**

The above theorem may simply be asked as follows. If  $X_1, X_2, \dots, X_p$  have a joint normal distribution, then the conditional distribution of a subset of r.v's given that the remaining r.v's is also having normal distribution.

## THE CHARACTERISTIC FUNCTION:-

### DEFINITION:-

The characteristic function of a random vector  $\underline{X}$  is

$$\phi(\underline{t}) = E(e^{i\underline{t}'\underline{X}})$$

defined for every real vector  $\underline{t}$ .

### RESULT:-

If the components of a random vector  $\underline{X}$  are independently distributed,

Then,

$$\begin{aligned} E(e^{i\underline{t}'\underline{X}}) &= E \left( e^{i \sum_{j=1}^p t_j X_j} \right) \\ &= \prod_{j=1}^p E(e^{i t_j X_j}) \end{aligned}$$

### THEOREM:-

The characteristic function of  $\underline{X}$  which is distributed according to

$$N(\underline{\mu}, \underline{\Sigma}) \text{ is } \phi(\underline{t}) = E(e^{i\underline{t}'\underline{X}}) = e^{i\underline{t}'\underline{\mu} - \frac{1}{2}\underline{t}'\underline{\Sigma}\underline{t}}$$

for every real vector  $\underline{t}$ .

### PROOF:-

We have given  $\underline{X} \sim N(\underline{\mu}, \underline{\Sigma})$ .

Since,  $\underline{\Sigma}$  and hence  $\underline{\Sigma}^{-1}$  is symmetric and positive definite matrix there exists a non-singular matrix  $\underline{C}'$  such that

$$\underline{C}'\underline{\Sigma}^{-1}\underline{C} = \underline{I} \quad \rightarrow (1)$$

$$\Rightarrow \underline{\Sigma}^{-1} = (\underline{C}\underline{C}')^{-1} \text{ or } \underline{\Sigma} = \underline{C}\underline{C}' \quad \rightarrow (1.a)$$

we have the p.d.f. of  $\underline{X}$  is

$$f(\underline{X}) = \frac{1}{(2\pi)^{p/2} |\underline{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})'\underline{\Sigma}^{-1}(\underline{x}-\underline{\mu})} \quad \rightarrow (2)$$

Let us make use of the linear transformation,

$$\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}} = \mathbf{C}\underline{\mathbf{Y}} \quad (\mathbf{C} \text{ is defined as in (1)}) \quad \rightarrow (3)$$

Then the p.d.f. of the new random vector  $\underline{\mathbf{Y}}$  is

$$g(\underline{\mathbf{Y}}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}|^{1/2}} e^{-\frac{1}{2} \underline{\mathbf{y}}' \mathbf{C}' \boldsymbol{\Sigma}^{-1} \mathbf{C} \underline{\mathbf{y}}} J(\underline{\mathbf{y}}) \quad \rightarrow (4)$$

where  $J(\underline{\mathbf{y}})$  is the jacobian transformation and is given by

$$\begin{aligned} J(\underline{\mathbf{y}}) &= \text{mod} \left| \frac{\partial \underline{\mathbf{x}}}{\partial \underline{\mathbf{y}}} \right| = \text{mod} |\mathbf{C}| \left( \because \frac{\partial \underline{\mathbf{x}}}{\partial \underline{\mathbf{y}}} = \frac{\partial \mathbf{C}\underline{\mathbf{y}} + \underline{\boldsymbol{\mu}}}{\partial \underline{\mathbf{y}}} = \mathbf{C} \right) \\ &= \text{mod} |\boldsymbol{\Sigma}|^{1/2} \quad (\text{from 1(a)} |\boldsymbol{\Sigma}| = |\mathbf{C}\mathbf{C}'| = |\mathbf{C}|^2) \\ &= |\boldsymbol{\Sigma}|^{1/2} \quad (\because |\boldsymbol{\Sigma}| > 0) \end{aligned}$$

Therefore (4) becomes from (1),

$$g(\underline{\mathbf{Y}}) = \frac{1}{(2\pi)^{p/2}} e^{-\frac{1}{2} \underline{\mathbf{y}}' \underline{\mathbf{y}}} = n(\underline{\mathbf{y}} | \mathbf{0}, I_p)$$

i.e.  $\underline{\mathbf{Y}} \sim N_p(\mathbf{0}, I_p)$ .

The characteristic function of  $\underline{\mathbf{Y}}$  is

$$\begin{aligned} \phi(\underline{\mathbf{u}}) &= E(e^{i \underline{\mathbf{u}}' \underline{\mathbf{y}}}) = E \left( e^{i \sum_{j=1}^p u_j y_j} \right) \\ &= \prod_{j=1}^p E(e^{i u_j y_j}) \quad (\because Y_j \text{'s are independent}) \\ &= \prod_{j=1}^p e^{-\frac{1}{2} u_j^2} \quad (\because \text{the characteristic function of the standard} \\ &= e^{-\frac{1}{2} \sum_{j=1}^p u_j^2} \quad \text{normal variate } Y_j \text{ is } e^{-\frac{1}{2} u_j^2}) \\ &= e^{-\frac{1}{2} \underline{\mathbf{u}}' \underline{\mathbf{u}}} \quad \rightarrow (5) \end{aligned}$$

Now,

$$\phi(t) = E(e^{i t' \underline{\mathbf{X}}}) = E(e^{i t' (\mathbf{C}\underline{\mathbf{Y}} + \underline{\boldsymbol{\mu}})}) \quad (\text{from (3)})$$

i.e.

$$\begin{aligned}
\phi(\underline{\mathbf{t}}) &= E(e^{i\underline{\mathbf{t}}'\underline{\boldsymbol{\mu}}}) = e^{i\underline{\mathbf{t}}'\underline{\boldsymbol{\mu}}} E\left(e^{i\underline{\mathbf{t}}'\mathbf{C}\underline{\mathbf{Y}}}\right) \\
&= e^{i\underline{\mathbf{t}}'\underline{\boldsymbol{\mu}}} E\left(e^{i\underline{\mathbf{u}}'\underline{\mathbf{Y}}}\right) \quad (\text{where } \underline{\mathbf{u}}' = \underline{\mathbf{t}}'\mathbf{C}) \\
&= e^{i\underline{\mathbf{t}}'\underline{\boldsymbol{\mu}}} e^{-\frac{1}{2}\underline{\mathbf{u}}'\underline{\mathbf{u}}} \quad (\text{from (5)}) \\
&= e^{i\underline{\mathbf{t}}'\underline{\boldsymbol{\mu}} - \frac{1}{2}\underline{\mathbf{t}}'\mathbf{C}\mathbf{C}'\underline{\mathbf{t}}} \quad (\because \underline{\mathbf{u}}' = \underline{\mathbf{t}}'\mathbf{C}) \\
&= e^{i\underline{\mathbf{t}}'\underline{\boldsymbol{\mu}} - \frac{1}{2}\underline{\mathbf{t}}'\underline{\boldsymbol{\Sigma}}\underline{\mathbf{t}}} \quad (\text{from 1.a})
\end{aligned}$$

Hence the proof .

**THEOREM :-**

If every linear combination of the components of a vector  $\underline{\mathbf{X}}$  is normally distributed, then  $\underline{\mathbf{X}}$  is normally distributed .

**PROOF:-**

Suppose  $\underline{\mathbf{X}}$  is a random vector of  $p$  random variables with mean vector  $\underline{\boldsymbol{\mu}}$  and co-variance matrix  $\underline{\boldsymbol{\Sigma}}$  .

Let us consider a linear combination of  $\underline{\mathbf{X}}$  viz...  $\underline{\mathbf{c}}'\underline{\mathbf{X}}$  , where  $\underline{\mathbf{c}}' = (c_1, c_2, \dots, c_p)$  .

We have given ,  $\underline{\mathbf{c}}'\underline{\mathbf{X}}$  is normal variate.

We have,

$$\begin{aligned}
E(\underline{\mathbf{c}}'\underline{\mathbf{X}}) &= \underline{\mathbf{c}}'E(\underline{\mathbf{X}}) \\
&= \underline{\mathbf{c}}'\underline{\boldsymbol{\mu}} \\
V(\underline{\mathbf{c}}'\underline{\mathbf{X}}) &= V(\underline{\mathbf{c}}'\underline{\mathbf{X}}) \\
&= \underline{\mathbf{c}}'V(\underline{\mathbf{X}})\underline{\mathbf{c}} \\
&= \underline{\mathbf{c}}'\underline{\boldsymbol{\Sigma}}\underline{\mathbf{c}} \quad (\text{variance})
\end{aligned}$$

It may be noted  $\underline{\mathbf{c}}'\underline{\boldsymbol{\mu}}$  &  $\underline{\mathbf{c}}'\underline{\boldsymbol{\Sigma}}\underline{\mathbf{c}}$  are scalars and they are respectively the mean & variances of the univariate random variables  $\underline{\mathbf{c}}'\underline{\mathbf{X}}$  . We have given that  $\underline{\mathbf{c}}'\underline{\mathbf{X}} \sim N(\underline{\mathbf{c}}'\underline{\boldsymbol{\mu}}, \underline{\mathbf{c}}'\underline{\boldsymbol{\Sigma}}\underline{\mathbf{c}})$  .

Let  $Y = \underline{\mathbf{c}}'\underline{\mathbf{X}}$  and from the univariate normal distribution theory.

The characteristic function of  $Y$  is given by

$$\psi(t) = E(e^{itY})$$

$$\begin{aligned}
&= e^{itE(Y) - \frac{1}{2}t^2V(Y)} \\
&= e^{it\tilde{\mathbf{c}}'\tilde{\boldsymbol{\mu}} - \frac{1}{2}t^2\tilde{\mathbf{c}}'\tilde{\boldsymbol{\Sigma}}\tilde{\mathbf{c}}}
\end{aligned}$$

If we write  $t=1$  then,  $\psi(t)$  becomes

$$\phi(\tilde{\mathbf{c}}) = e^{i\tilde{\mathbf{c}}'\tilde{\boldsymbol{\mu}} - \frac{1}{2}\tilde{\mathbf{c}}'\tilde{\boldsymbol{\Sigma}}\tilde{\mathbf{c}}} \quad \text{where, } \tilde{\mathbf{X}} \sim N(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$$

which is the characteristic function of a multivariate random vector  $\tilde{\mathbf{X}}$  whose mean vector is  $\tilde{\boldsymbol{\mu}}$  & variance-covariance matrix is of  $\tilde{\boldsymbol{\Sigma}}$ .

But the mean & variance-covariance matrix of  $\tilde{\mathbf{X}}$  respectively same as  $\tilde{\boldsymbol{\mu}}$  &  $\tilde{\boldsymbol{\Sigma}}$  and therefore,  $\tilde{\mathbf{X}} \sim N(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ .

Hence the proof.

**THEOREM :-**

If  $\tilde{\mathbf{X}} \sim N_p(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$  then,  $\tilde{\mathbf{c}}'\tilde{\mathbf{X}}$  is uni-normal variate with mean  $\tilde{\mathbf{c}}'\tilde{\boldsymbol{\mu}}$  and variance  $\tilde{\mathbf{c}}'\tilde{\boldsymbol{\Sigma}}\tilde{\mathbf{c}}$ .

(OR)

If  $X_1, X_2, \dots, X_p$  are jointly distributed as p-variate normal then its linear combination follows uni-variate normal distribution.

**PROOF:-**

$$\text{Let } \tilde{\mathbf{X}} = \begin{pmatrix} X_1 \\ \vdots \\ X_p \end{pmatrix} \sim N_p(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$$

Then its characteristic function is given by

$$\begin{aligned}
\phi(\tilde{\mathbf{t}}) &= E(e^{i\tilde{\mathbf{t}}'\tilde{\mathbf{X}}}) \\
&= e^{i\tilde{\mathbf{t}}'\tilde{\boldsymbol{\mu}} - \frac{1}{2}\tilde{\mathbf{t}}'\tilde{\boldsymbol{\Sigma}}\tilde{\mathbf{t}}} \quad \rightarrow (1)
\end{aligned}$$

$$\text{Let us write } \tilde{\mathbf{t}} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_p \end{pmatrix}$$

$$= \begin{pmatrix} t.c_1 \\ t.c_2 \\ \vdots \\ t.c_p \end{pmatrix} \\ = t\underset{\sim}{\mathbf{c}}$$

Then (1) becomes ,

$$\begin{aligned} \phi(\mathbf{t}) &= e^{it\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\mu}} - \frac{1}{2}t\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\Sigma}}t} \\ &= e^{it\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\mu}} - \frac{1}{2}t^2\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\Sigma}}\underset{\sim}{\mathbf{c}}} \\ &= E(e^{itY}) \\ &= \psi(t) \quad , \text{say} \end{aligned} \quad \rightarrow (2)$$

where Y is normal variate with mean  $\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\mu}}$  and variance  $\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\Sigma}}\underset{\sim}{\mathbf{c}}$  .

In other words (2) is the characteristic function  $\psi(t)$  of a uni-normal variate whose mean is  $\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\mu}}$  and variance is  $\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\Sigma}}\underset{\sim}{\mathbf{c}}$  .

If we consider the linear combination of the components of the normal random vector  $\underset{\sim}{\mathbf{X}}$  viz.,

$$\begin{aligned} Y &= \underset{\sim}{\mathbf{c}}'\underset{\sim}{\mathbf{X}} \\ &= c_1X_1 + c_2X_2 + \dots + c_pX_p \end{aligned}$$

its mean and variance are given by

$$\begin{aligned} E(Y) &= E(\underset{\sim}{\mathbf{c}}'\underset{\sim}{\mathbf{X}}) = \underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\mu}} \quad \& \quad V(Y) = V(\underset{\sim}{\mathbf{c}}'\underset{\sim}{\mathbf{X}}) \\ &= \underset{\sim}{\mathbf{c}}'V(\underset{\sim}{\mathbf{X}})\underset{\sim}{\mathbf{c}} \\ &= \underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\Sigma}}\underset{\sim}{\mathbf{c}} \end{aligned}$$

Thus, from the above explanation it follows that  $Y = \underset{\sim}{\mathbf{c}}'\underset{\sim}{\mathbf{X}}$  follows uni-variate normal distribution,

$$\text{i.e. , } Y = \underset{\sim}{\mathbf{c}}'\underset{\sim}{\mathbf{X}} \sim N(\underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\mu}}, \underset{\sim}{\mathbf{c}}'\underset{\sim}{\boldsymbol{\Sigma}}\underset{\sim}{\mathbf{c}})$$

Hence the proof.

## SAMPLING FROM MULTINORMAL DISTRIBUTION

The Multivariate normal likelihood:

Let us assume that the  $p \times 1$  vectors  $\underline{\mathbf{X}}_1, \underline{\mathbf{X}}_2, \dots, \underline{\mathbf{X}}_n$  represent a random sample from a multivariate normal population with mean vector  $\underline{\boldsymbol{\mu}}$  and covariance matrix  $\boldsymbol{\Sigma}$ . Since  $\underline{\mathbf{X}}_1, \underline{\mathbf{X}}_2, \dots, \underline{\mathbf{X}}_n$  are mutually independent (by virtue of randomization) and each has distributed as the joint p.d.f. of all the observations is the product of the marginal normal densities.

i.e.

$$f(\underline{\mathbf{x}}_1, \underline{\mathbf{x}}_2, \dots, \underline{\mathbf{x}}_n) = f(\underline{\mathbf{x}}_1) \cdot \dots \cdot f(\underline{\mathbf{x}}_n)$$

$$\begin{aligned}
 &= \prod_{j=1}^n \left\{ \frac{1}{(2\pi)^{\frac{p}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} e^{-\frac{1}{2}(\underline{\mathbf{x}}_j - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\underline{\mathbf{x}}_j - \underline{\boldsymbol{\mu}})} \right\} \\
 &\quad \left[ \because f(\underline{\mathbf{x}}_j) = n(\underline{\mathbf{x}}_j / \underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}) \right] \\
 &= \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2}} e^{-\frac{1}{2} \sum_{j=1}^n (\underline{\mathbf{x}}_j - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\underline{\mathbf{x}}_j - \underline{\boldsymbol{\mu}})} \quad \text{--- (1)}
 \end{aligned}$$

When the numerical value of the observations become available, there may be substituted for  $\underline{\mathbf{X}}_j$  in equation (1). The resulting expression, now considered as a function of  $\underline{\boldsymbol{\mu}}$  and  $\boldsymbol{\Sigma}$  and for a fixed set of observations  $\underline{\mathbf{x}}_1, \underline{\mathbf{x}}_2, \dots, \underline{\mathbf{x}}_n$ , it is called as “the likelihood function” and is denoted as  $L(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$ .

**ML ESTIMATION OF MEAN VECTOR  $\underline{\mu}$  AND VARIANCE – COVARIANCE MATRIX  $\Sigma$  :**

Consider the likelihood function  $L(\underline{\mu}, \Sigma)$  given by (1). i.e.

$$L(\underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{np/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \sum_{j=1}^n (\underline{x}_j - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \underline{\mu})} \quad \text{--- (2)}$$

Now the maximum likelihood estimates of  $\underline{\mu}$  and  $\Sigma$  can be obtained by maximizing  $L(\underline{\mu}, \Sigma)$ .

In order to obtain the MLE's of  $\underline{\mu}$  and  $\Sigma$ , let us consider logarithms of (2) and is given by

$$\log L(\underline{\mu}, \Sigma) = \frac{-np}{2} \log(2\pi) - \frac{n}{2} \log |\Sigma| - \frac{1}{2} \sum_{j=1}^n (\underline{x}_j - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \underline{\mu}) \quad \text{-->(3)}$$

Consider the last term of (3) and as if is a scalar we may write if as

$$\begin{aligned} & \sum_{j=1}^n (\underline{x}_j - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \underline{\mu}) \\ &= tr \left[ \sum_{j=1}^n (\underline{x}_j - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \underline{\mu}) \right] \\ &= \sum_{j=1}^n tr \left[ (\underline{x}_j - \underline{\mu})' \Sigma^{-1} (\underline{x}_j - \underline{\mu}) \right] \\ &= \sum_{j=1}^n tr \left[ \Sigma^{-1} (\underline{x}_j - \underline{\mu}) (\underline{x}_j - \underline{\mu})' \right] \\ & \quad (\because tr(\mathbf{A B}) = tr(\mathbf{B A})) \end{aligned}$$



$$= tr \left[ \Sigma^{-1} \left\{ \sum_{j=1}^n (\mathbf{x}_j - \underline{\boldsymbol{\mu}})(\mathbf{x}_j - \underline{\boldsymbol{\mu}})' \right\} \right] \quad \rightarrow (4)$$

Now consider

$$\sum_{j=1}^n (\mathbf{x}_j - \underline{\boldsymbol{\mu}})(\mathbf{x}_j - \underline{\boldsymbol{\mu}})'$$

$$= \sum_{j=1}^n \left[ \mathbf{x}_j - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \underline{\boldsymbol{\mu}} \right] \left[ \mathbf{x}_j - \bar{\mathbf{x}} + \bar{\mathbf{x}} - \underline{\boldsymbol{\mu}} \right]'$$

$$\text{Where } \bar{\mathbf{x}} = \frac{1}{n} (\mathbf{x}_1 + \mathbf{x}_2 + \cdots + \mathbf{x}_n)$$

$$= \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' + \sum_{j=1}^n (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})'$$

(Since the cross product terms

$$\sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})' = (n\bar{\mathbf{x}} - n\bar{\mathbf{x}})(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})' = 0$$

$$\left[ \because \sum_{j=1}^n \mathbf{x}_j = n\bar{\mathbf{x}} \right]$$

and similarly

$$\sum_{j=1}^n (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})(\mathbf{x}_j - \bar{\mathbf{x}})' = 0$$

Thus

$$\begin{aligned} & \sum_{j=1}^n (\mathbf{x}_j - \underline{\boldsymbol{\mu}})(\mathbf{x}_j - \underline{\boldsymbol{\mu}})' \\ &= \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' + n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \end{aligned}$$

Substituting this in (4) we get ,

$$\begin{aligned} & \sum_{j=1}^n (\mathbf{x}_j - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_j - \underline{\boldsymbol{\mu}}) \\ &= \text{tr} \left[ \boldsymbol{\Sigma}^{-1} \left\{ \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' + n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \right\} \right] \\ &= \text{tr} \left[ \boldsymbol{\Sigma}^{-1} \left\{ \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right\} \right] + n \left[ \text{tr} \left\{ (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}) \right\} \right] \end{aligned}$$

[Since  $\text{tr}(AB) = \text{tr}(BA)$ ]

$$= \text{tr} \left[ \boldsymbol{\Sigma}^{-1} \left\{ \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right\} \right] + n \left\{ (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}) \right\}$$

--- (5)

(Since trace of scalar is scalar)

Substituting (5) in (3),we get

$$\log L(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}) = \frac{-np}{2} \log(2\pi) + \frac{n}{2} \log |\boldsymbol{\Sigma}^{-1}|$$

$$\begin{aligned}
& -\frac{1}{2} \left[ \text{tr} \left\{ \Sigma^{-1} \left( \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right) \right\} \right] \\
& -\frac{n}{2} (\bar{\mathbf{x}} - \underline{\underline{\boldsymbol{\mu}}})' \Sigma^{-1} (\bar{\mathbf{x}} - \underline{\underline{\boldsymbol{\mu}}}) \quad \text{--- (6)}
\end{aligned}$$

Since  $\Sigma^{-1}$  is positive definite

$$\begin{aligned}
(\bar{\mathbf{x}} - \underline{\underline{\boldsymbol{\mu}}})' \Sigma^{-1} (\bar{\mathbf{x}} - \underline{\underline{\boldsymbol{\mu}}}) &> 0 \quad \forall \underline{\underline{\boldsymbol{\mu}}} \neq \bar{\mathbf{x}} \\
&= 0 \quad \text{if } \underline{\underline{\boldsymbol{\mu}}} = \bar{\mathbf{x}}
\end{aligned}$$

From (6), we can observe that if the last term is zero then (6) becomes maximum that is

$\log L(\underline{\underline{\boldsymbol{\mu}}}, \Sigma)$  can be maximized with respect to  $\underline{\underline{\boldsymbol{\mu}}}$  at  $\hat{\underline{\underline{\boldsymbol{\mu}}}} = \bar{\mathbf{x}}$

$\therefore$  The MLE of  $\underline{\underline{\boldsymbol{\mu}}}$  is  $\bar{\mathbf{x}}$  substituting the MLE of  $\underline{\underline{\boldsymbol{\mu}}}$  ( $\bar{\mathbf{x}}$ ) in (6)

We get

$$\begin{aligned}
\log L(\underline{\underline{\boldsymbol{\mu}}}, \Sigma) &= \frac{-np}{2} \log(2\pi) + \frac{n}{2} \log |\Sigma^{-1}| \\
& -\frac{1}{2} \left[ \text{tr} \left\{ \Sigma^{-1} \left( \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right) \right\} \right] \\
& \text{--- (7)}
\end{aligned}$$

Now we have to maximize (7) w.r.t.  $\Sigma$  as the equation is free of  $\underline{\underline{\boldsymbol{\mu}}}$

We can prove that (7) attains its maximum value at  $\Sigma = \hat{\Sigma}$ ,

$$\text{Where, } \hat{\Sigma} = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \quad \text{---(8)}$$

Thus  $\hat{\Sigma}$  (given by (8)) is the MLE of  $\Sigma$ .

The maximum value of the likelihood can be obtained by substituting the MLEs of  $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$  respectively given by

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$$

$$\text{and } \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})'$$

in (2) and it is given by

$$\begin{aligned} L(\hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\Sigma}}) &= (2\pi)^{-np/2} |\hat{\boldsymbol{\Sigma}}|^{-n/2} \\ &e^{-\frac{1}{2} \text{tr} \left[ n \left( \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right)^{-1} \left( \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}})(\mathbf{x}_j - \bar{\mathbf{x}})' \right) \right]} \\ &= (2\pi)^{-np/2} |\hat{\boldsymbol{\Sigma}}|^{-n/2} e^{-\frac{n}{2} \text{tr}(\mathbf{I}_p)} \\ &= (2\pi)^{-np/2} |\hat{\boldsymbol{\Sigma}}|^{-n/2} e^{-np/2} \\ &= \text{const. } x |\hat{\boldsymbol{\Sigma}}|^{-\frac{n}{2}} \\ &= \text{const. } x (\text{generalised variance})^{\frac{-n}{2}} \end{aligned}$$

since generalized variance is defined as  $|\hat{\boldsymbol{\Sigma}}|$ . The generalized variance determines the *peakedness* of the likelihood function and consequently is a natural measure of variability when the parent population is multivariate normal.

**NOTE:-**

1. MLE's possess an invariance property which means if  $\hat{\theta}$  is the MLE of  $\theta$  then  $h(\hat{\theta})$  is the MLE of  $h(\theta)$ , where  $h(\theta)$  is a function of  $\theta$ .

**For Example :-**

i. If  $\hat{\mu}$  is MLE of  $\mu$  and  $\hat{\Sigma}$  is MLE of  $\Sigma$  then  $\hat{\mu}^{-1}\hat{\Sigma}^{-1}\hat{\mu}$  is MLE of

$$\mu' \Sigma^{-1} \mu .$$

ii. If  $\sigma_{ij}$  is the  $ij^{th}$  element of  $\Sigma$  and  $\hat{\sigma}_{ij}$  is the  $ij^{th}$  element of

$\hat{\Sigma}$  where  $\hat{\Sigma}$  is the MLE of  $\Sigma$ .

$$\text{Where } \hat{\sigma}_{ij} = \frac{1}{n} \sum_{k=1}^n (X_{ik} - \bar{X}_i)(X_{jk} - \bar{X}_j)$$

$$= COV(X_i, X_j).$$

2. From equation (6) the log-likelihood and hence the joint p.d.f depends on the whole set of observations  $x_1, \dots, x_n$  only through the sample mean  $\bar{x}$  and the sum of squares and cross product matrix,

$$\sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}}) (\mathbf{x}_j - \bar{\mathbf{x}})' = n \hat{\Sigma}$$

We may express this fact by saying that  $\hat{\mu}$  (or  $\bar{x}$ ) and  $\hat{\Sigma}$  are sufficient statistics.

Thus the MLEs  $\hat{\mu}$  and  $\hat{\Sigma}$  are sufficient statistics of  $\mu$  and  $\Sigma$ .

3. The MLE of  $\Sigma$  is

$$\hat{\Sigma} = \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}}) (\mathbf{x}_j - \bar{\mathbf{x}})'$$

This formula is not convenient to compute  $\hat{\Sigma}$  and the following is the convenient formula for computation

$$\hat{\Sigma} = \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j' - \bar{\mathbf{x}} \bar{\mathbf{x}}'$$

Explanation:-

$$\begin{aligned} \hat{\Sigma} &= \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \bar{\mathbf{x}}) (\mathbf{x}_j - \bar{\mathbf{x}})' \\ &= \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j' - \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \bar{\mathbf{x}}' - \frac{1}{n} \sum_{j=1}^n \bar{\mathbf{x}} \mathbf{x}_j' + \frac{1}{n} \sum_{j=1}^n \bar{\mathbf{x}} \bar{\mathbf{x}}' \\ &= \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j' - \bar{\mathbf{x}} \bar{\mathbf{x}}' - \bar{\mathbf{x}} \bar{\mathbf{x}}' + \bar{\mathbf{x}} \bar{\mathbf{x}}' \\ \hat{\Sigma} &= \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j \mathbf{x}_j' - \bar{\mathbf{x}} \bar{\mathbf{x}}' \end{aligned}$$

Sampling Distribution of the MLE's  $\hat{\mu}$  and  $\hat{\Sigma}$  and the independence of  $\hat{\mu}$  and

$\hat{\Sigma}$  :- Before going to obtain the sampling distribution of  $\hat{\mu}$  and  $\hat{\Sigma}$ , let us prove

the following result which is useful in obtaining the sampling distributions of  $\hat{\underline{\mu}}$  and  $\hat{\underline{\Sigma}}$ .

**Result:-**

Suppose  $\underline{X}_1, \dots, \underline{X}_n$  are independent where  $\underline{X}_a \sim Np(\underline{\mu}_a, \underline{\Sigma})$ .

Let

$\underline{C} = (c_{\alpha j})_{n \times n}$  be an orthogonal matrix then

$$\underline{Y}_a = \sum_{j=1}^n c_{\alpha j} \underline{X}_j \sim Np(\underline{v}_a, \underline{\Sigma})$$

Where  $\underline{v}_a = \sum_{j=1}^n c_{\alpha j} \underline{\mu}_j$  and

$\underline{Y}_1, \dots, \underline{Y}_n$  are independent.

**Proof:-** Since 'C' is orthogonal matrix,

$$\text{We have } \sum_{j=1}^n c_{\alpha j}^2 = 1 \quad \text{and} \quad \sum_{j=1}^n c_{\alpha j} c_{\beta j} = 0 \quad \text{---- (1)}$$

In order to prove  $\underline{Y}_a \sim Np(\underline{v}_a, \underline{\Sigma})$ , Let us consider the characteristic

function of  $\underline{Y}_a$  is

$$E \left( e^{i \underline{t}' \underline{Y}_a} \right) = E \left( e^{i \underline{t}' \sum_{j=1}^n c_{\alpha j} \underline{X}_j} \right)$$

$$= E \left( \prod_{j=1}^n e^{i C_{\alpha j} \mathbf{t}' \mathbf{X}_j} \right)$$

$$= \prod_{j=1}^n E \left( e^{i \mathbf{u}' \mathbf{X}_j} \right)$$

(Since  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are independent and where  $\mathbf{u} = C_{\alpha j} \mathbf{t}$ )

$$= \prod_{j=1}^n e^{i \mathbf{u}' \boldsymbol{\mu}_j - \frac{1}{2} \mathbf{u}' \boldsymbol{\Sigma} \mathbf{u}}$$

$$\left( \because \mathbf{X}_j \sim Np(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}) \right)$$

$$= e^{i \sum_{j=1}^n \mathbf{u}' \boldsymbol{\mu}_j - \frac{1}{2} \sum_{j=1}^n \mathbf{u}' \boldsymbol{\Sigma} \mathbf{u}}$$

$$= e^{i \mathbf{t}' \sum_{j=1}^n C_{\alpha j} \boldsymbol{\mu}_j - \frac{1}{2} \mathbf{t}' \boldsymbol{\Sigma} \mathbf{t} \sum_{j=1}^n C_{\alpha j}^2}$$

(Substituting  $\mathbf{u} = C_{\alpha j} \mathbf{t}$ )

$$= e^{i \mathbf{t}' \sum_{j=1}^n C_{\alpha j} \boldsymbol{\mu}_j - \frac{1}{2} \mathbf{t}' \boldsymbol{\Sigma} \mathbf{t}}$$

[Using (1)]

Which is the characteristic function of a multivariate normal vector whose mean vector



$\sum_{j=1}^n C_{\alpha j} \boldsymbol{\mu}_{\tilde{j}}$  & variance-covariance matrix  $\Sigma$ .

Thus  $\mathbf{Y}_{\tilde{\alpha}} \sim Np(\boldsymbol{\nu}_{\alpha}, \Sigma)$ , Where  $\boldsymbol{\nu}_{\alpha} = \sum_{j=1}^n C_{\alpha j} \boldsymbol{\mu}_{\tilde{j}}$

Now if remains to prove that  $\mathbf{Y}_{\tilde{1}}, \dots, \mathbf{Y}_{\tilde{n}}$  are independent.

In order to prove  $\mathbf{Y}_{\tilde{1}}, \dots, \mathbf{Y}_{\tilde{n}}$  are independent we have to prove that (as  $\mathbf{Y}_{\tilde{j}}$ 's are multivariate normal vector).

$$\text{COV}(\mathbf{Y}_{\tilde{\alpha}}, \mathbf{Y}_{\tilde{\beta}}) = \mathbf{O}_{p \times p} \quad (\text{Zero matrix})$$

The covariance matrix between  $\mathbf{Y}_{\tilde{\alpha}}$  and  $\mathbf{Y}_{\tilde{\beta}}$  is

$$\begin{aligned} \text{cov}(\mathbf{Y}_{\tilde{\alpha}}, \mathbf{Y}_{\tilde{\beta}}) &= E\left[(\mathbf{Y}_{\tilde{\alpha}} - E(\mathbf{Y}_{\tilde{\alpha}}))(\mathbf{Y}_{\tilde{\beta}} - E(\mathbf{Y}_{\tilde{\beta}}))'\right] \\ &= E\left[(\mathbf{Y}_{\tilde{\alpha}} - \boldsymbol{\nu}_{\alpha})(\mathbf{Y}_{\tilde{\beta}} - \boldsymbol{\nu}_{\beta})'\right] \\ &= E\left[\left(\sum_{i=1}^n C_{\alpha i} \mathbf{X}_{\tilde{i}} - \sum_{i=1}^n C_{\alpha i} \boldsymbol{\mu}_{\tilde{i}}\right)\left(\sum_{j=1}^n C_{\beta j} \mathbf{X}_{\tilde{j}} - \sum_{j=1}^n C_{\beta j} \boldsymbol{\mu}_{\tilde{j}}\right)'\right] \end{aligned}$$

$$\begin{aligned}
&= E \left[ \left( \sum_{i=1}^n C_{\alpha i} (\mathbf{X}_{\tilde{i}} - \boldsymbol{\mu}_{\tilde{i}}) \sum_{j=1}^n C_{\beta j} (\mathbf{X}_{\tilde{j}} - \boldsymbol{\mu}_{\tilde{j}}) \right)' \right] \\
&= \sum_{i=1}^n \sum_{j=1}^n C_{\alpha i} C_{\beta j} E \left[ (\mathbf{X}_{\tilde{i}} - \boldsymbol{\mu}_{\tilde{i}}) (\mathbf{X}_{\tilde{j}} - \boldsymbol{\mu}_{\tilde{j}})' \right] \\
&= \sum_{i=1}^n \sum_{j=1}^n C_{\alpha i} C_{\beta j} \text{cov}(\mathbf{X}_{\tilde{i}}, \mathbf{X}_{\tilde{j}}') \\
&= \sum_{i=1}^n c_{\alpha i} c_{\beta i} V(X_{\tilde{i}}) + \sum_{i=1}^n \sum_{j=1}^n c_{\alpha i} c_{\beta j} \text{cov}(X_{\tilde{i}}, X_{\tilde{j}})
\end{aligned}$$

since  $X_{\tilde{i}}$ 's are independent  $\text{cov}(X_{\tilde{i}}, X_{\tilde{j}}) = 0$

Therefore the second term of the above will vanish.

Therefore now

$$\text{cov}(Y_{\tilde{\alpha}}, Y_{\tilde{\beta}}') = \sum_{i=1}^n c_{\alpha i} c_{\beta i} V(X_{\tilde{i}})$$

$$= \sum_{i=1}^n \sum_{i=1}^n c_{\alpha i} c_{\beta i}$$

$$(\because V(X_{\tilde{i}}) = \sum)$$

But from equation (1)

$$\sum_{i=1}^n c_{\alpha i} c_{\beta i} = 0$$

$$\therefore \text{cov}(Y_{\tilde{\alpha}}, Y_{\tilde{\beta}}') = 0$$

Thus  $\mathbf{Y}_{\tilde{\alpha}}$  and  $\mathbf{Y}_{\tilde{\beta}}$  are independent and consequently  $\mathbf{Y}_{\tilde{1}}, \dots, \mathbf{Y}_{\tilde{n}}$  are independent. Hence the theorem.

**Theorem:-**

Let  $\underline{\mathbf{X}}_1, \dots, \underline{\mathbf{X}}_n$  be an independent random sample from  $N_p(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$ .

Then the MLE of  $\underline{\boldsymbol{\mu}}$  say  $\hat{\underline{\boldsymbol{\mu}}}$  (also the sample mean) is distributed according to  $N_p(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}/n)$  and is independent of the MLE of  $\boldsymbol{\Sigma}$  given by

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{\alpha=1}^n (\underline{\mathbf{X}}_{\alpha} - \bar{\underline{\mathbf{X}}})(\underline{\mathbf{X}}_{\alpha} - \bar{\underline{\mathbf{X}}})'$$

and  $n\hat{\boldsymbol{\Sigma}}$  is distributed as  $\sum_{\alpha=1}^n \underline{\mathbf{z}}_{\alpha} \underline{\mathbf{z}}_{\alpha}'$  where  $\underline{\mathbf{z}}_{\alpha} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma})$  and is  $\underline{\mathbf{z}}_1, \dots, \underline{\mathbf{z}}_n$  are independent.

**Proof:-**

We have given a random sample  $\underline{\mathbf{X}}_1, \dots, \underline{\mathbf{X}}_n$  where  $\underline{\mathbf{X}}_{\alpha} \sim N_p(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$

and is independent of  $\underline{\mathbf{X}}_{\beta}$  for  $\alpha \neq \beta$ . We have the MLE's of  $\underline{\boldsymbol{\mu}}$  and  $\boldsymbol{\Sigma}$  are respectively given by

$$\hat{\underline{\boldsymbol{\mu}}} = 1/n \sum_{\alpha=1}^n \underline{\mathbf{X}}_{\alpha} = \bar{\underline{\mathbf{X}}}$$

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{\alpha=1}^n (\underline{\mathbf{X}}_{\alpha} - \bar{\underline{\mathbf{X}}})(\underline{\mathbf{X}}_{\alpha} - \bar{\underline{\mathbf{X}}})' \quad \dots$$

(1)

Now there exists an  $n \times n$  orthogonal matrix  $\mathbf{B} = (b_{\alpha\beta})$  with the last row i.e.

$$b_{n\beta} = \frac{1}{\sqrt{n}} \quad \forall \beta$$

$$\left( \frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}} \right) \quad \text{--- (1.a)}$$

Let us define a new random sample  $\underline{\tilde{z}}_1, \dots, \underline{\tilde{z}}_n$  from the given random sample  $\underline{X}_1, \dots, \underline{X}_n$  using the orthogonal transformation from the orthogonal matrix  $\mathbf{B}$ . Thus

$$\underline{\tilde{z}}_\alpha = \sum_{\beta=1}^n b_{\alpha\beta} \underline{X}_\beta \quad \text{for } \alpha = 1, 2, \dots, n \quad \text{---}$$

(2)  
In particular,

$$\underline{\tilde{z}}_n = \sum_{\beta} b_{n\beta} \underline{X}_\beta$$

$$= \sum_{\beta} \frac{1}{\sqrt{n}} \underline{X}_\beta \quad \text{[The last row of } \mathbf{B} \text{ is as given in (1.a)]}$$

$$= \sqrt{n} \bar{\underline{X}} \quad \text{[From (1)]} \quad \text{--- (3)}$$

Let us consider

$$\sum_{\alpha=1}^n \underline{\tilde{z}}_\alpha \underline{\tilde{z}}_\alpha' = \sum_{\alpha=1}^n \left( \sum_{i=1}^n b_{\alpha i} \underline{X}_i \right) \left( \sum_{j=1}^n b_{\alpha j} \underline{X}_j \right)' \quad \text{[ Using (2) ]}$$

i.e.

$$\begin{aligned}
&= \sum_{\alpha=1}^n \sum_{i=1}^n \sum_{j=1}^n b_{\alpha i} b_{\alpha j} \underline{\underline{X}}_i \underline{\underline{X}}_j' \\
&= \sum_{i=j=1}^n \sum_{\alpha=1}^n b_{\alpha i}^2 \underline{\underline{X}}_i \underline{\underline{X}}_i' + \sum_{i \neq j=1}^n \sum_{\alpha=1}^n b_{\alpha i} b_{\alpha j} \underline{\underline{X}}_i \underline{\underline{X}}_j' \\
&= \sum_{i=1}^n \underline{\underline{X}}_i \underline{\underline{X}}_i' \left( \sum_{\alpha=1}^n b_{\alpha i}^2 \right) + \sum_{i \neq j=1}^n \underline{\underline{X}}_i \underline{\underline{X}}_j' \left( \sum_{\alpha=1}^n b_{\alpha i} b_{\alpha j} \right) \\
&= \sum_{\alpha=1}^n \underline{\underline{X}}_{\alpha} \underline{\underline{X}}_{\alpha}' \quad \text{--- (4)}
\end{aligned}$$

[ $\because$   $\mathbf{B}$  is the orthogonal matrix and as a consequence

$$\sum_{\alpha=1}^n b_{\alpha i} b_{\alpha j} = 0 \quad \text{and} \quad \sum_{\alpha=1}^n b_{\alpha i}^2 = 1$$

Now consider  $n\hat{\Sigma}$  from (1) i.e.

$$\begin{aligned}
n\hat{\Sigma} &= \sum_{\alpha=1}^n \left( \underline{\underline{X}}_{\alpha} - \bar{\underline{\underline{X}}} \right) \left( \underline{\underline{X}}_{\alpha} - \bar{\underline{\underline{X}}} \right)' \\
&= \sum_{\alpha=1}^n \underline{\underline{X}}_{\alpha} \underline{\underline{X}}_{\alpha}' - \sum_{\alpha=1}^n \bar{\underline{\underline{X}}} \bar{\underline{\underline{X}}}' \\
&\left( \because \bar{\underline{\underline{X}}} \sum_{\alpha=1}^n \left( \underline{\underline{X}}_{\alpha} - \bar{\underline{\underline{X}}} \right)' = 0 \right)
\end{aligned}$$

$$\begin{aligned}
&= \sum_{\alpha=1}^n \tilde{\mathbf{X}}_{\alpha} \tilde{\mathbf{X}}_{\alpha}' - n \bar{\tilde{\mathbf{X}}} \bar{\tilde{\mathbf{X}}}' \\
&= \sum_{\alpha=1}^n \tilde{\mathbf{z}}_{\alpha} \tilde{\mathbf{z}}_{\alpha}' - \tilde{\mathbf{z}}_n \tilde{\mathbf{z}}_n' \quad \text{[Using (3) and (4)]} \\
&= \sum_{\alpha=1}^{n-1} \tilde{\mathbf{z}}_{\alpha} \tilde{\mathbf{z}}_{\alpha}' \quad \text{--- (5)}
\end{aligned}$$

From (3) and (5) we observe that  $\bar{\tilde{\mathbf{X}}} (\hat{\boldsymbol{\mu}})$  is distributed according to the distribution of  $\tilde{\mathbf{z}}_n$  and  $n\hat{\boldsymbol{\Sigma}}$  (and hence  $\hat{\boldsymbol{\Sigma}}$ ) is distributed according to the distribution of  $\tilde{\mathbf{z}}_1, \dots, \tilde{\mathbf{z}}_{n-1}$ .

Also, since  $\tilde{\mathbf{z}}_1, \dots, \tilde{\mathbf{z}}_n$  are obtained from  $\tilde{\mathbf{X}}_1, \dots, \tilde{\mathbf{X}}_n$  using the orthogonal linear transformation (using orthogonal matrix  $\mathbf{B}$ )  $\tilde{\mathbf{z}}_1, \dots, \tilde{\mathbf{z}}_n$  are independently distributed as Multivariate normal distribution with common covariance matrix '  $\boldsymbol{\Sigma}$  '. Therefore  $\hat{\boldsymbol{\mu}}$  and  $\hat{\boldsymbol{\Sigma}}$  are independently distributed.

Now let us obtain the mean vector of  $\tilde{\mathbf{z}}_1, \dots, \tilde{\mathbf{z}}_n$

From (3)

$$\begin{aligned}
E(\tilde{\mathbf{z}}_n) &= \sqrt{n} E(\bar{\tilde{\mathbf{X}}}) \\
&= \sqrt{n} \frac{1}{n} \left( E(\tilde{\mathbf{X}}_1 + \tilde{\mathbf{X}}_2 + \dots + \tilde{\mathbf{X}}_n) \right) \\
&\quad [ \because \tilde{\mathbf{X}}_i \text{'s are independent} ]
\end{aligned}$$

$$= \sqrt{n} \frac{1}{n} \sum_{i=1}^n E(X_i) = \sqrt{n} \frac{1}{n} \sum_{i=1}^n \mu = \sqrt{n} \mu$$

$$\left( \because \underline{\mathbf{X}}_i \sim Np(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}) \right)$$

$$\text{Thus } \underline{\mathbf{z}}_n \sim Np(\sqrt{n}\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$$

$$\text{i.e. } \sqrt{n}\bar{\underline{\mathbf{X}}} \sim Np(\sqrt{n}\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})$$

$$\text{i.e. } \bar{\underline{\mathbf{X}}} \sim Np\left(\underline{\boldsymbol{\mu}}, \frac{\boldsymbol{\Sigma}}{n}\right)$$

From (2), we have

$$E(\underline{\mathbf{z}}_\alpha) = \sum_{\beta=1}^n b_{\alpha\beta} E(\underline{\mathbf{X}}_\beta) \quad [ \because \underline{\mathbf{X}}_\beta \text{ 's are independent } ]$$

$$= \sum_{\beta=1}^n b_{\alpha\beta} \underline{\boldsymbol{\mu}}$$

$$= \sum_{\beta=1}^n b_{\alpha\beta} \frac{1}{\sqrt{n}} \sqrt{n} \underline{\boldsymbol{\mu}}$$

$$= \sqrt{n} \underline{\boldsymbol{\mu}} \sum_{\beta=1}^n b_{\alpha\beta} b_{n\beta} \quad [ \because b_{n\beta} = \frac{1}{\sqrt{n}} ]$$

$$= \underline{\mathbf{0}} \quad \forall \alpha \neq n$$

Thus each of  $\underline{z}_1, \dots, \underline{z}_{n-1}$  are distributed as  $Np(\underline{0}, \Sigma)$ . Therefore

from (5)  $n\hat{\Sigma}$  is distributed as  $\sum_{\alpha=1}^{n-1} \underline{z}_\alpha \underline{z}'_\alpha$ , where  $\underline{z}_\alpha \sim N(\underline{0}, \Sigma)$

and is independent of  $\underline{z}_\beta$  ( $\beta \neq \alpha$ )

Thus the MLE's of  $\underline{\mu}$  and  $\Sigma$  are independently distributed.

Hence the proof.

**NOTE:-**

➤ Since  $\underline{X}_1, \dots, \underline{X}_n$  is a random sample

$$E(\bar{\underline{X}}) = \frac{1}{n} (E(\underline{X}_1) + E(\underline{X}_2) + \dots + E(\underline{X}_n))$$

$$= \frac{n \underline{\mu}}{n} = \underline{\mu}$$

Thus  $\bar{\underline{X}}$  is an unbiased estimator of  $\underline{\mu}$ . Thus sample mean is an unbiased estimator of the population mean vector  $\underline{\mu}$ .

$$\begin{aligned} \text{➤ } E(\hat{\Sigma}) &= \frac{1}{n} E\left(\sum_{\alpha=1}^{n-1} \underline{z}_\alpha \underline{z}'_\alpha\right) \\ &= \frac{1}{n} \sum_{\alpha=1}^{n-1} E(\underline{z}_\alpha \underline{z}'_\alpha) \quad [ \because \underline{z}_\alpha \text{'s are} \end{aligned}$$

independent]



$$\begin{aligned}
&= \frac{1}{n} \sum_{\alpha=1}^{n-1} V(\mathbf{z}_{\alpha}) \quad [\because E(\mathbf{z}_{\alpha} = \mathbf{0})] \\
&= \frac{1}{n} \sum_{\alpha=1}^{n-1} \Sigma = \frac{n-1}{n} \Sigma
\end{aligned}$$

Thus  $\hat{\Sigma}$  is not an unbiased estimator of  $\Sigma$ . But

$$\frac{n}{n-1} \hat{\Sigma} = \frac{1}{n-1} \sum_{\alpha=1}^n (\mathbf{X}_{\alpha} - \bar{\mathbf{X}}) (\mathbf{X}_{\alpha} - \bar{\mathbf{X}})' = \mathbf{S}$$

(say) is an unbiased estimator of  $\Sigma$   $[\because E\left(\frac{n}{n-1} \hat{\Sigma}\right) = \Sigma]$

Hence  $\mathbf{S} = \frac{1}{n-1} \sum_{\alpha=1}^n (\mathbf{X}_{\alpha} - \bar{\mathbf{X}}) (\mathbf{X}_{\alpha} - \bar{\mathbf{X}})'$  is

called the sample covariance matrix and is an unbiased estimator of  $\Sigma$ .

- From the above theorem, it obviously follows that the sample mean ( $\bar{\mathbf{X}}$ ) and the sample covariance matrix

$$\left( \mathbf{S} = \frac{1}{n-1} \sum_{\alpha=1}^n (\mathbf{X}_{\alpha} - \bar{\mathbf{X}}) (\mathbf{X}_{\alpha} - \bar{\mathbf{X}})' \right)$$

are independently distributed. Also it may be seen that

$$\bar{\mathbf{X}} \sim N\left(\boldsymbol{\mu}, \frac{\Sigma}{n}\right)$$

$$S = \frac{1}{n-1} \sum_{\alpha=1}^{n-1} \mathbf{z}_{\alpha} \mathbf{z}_{\alpha}' , \text{ where } \mathbf{z}_{\alpha} \sim N_p(\mathbf{0}, \Sigma)$$

i.e.  $(n-1)S$  is distributed according to the distribution of  $\sum_{\alpha=1}^{n-1} \mathbf{z}_{\alpha} \mathbf{z}_{\alpha}'$ , where

$\mathbf{z}_1, \dots, \mathbf{z}_n$  are independently distributed as  $N_p(\mathbf{0}, \Sigma)$ . The matrix

$\sum_{\alpha=1}^{n-1} \mathbf{z}_{\alpha} \mathbf{z}_{\alpha}'$  is called “Wishart random matrix” and it is distributed according to “wishart distribution” with  $(n-1)$  degrees of freedom.

And is denoted as  $W_{n-1}(\Sigma)$ , where  $\Sigma$  is the covariance matrix of Wishart random matrix .

Hence it may be noted that

$$\begin{aligned} \Sigma &= E \left( \frac{1}{n-1} \sum_{\alpha=1}^{n-1} \mathbf{z}_{\alpha} \mathbf{z}_{\alpha}' \right) \\ &= E \left( \frac{\text{Wishart random matrix}}{\text{degrees of freedom}} \right) \end{aligned}$$

Thus  $(n-1)S$  (and hence  $S$ ) provides independent information about  $\Sigma$  and the distribution of  $S$  does not depend on  $\mu$ . This allows us to construct a statistics for making inferences about  $\mu$  as we shall see in the later section.

**Properties of the Wishart distribution:-**

➤ If  $\mathbf{A}_1$  is distributed as  $W_{m_1}(\Sigma)$  independently of  $\mathbf{A}_2$ , which is distributed as  $W_{m_2}(\Sigma)$ , then  $\mathbf{A}_1 + \mathbf{A}_2$  is distributed as  $W_{m_1+m_2}(\Sigma)$ . That is the degrees of freedom are added.

**Proof:** - Since  $\mathbf{A}_1 \sim W_{m_1}(\Sigma)$ .

$\mathbf{A}_1$  may be written as

$$\mathbf{A}_1 = \sum_{\alpha=1}^{m_1} \mathbf{z}_\alpha \mathbf{z}'_\alpha, \text{ where } \mathbf{z}_\alpha \sim Np(\mathbf{0}, \Sigma)$$

Also since  $\mathbf{A}_2$  is independently distributed as  $W_{m_2}(\Sigma)$ ,

We may write

$$\mathbf{A}_2 = \sum_{\alpha=m_1+1}^{m_1+m_2} \mathbf{z}_\alpha \mathbf{z}'_\alpha, \text{ where } \mathbf{z}_\alpha \sim Np(\mathbf{0}, \Sigma)$$

Since  $\mathbf{A}_1$  and  $\mathbf{A}_2$  are independent,  $\mathbf{z}_1, \dots, \mathbf{z}_{m_1+m_2}$  are independent and as a consequence .

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2 = \sum_{\alpha=1}^{m_1+m_2} \mathbf{z}_\alpha \mathbf{z}'_\alpha \sim W_{m_1+m_2}(\Sigma)$$

Hence the proof.

➤ If  $\mathbf{A} \sim W_m(\Sigma)$ , then  $\mathbf{CAC}' \sim W_m(\mathbf{C}\Sigma\mathbf{C}')$

**Proof:**-Given  $\mathbf{A} \sim W_m(\Sigma)$

$$\therefore \mathbf{A} = \sum_{\alpha=1}^m \mathbf{z}_{\alpha} \mathbf{z}_{\alpha}' , \text{ where } \mathbf{z}_{\alpha} \sim N_p(\mathbf{0}, \mathbf{\Sigma})$$

$$\mathbf{CAC}' = \mathbf{W}_m \sum_{\alpha=1}^m \mathbf{Cz}_{\alpha} \mathbf{z}_{\alpha}' \mathbf{C}' = \sum_{\alpha=1}^m \mathbf{Y}_{\alpha} \mathbf{Y}_{\alpha}'$$

$$\text{where } \mathbf{Y}_{\alpha} = \mathbf{Cz}_{\alpha} \sim N_p(\mathbf{0}, \mathbf{C}\mathbf{\Sigma}\mathbf{C}')$$

$$\therefore E(\mathbf{Y}_{\alpha}) = \mathbf{C}E(\mathbf{z}_{\alpha}) = \mathbf{0}$$

$$V(\mathbf{Y}_{\alpha}) = V(\mathbf{Cz}_{\alpha}) = \mathbf{C}\mathbf{\Sigma}\mathbf{C}' \text{ and } \mathbf{Y}_{\alpha} \text{ is normal random}$$

vector]

$$\mathbf{CAC}' \sim \mathbf{W}_m(\mathbf{C}\mathbf{\Sigma}\mathbf{C}')$$

Hence the proof.

#### P.D.F. Of Wishart Distribution :-

The p.d.f. of  $\mathbf{A} \sim \mathbf{W}_n(\mathbf{\Sigma})$  is given by

$$W_n \left( \frac{\mathbf{A}}{\mathbf{\Sigma}} \right) = \frac{|\mathbf{A}|^{(n-p-1)/2} e^{-\text{tr}(\mathbf{A}\mathbf{\Sigma}^{-1})/2}}{2^{np/2} \pi^{p(p-1)/4} |\mathbf{\Sigma}|^{n/2} \prod_{i=1}^p \left[ \left( \frac{1}{2} (n+1-i) \right) \right]}$$

$\mathbf{A}$  is positive definite and  $\Gamma(\cdot)$  is gamma function.

## INFERENCES ABOUT MULTIVARIATE NORMAL MEAN VECTOR(S)

One Sample Problem:

Suppose  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  is a random sample from a multivariate normal population. Now, our statistical problem is whether the given sample has come from the multivariate normal population, whose mean vector is given by  $\underline{\mu} = \underline{\mu}_0$ . In other words, we have to test

$$H_0 : \underline{\mu} = \underline{\mu}_0 \text{ vs } H_1 : \underline{\mu} \neq \underline{\mu}_0$$

based on the given random sample  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ .

Two Sample Problem:

Suppose we have two different samples from two different multivariate normal populations  $N_p(\underline{\mu}^{(1)}, \underline{\Sigma})$  and  $N_p(\underline{\mu}^{(2)}, \underline{\Sigma})$  with common variance-covariance matrix  $\underline{\Sigma}$ . Now, our statistical problem is whether the two normal populations have the same mean vector or not. In other words, our problem is equivalent to test the hypothesis

$$H_0 : \underline{\mu}^{(1)} = \underline{\mu}^{(2)} \text{ vs } H_1 : \underline{\mu}^{(1)} \neq \underline{\mu}^{(2)}$$

based on the given two samples.

For developing the test statistics in the above two problems, we have to consider whether the common covariance matrix  $\underline{\Sigma}$  is known or not. First, let us develop the test statistics for the above one-sample case as well as two-sample case assuming the population variance-covariance matrix  $\underline{\Sigma}$  is known.

**Developing Test Statistics When  $\underline{\Sigma}$  is known :**

**One-Sample problem:**

If  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  is a random sample of size  $n$  drawn from a multivariate normal population with known variance-covariance matrix  $\underline{\Sigma}$ , then obtain the test statistic for testing

$$H_0 : \underline{\mu} = \underline{\mu}_0 \text{ vs } H_1 : \underline{\mu} \neq \underline{\mu}_0$$

and the critical region of size ' $\alpha$ ' as well as the confidence region for  $\underline{\mu}$  of confidence  $1 - \alpha$ .

**Proof :-**

In order to obtain the above result, let us prove the following theorem.

**Theorem(1) :**

If a  $p$ -component vector  $\underline{Y} \sim N_p(\underline{0}, \underline{\Sigma})$ , where  $\underline{\Sigma}$  is non-singular (positive definite), then

$$\underline{Y}' \underline{\Sigma}^{-1} \underline{Y} \sim \chi_p^2 \quad \rightarrow (1)$$

where  $\chi_p^2$  is Chi-square distribution with  $p$  d.f.

**Proof:** We have given  $\underline{Y} \sim N_p(\underline{0}, \underline{\Sigma})$ .

Since,  $\underline{\Sigma}$  is p.d.f  $\exists$  a non-singular matrix  $\underline{C}$  such that,

$$\underline{C} \underline{\Sigma} \underline{C}' = \underline{I}$$

$$\Rightarrow \Sigma = \mathbf{C}^{-1} \mathbf{I} (\mathbf{C}')^{-1} = (\mathbf{C}' \mathbf{C})^{-1} \quad \rightarrow (2)$$

Let us define the linear transformation,

$$\underline{\mathbf{Z}} = \mathbf{C} \underline{\mathbf{Y}} \quad \rightarrow (3)$$

Then,  $E(\underline{\mathbf{Z}}) = \mathbf{C} E(\underline{\mathbf{Y}}) = \underline{\mathbf{0}}$

$$V(\underline{\mathbf{Z}}) = V(\mathbf{C} \underline{\mathbf{Y}}) = \mathbf{C} V(\underline{\mathbf{Y}}) \mathbf{C}' = \mathbf{C} \Sigma \mathbf{C}' = \mathbf{I} \quad (\text{from (2)})$$

Since the transformation is linear,

$\underline{\mathbf{Z}} \sim N_p(\underline{\mathbf{0}}, \mathbf{I})$  i.e.,  $Z_1, Z_2, \dots, Z_p$ , the individual components of  $\underline{\mathbf{Z}}$  are distributed as  $N(0,1)$ .

Further, since the covariances are zeros,  $Z_1, Z_2, \dots, Z_p$  are independent which follows from the normality of the components.

$$\begin{aligned} \therefore \underline{\mathbf{Z}}' \underline{\mathbf{Z}} &= Z_1^2 + Z_2^2 + \dots + Z_p^2 \sim \chi_p^2 \\ \Rightarrow \underline{\mathbf{Y}}' \mathbf{C}' \mathbf{C} \underline{\mathbf{Y}} &\sim \chi_p^2 \quad (\text{from (3)}) \\ \Rightarrow \underline{\mathbf{Y}}' \Sigma^{-1} \underline{\mathbf{Y}} &\sim \chi_p^2 \quad (\text{from (2)}) \end{aligned}$$

Hence the result (1).

### Proof of the Result:-

We have given the random sample  $\underline{\mathbf{X}}_1, \underline{\mathbf{X}}_2, \dots, \underline{\mathbf{X}}_n$  from  $N_p(\underline{\boldsymbol{\mu}}, \Sigma)$ , where  $\Sigma$  is known.

Now, we know that, the sample mean,  $\bar{\underline{\mathbf{X}}} \sim N_p(\underline{\boldsymbol{\mu}}, \Sigma/n)$ .

Define the random vector,  $\underline{\mathbf{Y}} = \sqrt{n} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}})$   $\rightarrow (4)$

With  $E(\underline{\mathbf{Y}}) = \sqrt{n} E(\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}) = \sqrt{n} (\underline{\boldsymbol{\mu}} - \underline{\boldsymbol{\mu}}) = \underline{\mathbf{0}}$ .

$$V(\underline{\mathbf{Y}}) = V(\sqrt{n} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}})) = n V(\bar{\underline{\mathbf{X}}}) = n \Sigma/n = \Sigma.$$

Thus, the mean vector of  $\underline{\mathbf{Y}}$  is  $\underline{\mathbf{0}}$  and covariance matrix is  $\Sigma$ .

Further, since the transformation in (4) is linear,

$$\underline{\mathbf{Y}} \sim N_p(\underline{\mathbf{0}}, \Sigma)$$

Now, from the above **theorem (1)**, it immediately follows

$$\begin{aligned} \underline{\mathbf{Y}}' \Sigma^{-1} \underline{\mathbf{Y}} &\sim \chi_p^2 \\ \Rightarrow \sqrt{n} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}})' \Sigma^{-1} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}) \sqrt{n} &\sim \chi_p^2 \quad (\text{from (4)}) \\ \Rightarrow n (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}})' \Sigma^{-1} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}) &\sim \chi_p^2 \end{aligned}$$

Thus, the test statistic for  $H_0 : \underline{\boldsymbol{\mu}} = \underline{\boldsymbol{\mu}}_0$  is given by

$$n (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}_0)' \Sigma^{-1} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}_0) \quad \rightarrow (5)$$

which follows  $\chi^2$  distribution with  $p$  d.f.

Let  $\chi_p^2(\alpha)$  be the number such that  $\Pr \{ \chi_p^2 \geq \chi_p^2(\alpha) \} = \alpha$ .

Thus,  $\Pr \{ n (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}_0)' \Sigma^{-1} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}_0) \geq \chi_p^2(\alpha) \} = \alpha$

and to test  $H_0 : \underline{\boldsymbol{\mu}} = \underline{\boldsymbol{\mu}}_0$  (given), we use

$$n (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}_0)' \Sigma^{-1} (\bar{\underline{\mathbf{X}}} - \underline{\boldsymbol{\mu}}_0) \geq \chi_p^2(\alpha) \quad \rightarrow (6)$$

as critical region.

Similarly, we use the inequality,

$$(\bar{\mathbf{X}} - \underline{\boldsymbol{\mu}}^*)' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}} - \underline{\boldsymbol{\mu}}^*) \leq \chi_p^2(\alpha) \quad \rightarrow (7)$$

for obtaining the confidence region for  $\underline{\boldsymbol{\mu}}$  (the set of all  $\underline{\boldsymbol{\mu}}^*$  satisfying (7)) with confidence  $1 - \alpha$ .

Hence the result.

### Two Sample Problem :-

Suppose we have a sample  $\underline{\mathbf{X}}_1^{(1)}, \underline{\mathbf{X}}_2^{(1)}, \dots, \underline{\mathbf{X}}_{n_1}^{(1)}$  from  $N_p(\underline{\boldsymbol{\mu}}^{(1)}, \boldsymbol{\Sigma})$  and another sample  $\underline{\mathbf{X}}_1^{(2)}, \underline{\mathbf{X}}_2^{(2)}, \dots, \underline{\mathbf{X}}_{n_2}^{(2)}$  from  $N_p(\underline{\boldsymbol{\mu}}^{(2)}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\Sigma}$  is known. Now, under the null hypothesis

$$H_0: \underline{\boldsymbol{\mu}}^{(1)} = \underline{\boldsymbol{\mu}}^{(2)}$$

$$\frac{n_1 n_2}{n_1 + n_2} \left[ (\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)})' \boldsymbol{\Sigma}^{-1} (\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) \right] \sim \chi_p^2,$$

where,  $\bar{\mathbf{X}}^{(1)}$  = mean of the random sample  $\underline{\mathbf{X}}_1^{(1)}, \underline{\mathbf{X}}_2^{(1)}, \dots, \underline{\mathbf{X}}_{n_1}^{(1)}$

and  $\bar{\mathbf{X}}^{(2)}$  = mean of the random sample  $\underline{\mathbf{X}}_1^{(2)}, \underline{\mathbf{X}}_2^{(2)}, \dots, \underline{\mathbf{X}}_{n_2}^{(2)}$ .

### Solution:

From the given hypothesis, we have

$$\bar{\mathbf{X}}^{(1)} \sim N_p(\underline{\boldsymbol{\mu}}^{(1)}, \boldsymbol{\Sigma}/n_1) \quad \& \quad \bar{\mathbf{X}}^{(2)} \sim N_p(\underline{\boldsymbol{\mu}}^{(2)}, \boldsymbol{\Sigma}/n_2) \quad \rightarrow (1)$$

$$\text{Now define, } \underline{\mathbf{Y}} = \bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)} \quad \rightarrow (2)$$

$$\text{with mean vector, } E(\underline{\mathbf{Y}}) = \underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)} \quad \rightarrow (3)$$

and variance- covariance matrix,

$$\begin{aligned} V(\underline{\mathbf{Y}}) &= V(\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) \\ &= V(\bar{\mathbf{X}}^{(1)}) + V(\bar{\mathbf{X}}^{(2)}) - \text{Cov}(\bar{\mathbf{X}}^{(1)}, \bar{\mathbf{X}}^{(2)}) - \text{Cov}(\bar{\mathbf{X}}^{(2)}, \bar{\mathbf{X}}^{(1)}) \\ &= \frac{\boldsymbol{\Sigma}}{n_1} + \frac{\boldsymbol{\Sigma}}{n_2} = \boldsymbol{\Sigma} \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \end{aligned} \quad \rightarrow (4)$$

(since the two samples are independent and as a consequence the covariance matrices  $\text{Cov}(\bar{\mathbf{X}}^{(1)}, \bar{\mathbf{X}}^{(2)}) = \mathbf{0}$  &  $\text{Cov}(\bar{\mathbf{X}}^{(2)}, \bar{\mathbf{X}}^{(1)}) = \mathbf{0}$ )

Since the transformation used in (2) is linear, we have

$$\begin{aligned} \underline{\mathbf{Y}} &\sim N_p \left( \underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)}, \boldsymbol{\Sigma} \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right) \\ \Rightarrow \underline{\mathbf{Y}} - (\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)}) &\sim N_p \left( \mathbf{0}, \boldsymbol{\Sigma} \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right) \end{aligned}$$

Now, from the above theorem (1), it immediately follows

$$\left[ \underline{\mathbf{Y}} - (\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)}) \right]' \left[ \boldsymbol{\Sigma} \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right]^{-1} \left[ \underline{\mathbf{Y}} - (\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)}) \right] \sim \chi_p^2$$

But, under the null hypothesis,  $H_0: \underline{\boldsymbol{\mu}}^{(1)} = \underline{\boldsymbol{\mu}}^{(2)}$

$$\begin{aligned} & (\bar{X}^{(1)} - \bar{X}^{(2)})' \left[ \Sigma \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \right]^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)}) \sim \chi_p^2 \quad (\text{from(2)}) \\ \Rightarrow & \frac{n_1 n_2}{n_1 + n_2} (\bar{X}^{(1)} - \bar{X}^{(2)})' \Sigma^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)}) \sim \chi_p^2 \end{aligned}$$

Hence the proof .

**NOTE :-**

(i) The critical region for testing hypothesis,  $H_0 : \underline{\mu}^{(1)} = \underline{\mu}^{(2)}$  at  $\alpha$  level of significance is given by  $\left( \frac{n_1 n_2}{n_1 + n_2} \right) (\bar{X}^{(1)} - \bar{X}^{(2)})' \Sigma^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)}) \geq \chi_p^2(\alpha)$

where,  $\chi_p^2(\alpha)$  is a number such that  $P(\chi_p^2 \leq \chi_p^2(\alpha)) = \alpha$  .

(ii) Mahalonobis (1930) suggested ,  $(\underline{\mu}^{(1)} - \underline{\mu}^{(2)})' \Sigma^{-1} (\underline{\mu}^{(1)} - \underline{\mu}^{(2)})$  as a measure of the distance between two populations .

(iii) The  $(1 - \alpha)\%$  confidence interval for  $\underline{\mu}^{(1)} - \underline{\mu}^{(2)}$  is given by the set of  $\underline{y}^*$  all which are satisfying the inequality,

$$\left( \frac{n_1 n_2}{n_1 + n_2} \right) (\bar{X}^{(1)} - \bar{X}^{(2)} - \underline{y}^*)' \Sigma^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)} - \underline{y}^*) \leq \chi_p^2(\alpha).$$

**When ‘ $\Sigma$ ’ Is Unknown:-**

**Introduction:-**

One of the most important groups of problems in Univariate theory relates to questions concerning the mean ( $\mu$ ) of a give distribution, when the variance ( $\sigma^2$ ) of the distribution is unknown . On the basis of a sample, one may wish to decide whether the mean is equal to a specified number ( $\mu_0$ ) or one may wish to give an interval with in which ‘ $\mu$ ’ lies . The following is the mathematical formulation of one sample t-test.

Suppose  $X_1, X_2, \dots, X_n$  denote a random sample from a normal population  $N(\mu, \sigma)$ . Now the well – known test statistic for testing the hypothesis  $H_0 : \mu = \mu_0$  against  $H_1 : \mu \neq \mu_0$  is given by

$$t = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} \quad \text{where, } \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

This test statistic follows students’ t-distribution with n-1 degrees of freedom . We reject  $H_0$  is  $|t|$  exceeds a specified percentage point of a t-distribution with n-1 degrees of freedom .

$$\begin{aligned} \frac{(n-1) S^2}{\sigma^2} & \sim \chi_{n-1}^2 \\ (\bar{X} - \mu) & \sim N\left(0, \frac{\sigma^2}{n}\right) \\ \text{i.e., } \sqrt{n} (\bar{X} - \mu) & \sim N(0, \sigma) \end{aligned}$$



Rejecting  $H_0$  when  $|t|$  is large is equivalently to rejecting  $H_0$  is its square.

$$t^2 = \frac{(\bar{X} - \mu_0)^2}{s^2/n} = n(\bar{X} - \mu_0)(S^2)^{-1}(\bar{X} - \mu_0) \rightarrow (1)$$

is large. Thus the test becomes.

Reject  $H_0$  in favour of  $H_1$  at  $\alpha$  level of significance, if

$$n(\bar{X} - \mu_0)(S^2)^{-1}(\bar{X} - \mu_0) > t_{n-1}^2(\alpha/2) \rightarrow (2)$$

where,  $t_{n-1}(\alpha/2)$  denotes the upper  $100(\alpha/2)^{th}$  percentile of the t-distribution with  $n-1$ d.f and  $\bar{X}$  is the mean of the given sample  $X_1, X_2, \dots, X_n$  and  $S^2$  is the sample variance of the given sample.

Now, the multivariate analogue of  $t^2$  given in (1) is

$$T^2 = n(\bar{\underline{X}} - \underline{\mu}_0)' S^{-1}(\bar{\underline{X}} - \underline{\mu}_0) \rightarrow (3)$$

Where,  $\bar{\underline{X}}$  &  $S$  are respectively the sample mean and sample covariance matrix obtained. Using the multivariate normal sample  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  and are given by

$$\begin{aligned} \bar{\underline{X}} &= \frac{1}{n} \sum_{i=1}^n \underline{X}_i \\ S &= \frac{1}{n-1} \sum_{i=1}^n (\underline{X}_i - \bar{\underline{X}})(\underline{X}_i - \bar{\underline{X}})' \end{aligned} \rightarrow (4)$$

and  $\underline{\mu}_0 = \begin{pmatrix} \mu_{10} \\ \mu_{20} \\ \vdots \\ \mu_{p0} \end{pmatrix}_{p \times 1}$  is the specified value of population mean vector  $\underline{\mu}$ .

The statistic  $T^2$  is called Hotelling's  $T^2$  in honour of Harnold Hotelling, a pioneer in multivariate analysis, who first obtained its sampling distribution.

$$\text{It may be proved that } T^2 \sim \frac{(n-1)p}{(n-p)} F_{p, n-p} \rightarrow (5)$$

where,  $F_{p, n-p}$  denotes F-distribution with  $p, n-p$  degrees of freedom.

The above explanation about  $T^2$  statistic can be summarized as follows.

Let  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  be a random sample from  $N_p(\underline{\mu}, \Sigma)$  population ( $\Sigma$  is unknown), then Hotellings  $T^2$ - statistics for testing.

$H_0 : \underline{\mu} = \underline{\mu}_0$  against  $H_1 : \underline{\mu} \neq \underline{\mu}_0$ , is given by (3) and  $H_0$  may be selected at  $\alpha$ -level of significance infavour of  $H_1$  is

$$\begin{aligned} T^2 &= n(\bar{\underline{X}} - \underline{\mu}_0)' S^{-1}(\bar{\underline{X}} - \underline{\mu}_0) \\ &> \frac{(n-1)p}{(n-p)} F_{p, n-p}(\alpha) \end{aligned} \rightarrow (6)$$

where,  $F_{p,n-p}(\alpha)$  is the upper  $100\alpha^{\text{th}}$  percentile of the  $F_{p,n-p}$  distribution and  $\bar{\mathbf{X}}$  &  $\mathbf{S}$  are respectively. The mean vector and covariance matrix of the given sample  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ .

**Definition of Hotelling's  $T^2$ -distribution(statistic):**

Suppose  $\mathbf{Y}$  is a  $p$ -variate random vector distributed according to  $N_p(\mathbf{0}, \Sigma)$ .

Suppose  $\mathbf{B} = \sum_{i=1}^n \mathbf{Z}_i \mathbf{Z}_i'$  (where each  $\mathbf{Z}_i \sim N_p(\mathbf{0}, \Sigma)$  and are independent) is a Wishart random matrix and is distributed as Wishart distribution with  $n$  degrees of freedom i.e.  $\mathbf{B} \sim \mathbf{W}_n(\Sigma)$ .

Now, if  $\mathbf{Y}$  and  $\mathbf{B}$  are independent then the quantity

$$T^2 = \mathbf{Y}' \left( \frac{\mathbf{B}}{n} \right)^{-1} \mathbf{Y}$$

is called as Hotelling's  $T^2$  statistic and its distribution is called as Hotelling's  $T^2$ -distribution with  $n$  d.f. and is denoted as

$$T^2 \sim T_n^2$$

**NOTE :-**

(1) the statistic  $\frac{T^2}{n}$  is the ratio of two independent  $\chi^2$ - variates with  $p$  d.f.

and  $n-p+1$  d.f. respectively i.e.,  $\frac{T^2}{n} = \frac{\chi_p^2}{\chi_{n-p+1}^2}$ .

(2)  $\frac{T^2}{n} \left( \frac{n-p+1}{p} \right) \sim F_{p,n-p+1}$ , where  $F_{p,n-p+1}$  is the  $F$ -distribution with  $p$  and  $n-p+1$  degrees of freedom.

**THE LIKELIHOOD RATION PRINCIPLE**

**Deriving  $T^2$ -statistic as the Likelihood Ratio Test of  $H_0: \boldsymbol{\mu} = \boldsymbol{\mu}_0$ :**

There is a general principle for constructing test procedures called the Likelihood Ratio (LR) principle method and the  $T^2$ -statistic can be derived as the LR test of  $H_0: \boldsymbol{\mu} = \boldsymbol{\mu}_0$  as explained below.

Suppose  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  ( $n > p$ ) is given random sample from  $N_p(\boldsymbol{\mu}, \Sigma)$ , the likelihood function is

$$L(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2}} e^{-\frac{1}{2} \sum_{\alpha=1}^n (\mathbf{x}_{\alpha} - \underline{\boldsymbol{\mu}})' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{\alpha} - \underline{\boldsymbol{\mu}})} \rightarrow (1)$$

Under the hypothesis,  $H_0: \underline{\boldsymbol{\mu}} = \underline{\boldsymbol{\mu}}_0$ , the likelihood becomes,

$$L(\underline{\boldsymbol{\mu}}_0, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{n/2}} e^{-\frac{1}{2} \sum_{\alpha=1}^n (\mathbf{x}_{\alpha} - \underline{\boldsymbol{\mu}}_0)' \boldsymbol{\Sigma}^{-1} (\mathbf{x}_{\alpha} - \underline{\boldsymbol{\mu}}_0)} \rightarrow (2)$$

The likelihood ratio criterion is

$$\lambda = \frac{\max_{\boldsymbol{\Sigma}} L(\underline{\boldsymbol{\mu}}_0, \boldsymbol{\Sigma})}{\max_{\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}} L(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma})} \rightarrow (3)$$

i.e., the numerator is the maximum of the likelihood function for  $\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}$  is the parameter space restricted by the null hypothesis ( $\underline{\boldsymbol{\mu}} = \underline{\boldsymbol{\mu}}_0$ ) and  $\boldsymbol{\Sigma}$  is positive definite and the denominator is the maximum over the entire parameter space ( $\boldsymbol{\Sigma}$  is positive definite).

When the parameters are unrestricted the MLE's of  $\underline{\boldsymbol{\mu}}$  and  $\boldsymbol{\Sigma}$  from (1) are given by

$$\begin{aligned} \hat{\underline{\boldsymbol{\mu}}}_{\Omega} &= \bar{\mathbf{x}} \\ \hat{\boldsymbol{\Sigma}}_{\Omega} &= \frac{1}{n} \sum_{\alpha} (\mathbf{x}_{\alpha} - \bar{\mathbf{x}}) (\mathbf{x}_{\alpha} - \bar{\mathbf{x}})' \end{aligned} \rightarrow (4)$$

When  $\underline{\boldsymbol{\mu}} = \underline{\boldsymbol{\mu}}_0$ , the likelihood function given by (2), minimizes at

$$\hat{\boldsymbol{\Sigma}}_0 = \frac{1}{n} \sum_{\alpha} (\mathbf{x}_{\alpha} - \underline{\boldsymbol{\mu}}_0) (\mathbf{x}_{\alpha} - \underline{\boldsymbol{\mu}}_0)' \rightarrow (5)$$

Substituting (4) in (1), we get ( after simplification),

$$\max_{\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}} L(\underline{\boldsymbol{\mu}}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{np/2} |\hat{\boldsymbol{\Sigma}}_{\Omega}|^{n/2}} e^{-np/2} \rightarrow (6)$$

Similarly, substituting (5) in (2), we get

$$\max_{\boldsymbol{\Sigma}} L(\underline{\boldsymbol{\mu}}_0, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{np/2} |\hat{\boldsymbol{\Sigma}}_0|^{n/2}} e^{-np/2} \rightarrow (7)$$

Substituting (6) & (7) in (3), we get,

$$\lambda = \left( \frac{|\hat{\boldsymbol{\Sigma}}_{\Omega}|}{|\hat{\boldsymbol{\Sigma}}_0|} \right)^{n/2}$$

$$\Rightarrow \lambda^{2/n} = \frac{|\hat{\Sigma}_{\Omega}|}{|\hat{\Sigma}_0|} = \frac{\left| \sum_{\alpha=1}^n (\mathbf{x}_\alpha - \bar{\mathbf{x}})(\mathbf{x}_\alpha - \bar{\mathbf{x}})' \right|}{\left| \sum_{\alpha=1}^n (\mathbf{x}_\alpha - \underline{\boldsymbol{\mu}}_0)(\mathbf{x}_\alpha - \underline{\boldsymbol{\mu}}_0)' \right|} \rightarrow (7a)$$

$$\Rightarrow \lambda^{2/n} = \frac{|\mathbf{A}|}{|\mathbf{A} + n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)'|}$$

$$\text{Where } \mathbf{A} = \sum_{\alpha=1}^n (\mathbf{x}_\alpha - \bar{\mathbf{x}})(\mathbf{x}_\alpha - \bar{\mathbf{x}})' \rightarrow (8)$$

$$\text{Consider the matrix, } \mathbf{B}_{(p+1) \times (p+1)} = \begin{bmatrix} \mathbf{A} & \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0) \\ \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' & -1 \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{B}_{11} & \vdots & \mathbf{B}_{12} \\ \dots & \vdots & \dots \\ \mathbf{B}_{21} & \vdots & \mathbf{B}_{22} \end{bmatrix}.$$

$$\text{We have, } |\mathbf{B}| = |\mathbf{B}_{11}| |\mathbf{B}_{22} - \mathbf{B}_{21} \mathbf{B}_{11}^{-1} \mathbf{B}_{12}|$$

$$= |\mathbf{B}_{22}| |\mathbf{B}_{11} - \mathbf{B}_{12} \mathbf{B}_{22}^{-1} \mathbf{B}_{21}|$$

$$\therefore |\mathbf{B}| = |-1| \left| \mathbf{A} - \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)(-1)^{-1} \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \right|$$

$$= \left| \mathbf{A} + n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \right|.$$

$$\therefore \lambda^{2/n} = \frac{|\mathbf{A}|}{\begin{vmatrix} \mathbf{A} & \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0) \\ \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' & -1 \end{vmatrix}}$$

$$= \frac{|\mathbf{A}|}{-|\mathbf{A}| \left| -1 - \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{A}^{-1} \sqrt{n}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0) \right|}$$

$$= \frac{|\mathbf{A}|}{|\mathbf{A}| \left| 1 + n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{A}^{-1} (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0) \right|}$$

$$= \frac{1}{1 + n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{A}^{-1} (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)}.$$

Where  $\mathbf{A}$  is as given in (8). But we have,

$$\mathbf{S} = \frac{1}{n-1} \mathbf{A} = \frac{1}{n-1} \sum_a (\mathbf{x}_a - \bar{\mathbf{x}})(\mathbf{x}_a - \bar{\mathbf{x}})'$$

$$\Rightarrow \mathbf{A} = (n-1) \mathbf{S}$$

$$\begin{aligned} \therefore \lambda^{2/n} &= \frac{1}{1+n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' [(\mathbf{n} - 1)\mathbf{S}]^{-1}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)} \\ &= \frac{1}{1 + \frac{n}{(\mathbf{n} - 1)}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{S}^{-1}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)} = \frac{1}{1 + T^2/(\mathbf{n} - 1)} \quad \rightarrow (9) \end{aligned}$$

where,  $T^2 = n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{S}^{-1}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)$  is Hotelling's  $T^2$ -statistic.

Now, from (7a) & (9), we can see

$$\begin{aligned} 1 + \frac{T^2}{(\mathbf{n} - 1)} &= \frac{\left| \sum_{\alpha=1}^n (\mathbf{x}_\alpha - \underline{\boldsymbol{\mu}}_0)(\mathbf{x}_\alpha - \underline{\boldsymbol{\mu}}_0)' \right|}{\left| \sum_{\alpha=1}^n (\mathbf{x}_\alpha - \bar{\mathbf{x}})(\mathbf{x}_\alpha - \bar{\mathbf{x}})' \right|} \\ \Rightarrow T^2 &= (\mathbf{n} - 1) \left[ \frac{\left| \sum_{\alpha=1}^n (\mathbf{x}_\alpha - \underline{\boldsymbol{\mu}}_0)(\mathbf{x}_\alpha - \underline{\boldsymbol{\mu}}_0)' \right|}{\left| \sum_{\alpha=1}^n (\mathbf{x}_\alpha - \bar{\mathbf{x}})(\mathbf{x}_\alpha - \bar{\mathbf{x}})' \right|} - 1 \right] \\ &= (\mathbf{n} - 1) \left[ \frac{\left| \hat{\Sigma}_0 \right|}{\left| \hat{\Sigma}_\Omega \right|} - 1 \right] \quad \rightarrow (10) \end{aligned}$$

In this formula, we need not find the inverse of a matrix, where as in the original formula we have to evaluate  $\mathbf{S}^{-1}$ .

**Theorem** :  $T^2$ - statistic is invariant (unchanged) under changes in the units of measurements for  $\mathbf{X}$  of the form,

$$\underline{\mathbf{Y}} = \mathbf{C}\mathbf{X} + \underline{\mathbf{d}}, \text{ where } \mathbf{C} \text{ is non-singular} \quad \rightarrow (1)$$

**Proof** :- We have,  $\mathbf{X} \sim N_p(\underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}})$ ,

$$\text{i.e., } E(\mathbf{X}) = \underline{\boldsymbol{\mu}} \Rightarrow E(\underline{\mathbf{Y}}) = \mathbf{C}\underline{\boldsymbol{\mu}} + \underline{\mathbf{d}} \quad (\because \text{from (1)}) \quad \rightarrow (2)$$

Now, we have the  $T^2$ - statistic for testing,  $H_0: \underline{\boldsymbol{\mu}} = \underline{\boldsymbol{\mu}}_0$  vs  $H_1: \underline{\boldsymbol{\mu}} \neq \underline{\boldsymbol{\mu}}_0$  based on the given sample  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  is

$$T_x^2 = n(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{S}_x^{-1}(\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0) \quad \rightarrow (3)$$

$$\text{where } \mathbf{S}_x = \frac{1}{\mathbf{n} - 1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' \quad \rightarrow (3a)$$

From (1) we can see  $\underline{\mathbf{Y}} \sim N_p(\mathbf{C}\underline{\boldsymbol{\mu}} + \underline{\mathbf{d}}, \mathbf{C}\underline{\boldsymbol{\Sigma}}\mathbf{C}')$ .

Now, the  $T^2$ - statistic for testing,

$$H_0: \underline{\boldsymbol{\mu}}_Y = \underline{\boldsymbol{\mu}}_{Y_0} \text{ vs } H_1: \underline{\boldsymbol{\mu}}_Y \neq \underline{\boldsymbol{\mu}}_{Y_0}$$

$$\text{where, } \underline{\boldsymbol{\mu}}_Y = \mathbf{C}\underline{\boldsymbol{\mu}} + \underline{\mathbf{d}} \quad \& \quad \underline{\boldsymbol{\mu}}_{Y_0} = \mathbf{C}\underline{\boldsymbol{\mu}}_0 + \underline{\mathbf{d}} \quad \rightarrow (4)$$

based on the sample  $\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_n$  is given by

$$T_y^2 = n (\bar{\mathbf{Y}} - \underline{\boldsymbol{\mu}}_{Y_0})' \mathbf{S}_Y^{-1} (\mathbf{Y} - \underline{\boldsymbol{\mu}}_{Y_0}) \quad \rightarrow (5)$$

where,  $\bar{\mathbf{Y}} = \mathbf{C} \bar{\mathbf{x}} + \mathbf{d}$  (from (1))

$$\mathbf{S}_y = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})' \quad \rightarrow (6)$$

and  $\underline{\boldsymbol{\mu}}_{Y_0}$  is given by (4).

In order to show that the Hotelling's  $T^2$  is invariant under the changes in the units of measurements, we have to show,  $T_y^2 = T_x^2$ .

For that, consider  $T_y^2$  given from (5),

$$\begin{aligned} T_y^2 &= n (\bar{\mathbf{y}} - \underline{\boldsymbol{\mu}}_{y_0})' \mathbf{S}_y^{-1} (\bar{\mathbf{y}} - \underline{\boldsymbol{\mu}}_{y_0}) \\ &= n (\mathbf{C}\bar{\mathbf{x}} - \mathbf{C}\underline{\boldsymbol{\mu}}_0)' \mathbf{S}_y^{-1} (\mathbf{C}\bar{\mathbf{x}} - \mathbf{C}\underline{\boldsymbol{\mu}}_0) \quad (\text{using (4)}) \\ &= n (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{C}' \mathbf{S}_y^{-1} \mathbf{C} (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0) \quad \rightarrow (7) \end{aligned}$$

$$\begin{aligned} \text{But, } \mathbf{S}_y &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})' \\ &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{C}\mathbf{x}_i - \mathbf{C}\bar{\mathbf{x}})(\mathbf{C}\mathbf{x}_i - \mathbf{C}\bar{\mathbf{x}})' \quad (\text{using (1)}) \\ &= \frac{1}{n-1} \sum_{i=1}^n \mathbf{C}(\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' \mathbf{C}' \\ &= \mathbf{C}\mathbf{S}_x\mathbf{C}' \quad (\text{from (3a)}) \end{aligned}$$

$$\Rightarrow \mathbf{S}_y^{-1} = (\mathbf{C}')^{-1} \mathbf{S}_x^{-1} \mathbf{C}^{-1}$$

$$\Rightarrow \mathbf{C}' \mathbf{S}_y^{-1} \mathbf{C} = \mathbf{S}_x^{-1}$$

Using this in (7), we get

$$T_y^2 = n (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0)' \mathbf{S}_x^{-1} (\bar{\mathbf{x}} - \underline{\boldsymbol{\mu}}_0) = T_x^2 \quad (\text{from (3)})$$

Thus,  $T^2$  is invariant under the changes in the units of measurements.

**NOTE** : The above theorem may be stated as “ The Hotellings  $T^2$  is invariant under linear transformation ( or under changes in the location and scale ) of the sample .

### Uses Of The $T^2$ -statistics :-

(i) For testing the significance of one sample mean vector  $\bar{\mathbf{x}}$  :

Suppose  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  is a random sample from a  $p$ -variate normal population  $N_p(\underline{\boldsymbol{\mu}}, \underline{\boldsymbol{\Sigma}})$ , where both  $\underline{\boldsymbol{\mu}}$  and  $\underline{\boldsymbol{\Sigma}}$  are assumed as unknown. Now, our statistical problem is whether the given sample has come from the multivariate normal population, whose mean vector is  $\underline{\boldsymbol{\mu}}_0$ . In other words, we want to test the hypothesis

$$H_0 : \underline{\boldsymbol{\mu}} = \underline{\boldsymbol{\mu}}_0 \quad \text{vs} \quad H_1 : \underline{\boldsymbol{\mu}} \neq \underline{\boldsymbol{\mu}}_0 \quad \rightarrow (1),$$

where  $\underline{\boldsymbol{\mu}}_0$  is the given mean vector.

For testing the above hypothesis, derive the test statistic.

Solution:

We have given a random sample of size  $n$  viz.,  $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$  from  $N_p(\underline{\mu}, \underline{\Sigma})$ , where both  $\underline{\mu}$  and  $\underline{\Sigma}$  are unknown.

Now, we know that the mean vector

$$\bar{\underline{X}} \sim N_p(\underline{\mu}, \underline{\Sigma}/n) \quad (\text{since } \bar{\underline{X}} \text{ is a linear function of the sample})$$

$$\text{Define the random vector, } \underline{Y} = \sqrt{n}(\bar{\underline{X}} - \underline{\mu}) \quad \rightarrow (2)$$

Whose population mean vector and population variance-covariance matrix are respectively given by

$$E(\underline{Y}) = \sqrt{n}E(\bar{\underline{X}} - \underline{\mu}) = \sqrt{n}(\underline{\mu} - \underline{\mu}) = \underline{0}.$$

$$V(\underline{Y}) = V(\sqrt{n}(\bar{\underline{X}} - \underline{\mu})) = n V(\bar{\underline{X}} - \underline{\mu}) = n V(\bar{\underline{X}}) = n \underline{\Sigma}/n = \underline{\Sigma}.$$

Thus, the mean vector of  $\underline{Y}$  is  $\underline{0}$  and covariance matrix is  $\underline{\Sigma}$ .

Further, since the transformation in (2) is linear, we have

$$\underline{Y} \sim N_p(\underline{0}, \underline{\Sigma}) \quad \rightarrow (3)$$

We have the sample variance-covariance matrix

$$\underline{S} = \frac{1}{n-1} \sum_{i=1}^n (\underline{X}_i - \bar{\underline{X}})(\underline{X}_i - \bar{\underline{X}})' \quad \rightarrow (4)$$

Now, we know that  $(n-1)\underline{S}$  follows Wishart distribution with  $n-1$  degrees freedom and parameter  $\underline{\Sigma}$  that is

$$(n-1)\underline{S} \sim W_{n-1}(\underline{\Sigma}) \quad \rightarrow (5)$$

Further, we know that the sample mean vector  $\bar{\underline{X}}$  and the sample variance-covariance matrix  $\underline{S}$  are independently distributed.

From (2), it immediately follows that the random vector  $\underline{Y}$  and the random matrix  $(n-1)\underline{S}$  distribute independently.

Now, by the definition of Hotelling's  $T^2$  distribution, the statistic

$$T^2 = \underline{Y}' \left( \frac{(n-1)\underline{S}}{(n-1)} \right)^{-1} \underline{Y} \quad \rightarrow (6)$$

follows Hotelling's  $T^2$  distribution with  $n-1$  d.f. i.e.

$$T^2 \sim T_{n-1}^2$$

Substituting (2) in (6), we can see that

$$T^2 = n(\bar{\underline{X}} - \underline{\mu})' \underline{S}^{-1} (\bar{\underline{X}} - \underline{\mu}) \sim T_{n-1}^2 \quad \rightarrow (7)$$

Now, under  $H_0 : \underline{\mu} = \underline{\mu}_0$  (7) becomes

$$T^2 = n(\bar{\underline{X}} - \underline{\mu}_0)' \underline{S}^{-1} (\bar{\underline{X}} - \underline{\mu}_0) \sim T_{n-1}^2 \quad \rightarrow (8)$$

$$\text{where, } \bar{\underline{X}} = \frac{1}{n} \sum_{i=1}^n \underline{x}_i \quad \text{and} \quad \underline{S} = \frac{1}{n-1} \sum_{i=1}^n (\underline{X}_i - \bar{\underline{X}})(\underline{X}_i - \bar{\underline{X}})'$$

Thus, the formula (8) gives us the Hotelling's  $T^2$  statistic which can be used to test (1) and follows  $T_{n-1}^2$

At the given  $\alpha$  level of significance,  $H_0$  may be rejected in favour of  $H_1$  if

$$\frac{T^2}{n-1} \left( \frac{n-p}{p} \right) > F_{p, n-p}(\alpha) \quad (\text{or}) \quad T^2 > T_0^2 \quad \rightarrow (5)$$

Where,  $T_0^2 = \frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)$  and  $F_{p, n-p}(\alpha)$  is the upper  $100\alpha^{\text{th}}$  percentile of the F-distribution and can be obtained from the F-tables.

### Nature Of $T^2$ - statistic ( Distribution ) :-

We can write,

$$T^2 = \sqrt{n} (\bar{\mathbf{X}} - \mu_0)' \left( \frac{\sum_{i=1}^n (\mathbf{X}_i - \bar{\mathbf{X}}) (\mathbf{X}_i - \bar{\mathbf{X}})'}{n-1} \right)^{-1} \sqrt{n} (\bar{\mathbf{X}} - \mu_0)$$

which is of the form,

$$\left( \begin{array}{c} \text{multivariate} \\ \text{normal r.v} \end{array} \right)' \left( \frac{\text{Wishart random matrix}}{\text{d.f}} \right)^{-1} \left( \begin{array}{c} \text{multivariate} \\ \text{normal r.v} \end{array} \right).$$

Since the multivariate normal random vector and the Wishart random matrix, given in  $T^2$  are independently distributed (  $\because \bar{\mathbf{X}}$  &  $S$  are independently distributed ). Their joint distribution is the product of the marginal normal and Wishart distributions and therefore  $T^2$ -distribution can be obtained from this.

### (ii) Computation Of Confidence Region For Mean Vector :-

The  $100(1-\alpha)$  percent confidence region for the population mean  $\underline{\mu}$  of  $N_p(\underline{\mu}, \underline{\Sigma})$  based on the given random sample  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  is given by the set of all  $\underline{\mu}$  satisfying the inequality,

$$n (\bar{\mathbf{X}} - \underline{\mu})' S^{-1} (\bar{\mathbf{X}} - \underline{\mu}) \leq T_0^2(\alpha) \quad , \text{where } T_0^2(\alpha) = \frac{(n-1)p}{n-p} F_{p, n-p}(\alpha).$$

### (iii) A Two Sample Problem (An application of Hotelling $T^2$ -statistic) :-

Another situation in which the  $T^2$ -statistic is used is that in which the null hypothesis is that the mean of one normal population is equal to the mean of the other, where the covariance matrices are assumed equal but unknown.

Suppose  $\mathbf{X}_1^{(1)}, \mathbf{X}_2^{(1)}, \dots, \mathbf{X}_n^{(1)}$  is a sample from  $N_p(\underline{\mu}^{(1)}, \underline{\Sigma})$  and  $\mathbf{X}_1^{(2)}, \mathbf{X}_2^{(2)}, \dots, \mathbf{X}_n^{(2)}$  is a another sample (independent of the first sample) from  $N_p(\underline{\mu}^{(2)}, \underline{\Sigma})$ .



Now, we wish to test the null hypothesis,

$$H_0 : \underline{\mu}^{(1)} = \underline{\mu}^{(2)} \text{ or } \underline{\mu}^{(1)} - \underline{\mu}^{(2)} = \mathbf{0}, \text{ against } H_1 : \underline{\mu}^{(1)} \neq \underline{\mu}^{(2)} \rightarrow (1)$$

The sample means from the hypothesis ,

$$\bar{\underline{X}}^{(1)} = \frac{1}{n_1} \sum_{i=1}^{n_1} \underline{X}_i^{(1)} \sim N_p(\underline{\mu}^{(1)}, \underline{\Sigma}/n_1)$$

$$\text{and } \bar{\underline{X}}^{(2)} = \frac{1}{n_2} \sum_{i=1}^{n_2} \underline{X}_i^{(2)} \sim N_p(\underline{\mu}^{(2)}, \underline{\Sigma}/n_2)$$

$$\text{Now define, } \underline{Y} = \bar{\underline{X}}^{(1)} - \bar{\underline{X}}^{(2)} \rightarrow (2)$$

$$\text{with mean, } E(\underline{Y}) = E(\bar{\underline{X}}^{(1)}) - E(\bar{\underline{X}}^{(2)}) = \underline{\mu}^{(1)} - \underline{\mu}^{(2)} \rightarrow (3)$$

and the variance- covariance matrix,

$$V(\underline{Y}) = V(\bar{\underline{X}}^{(1)}) + V(\bar{\underline{X}}^{(2)}) \quad (\because \text{The two samples are independent})$$

$$= \frac{1}{n_1} \underline{\Sigma} + \frac{1}{n_2} \underline{\Sigma}$$

$$= \left( \frac{1}{n_1} + \frac{1}{n_2} \right) \underline{\Sigma} \rightarrow (4)$$

Since the transformation in (2) is linear, from (3) and (4) it follows

$$\begin{aligned} \underline{Y} &\sim N_p\left(\underline{\mu}^{(1)} - \underline{\mu}^{(2)}, \left(\frac{1}{n_1} + \frac{1}{n_2}\right) \underline{\Sigma}\right) \\ \text{i.e., } \underline{Y} - (\underline{\mu}^{(1)} - \underline{\mu}^{(2)}) &\sim N_p\left(\mathbf{0}, \left(\frac{1}{n_1} + \frac{1}{n_2}\right) \underline{\Sigma}\right) \\ \text{i.e., } \frac{(\bar{\underline{X}}^{(1)} - \bar{\underline{X}}^{(2)}) - (\underline{\mu}^{(1)} - \underline{\mu}^{(2)})}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} &\sim N_p(\mathbf{0}, \underline{\Sigma}) \quad (\text{using (2)}) \rightarrow (5) \end{aligned}$$

The sample covariance matrix from sample 1, which is denoted by  $\mathbf{S}_1$  and is given by

$$\mathbf{S}_1 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (\underline{X}_i^{(1)} - \bar{\underline{X}}^{(1)}) (\underline{X}_i^{(1)} - \bar{\underline{X}}^{(1)})'$$

Similarly, the sample covariance matrix from sample 2, denoted by  $\mathbf{S}_2$  and is given by

$$\mathbf{S}_2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (\underline{X}_i^{(2)} - \bar{\underline{X}}^{(2)}) (\underline{X}_i^{(2)} - \bar{\underline{X}}^{(2)})'$$

$$\text{Let us denote, } \mathbf{S} = \frac{(n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2}{n_1 + n_2 - 2} \rightarrow (6)$$

We know that  $(n_1-1)\mathbf{S}_1$  and  $(n_2-1)\mathbf{S}_2$  are Wishart random matrices and are distributed as  $w_{n_1-1}(\boldsymbol{\Sigma})$  &  $w_{n_2-1}(\boldsymbol{\Sigma})$  respectively, where  $w_{n_1-1}(\boldsymbol{\Sigma})$  is Wishart distribution with  $(n_1 - 1)$  d.f and  $w_{n_2-1}(\boldsymbol{\Sigma})$  is Wishart distribution with  $(n_2 - 1)$  d.f. both have the parametric matrix  $\boldsymbol{\Sigma}$ .

By assumption, the samples are independent, so  $(n_1-1)\mathbf{S}_1$  and  $(n_2-1)\mathbf{S}_2$  are also independent. Therefore from (6),  $(n_1+n_2-2)\mathbf{S}$  is distributed as Wishart distribution with  $n_1 + n_2 - 2$  d.f and with the parametric matrix  $\boldsymbol{\Sigma}$ , i.e.  $(n_1+n_2-2)\mathbf{S} \sim w_{n_1+n_2-2}(\boldsymbol{\Sigma}) \rightarrow (7)$

Since, the sample variance-covariance matrix is independently distributed with the sample mean vector,  $\mathbf{S}_1$  is independently distributed with  $\bar{\mathbf{X}}^{(1)}$  and since the two samples are independent,  $\mathbf{S}_1$  is independently distributed with  $\bar{\mathbf{X}}^{(2)}$  and therefore  $\mathbf{S}_1$  is independently distributed with  $\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}$ . Similarly,  $\mathbf{S}_2$  is independently distributed with  $\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}$ .

Therefore,  $\mathbf{S} = \frac{(n_1-1)\mathbf{S}_1 + (n_2-1)\mathbf{S}_2}{n_1 + n_2 - 2}$  is independently distributed with  $\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}$ .

Thus, from the above explanation and from (5) & (6) and by the definition of  $T^2$ -distribution, we have

$$T^2 = \left[ \frac{(\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) - (\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)})}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right]' \mathbf{S}^{-1} \left[ \frac{(\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) - (\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)})}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \right]$$

$$= \left( \frac{n_1 n_2}{n_1 + n_2} \right) [(\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) - (\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)})]' \mathbf{S}^{-1} [(\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) - (\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)})] \rightarrow (8)$$

is distributed as  $T^2$ -distribution with  $n_1 + n_2 - 2$  d.f.

Now, by virtue of the relation between  $T^2$  and F – distribution, we have

$$\frac{T^2}{n_1 + n_2 - 2} \sim \frac{p}{n_1 + n_2 - 2 - (p-1)} F_{p, n_1 + n_2 - 2 - (p-1)}$$

$$\text{i.e., } \frac{T^2}{n_1 + n_2 - 2} \sim \frac{p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}$$

under  $H_0 : \underline{\boldsymbol{\mu}}^{(1)} = \underline{\boldsymbol{\mu}}^{(2)}$  i.e.,  $\underline{\boldsymbol{\mu}}^{(1)} - \underline{\boldsymbol{\mu}}^{(2)} = \mathbf{0}$ , (8) becomes

$$T^2 = \left( \frac{n_1 n_2}{n_1 + n_2} \right) (\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) \mathbf{S}^{-1} (\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)}) \rightarrow (9)$$

if  $T^2 > \frac{p(n_1 + n_2 - 2)}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}(\alpha)$ , where  $F_{p, n_1 + n_2 - p - 1}(\alpha)$  is table F-value at  $\alpha$  level of significance with  $(p, n_1 + n_2 - p - 1)$  d.f., then  $H_0 : \underline{\mu}^{(1)} = \underline{\mu}^{(2)}$  may be rejected .

### Confidence Interval :-

An  $100(1-\alpha)\%$  confidence region for estimation of  $\underline{\mu}^{(1)} - \underline{\mu}^{(2)}$  is given by the following set ( $\because$  from (8)),

$$\left\{ \underline{\mathbf{m}} / \left( \frac{n_1 n_2}{n_1 + n_2} \right) \left( \underline{\bar{\mathbf{X}}}^{(1)} - \underline{\bar{\mathbf{X}}}^{(2)} - \underline{\mathbf{m}} \right) \mathbf{S}^{-1} \left( \underline{\bar{\mathbf{X}}}^{(1)} - \underline{\bar{\mathbf{X}}}^{(2)} - \underline{\mathbf{m}} \right) < T_0^2 \right\}$$

where,  $T_0^2 = \frac{p(n_1 + n_2 - 2)}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}(\alpha)$ .

### The Two Sample Situation, When $\Sigma_1 \neq \Sigma_2$ :-

In the above problem, we have assumed that the covariance matrices of both the populations are assumed as equal i.e.,  $\Sigma_1 = \Sigma_2 = \Sigma$ .

Now, let us suppose that  $\Sigma_1 \neq \Sigma_2$  i.e., the population covariance matrices are not equal.

In this case, no tests are available for making inferences about  $\underline{\mu}^{(1)} - \underline{\mu}^{(2)}$ , when the sizes of the samples are small. However, if  $n_1$  &  $n_2$  are large i.e., in case of large samples, we have the following result.

### Result :-

Let the sample sizes be such that  $n_1 - p$  and  $n_2 - p$  are large. An approximation  $100(1 - \alpha)\%$  confidence region for  $\underline{\mu}^{(1)} - \underline{\mu}^{(2)}$  is given by all  $\underline{\mu}^{(1)} - \underline{\mu}^{(2)}$  satisfying,

$$\left[ \left( \underline{\bar{\mathbf{X}}}^{(1)} - \underline{\bar{\mathbf{X}}}^{(2)} \right) - \left( \underline{\mu}^{(1)} - \underline{\mu}^{(2)} \right) \right]' \left[ \frac{1}{n_1} \mathbf{S}_1 + \frac{1}{n_2} \mathbf{S}_2 \right]^{-1} \left[ \left( \underline{\bar{\mathbf{X}}}^{(1)} - \underline{\bar{\mathbf{X}}}^{(2)} \right) - \left( \underline{\mu}^{(1)} - \underline{\mu}^{(2)} \right) \right] \leq \chi_p^2(\alpha)$$

where,  $\chi_p^2(\alpha)$  is  $\chi^2$ -table values with p.d.f at  $100\alpha\%$  level of significance.

**Proof:-**  $E(\underline{\bar{\mathbf{X}}}^{(1)} - \underline{\bar{\mathbf{X}}}^{(2)}) = \underline{\mu}^{(1)} - \underline{\mu}^{(2)}$

$$\& V(\underline{\bar{\mathbf{X}}}^{(1)} - \underline{\bar{\mathbf{X}}}^{(2)}) = V(\underline{\bar{\mathbf{X}}}^{(1)}) + V(\underline{\bar{\mathbf{X}}}^{(2)}) = \frac{1}{n_1} \Sigma_1 + \frac{1}{n_2} \Sigma_2.$$

By the central limit theorem,

$$\bar{\tilde{X}}^{(1)} - \bar{\tilde{X}}^{(2)} \sim N_p \left( \tilde{\mu}^{(1)} - \tilde{\mu}^{(2)}, \frac{1}{n_1} \Sigma_1 + \frac{1}{n_2} \Sigma_2 \right).$$

If  $\Sigma_1$  &  $\Sigma_2$  are known,

$$\left[ (\bar{\tilde{X}}^{(1)} - \bar{\tilde{X}}^{(2)}) - (\tilde{\mu}^{(1)} - \tilde{\mu}^{(2)}) \right] \left[ \frac{1}{n_1} \Sigma_1 + \frac{1}{n_2} \Sigma_2 \right]^{-1}$$

$$\left[ (\bar{\tilde{X}}^{(1)} - \bar{\tilde{X}}^{(2)}) - (\tilde{\mu}^{(1)} - \tilde{\mu}^{(2)}) \right] \sim \chi_p^2(\alpha),$$

approximately, when  $n_1$  &  $n_2$  are large, with high probability  $S_1 \rightarrow \Sigma_1$  and  $S_2 \rightarrow \Sigma_2$ . Consequently, the approximation holds with  $S_1$  &  $S_2$ , in place of  $\Sigma_1$  and  $\Sigma_2$  respectively.

Hence the theorem.

### NULL DISTRIBUTION OF HOTELLING $T^2$

Statement:- Let  $\tilde{\mathbf{Y}} \sim N_p(\tilde{\mathbf{y}}, \Sigma)$  and let  $\mathbf{A}$  be a Wishart random matrix

independently distributed as  $\sum_{\alpha=1}^m \tilde{\mathbf{Z}}_{\alpha} \tilde{\mathbf{Z}}_{\alpha}'$ , where  $\tilde{\mathbf{Z}}_{\alpha}$ 's are

i.i.d  $\sim N_p(\mathbf{0}, \Sigma)$ . Also let

$$T^2 = m \tilde{\mathbf{Y}}' \mathbf{A}^{-1} \tilde{\mathbf{Y}} \quad \rightarrow (1)$$

then,  $\frac{T^2}{m} \left( \frac{m-p+1}{p} \right)$  is distributed as a non-central  $F$  with

$p$  and  $m-p+1$  d.f. and non-centrality parameter  $\tilde{\mathbf{y}}' \Sigma^{-1} \tilde{\mathbf{y}}$ . Further, if  $\tilde{\mathbf{y}} = \mathbf{0}$ , then

$$\frac{T^2}{m} \left( \frac{m-p+1}{p} \right) \sim F_{p, m-p+1} \quad \rightarrow (2)$$

and the distribution of  $T^2$  is called  $T^2$ -distribution.

Proof :- Since  $\Sigma$  is positive definite, there exists a non-singular  $\mathbf{C}$  such that

$$\mathbf{C} \Sigma \mathbf{C}' = \mathbf{I}_p \text{ so that, } \Sigma = (\mathbf{C}' \mathbf{C})^{-1} \quad \rightarrow (3)$$

Define,  $\tilde{\mathbf{Y}}^* = \mathbf{C} \tilde{\mathbf{Y}}$  and  $\mathbf{A}^* = \mathbf{C} \mathbf{A} \mathbf{C}'$   $\rightarrow (4)$

We can see that,  $\mathbf{E}(\tilde{\mathbf{Y}}^*) = \mathbf{C} \tilde{\mathbf{y}} = \tilde{\mathbf{y}}^*$  (say)

$$V(\tilde{\mathbf{Y}}^*) = V(\mathbf{C} \tilde{\mathbf{Y}}) = \mathbf{C} \Sigma \mathbf{C}' = \mathbf{I} \quad (\text{using (3)}) \quad \rightarrow (5)$$

Thus,  $\tilde{\mathbf{Y}}^* \sim N_p(\tilde{\mathbf{y}}^*, \mathbf{I})$ .

Since  $\mathbf{A}$  is distributed as  $\sum_{\alpha=1}^m \mathbf{Z}_{\alpha} \mathbf{Z}_{\alpha}'$ ,  $\mathbf{A}^* = \mathbf{C} \mathbf{A} \mathbf{C}'$  is distributed as

$$\mathbf{C} \sum_{\alpha=1}^m \mathbf{Z}_{\alpha} \mathbf{Z}_{\alpha}' \mathbf{C}' = \sum_{\alpha=1}^m \mathbf{Z}_{\alpha}^* \mathbf{Z}_{\alpha}^{*'} \quad \rightarrow (6)$$

where,  $\mathbf{Z}_{\alpha}^* = \mathbf{C} \mathbf{Z}_{\alpha} \sim N_p(\mathbf{C} \mathbf{0}, \mathbf{C} \Sigma \mathbf{C}') = N_p(\mathbf{0}, \mathbf{I})$ .

Eq (1) can be written as

$$\begin{aligned} T^2 &= m \mathbf{Y}' \mathbf{A}^{-1} \mathbf{Y} \\ &= m \mathbf{Y}' \mathbf{C}' (\mathbf{C}')^{-1} \mathbf{A}^{-1} (\mathbf{C})^{-1} \mathbf{C} \mathbf{Y} \\ &= m (\mathbf{C} \mathbf{Y})' (\mathbf{C} \mathbf{A} \mathbf{C}')^{-1} (\mathbf{C} \mathbf{Y}) \\ &= m \mathbf{Y}^{*'} \mathbf{A}^{*'} \mathbf{Y}^* \end{aligned} \quad \rightarrow (7)$$

where,  $\mathbf{Y}^* \sim N_p(\mathbf{0}, \mathbf{I})$  and  $\mathbf{A}^*$  is independently distributed as  $\sum_{\alpha=1}^m \mathbf{Z}_{\alpha}^* \mathbf{Z}_{\alpha}^{*'}$  in

which  $\mathbf{Z}_{\alpha}^*$ 's are i.i.d  $\sim N_p(\mathbf{0}, \mathbf{I})$ .

Also, since  $\mathbf{Y}$  and  $\mathbf{A}$  are independently distributed, from eq. (4),

$\mathbf{Y}^*$  and  $\mathbf{A}^*$  are also independently distributed.

Let  $\mathbf{\Omega} = (\omega_{ij})_{p \times p}$  is an orthogonal matrix in which first row is defined by

$$\omega_{1j} = \frac{Y_j^*}{\sqrt{\mathbf{Y}^{*'} \mathbf{Y}^*}}, \quad j=1,2,\dots,p \quad \rightarrow (8)$$

where,  $Y_j^*$  is  $j^{\text{th}}$  component of  $\mathbf{Y}^*$ .

Now define,  $\mathbf{U} = \mathbf{\Omega} \mathbf{Y}^*$

$$\mathbf{B} = \mathbf{\Omega} \mathbf{A}^* \mathbf{\Omega}' \quad \rightarrow (9)$$

The  $i^{\text{th}}$  component of  $\mathbf{U}$  is given by

$$\begin{aligned}
U_i &= \sum_{j=1}^p \omega_{ij} Y_j^* \\
&= \sqrt{\tilde{\mathbf{Y}}^{*'} \tilde{\mathbf{Y}}^*} \sum_{j=1}^p \omega_{ij} \omega_{1j} \quad [\text{using (8)}] \\
&= \begin{cases} \sqrt{\tilde{\mathbf{Y}}^{*'} \tilde{\mathbf{Y}}^*} & \text{if } i = 1 \\ 0 & \text{otherwise} \end{cases} \\
&\quad (\text{Since , } \mathbf{\Omega} \text{ is orthogonal matrix})
\end{aligned}$$

$$\text{Thus ,} \quad \tilde{\mathbf{U}} = \begin{bmatrix} \sqrt{\tilde{\mathbf{Y}}^{*'} \tilde{\mathbf{Y}}^*} \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \rightarrow (10)$$

From equation (7) ,

$$\begin{aligned}
\frac{T^2}{m} &= \tilde{\mathbf{Y}}^{*'} \mathbf{I} \mathbf{A}^* \mathbf{I} \tilde{\mathbf{Y}}^* \\
&= \tilde{\mathbf{Y}}^{*'} \mathbf{\Omega}' \mathbf{\Omega} \mathbf{A}^* \mathbf{\Omega}' \mathbf{\Omega} \tilde{\mathbf{Y}}^* \quad (\because \mathbf{\Omega} \text{ is orthogonal}) \\
&= (\mathbf{\Omega} \tilde{\mathbf{Y}}^*)' (\mathbf{\Omega} \mathbf{A}^* \mathbf{\Omega}')^{-1} (\mathbf{\Omega} \tilde{\mathbf{Y}}^*) \quad (\because \mathbf{\Omega}^{-1} = \mathbf{\Omega}') \\
&= \tilde{\mathbf{U}}' \mathbf{B}^{-1} \tilde{\mathbf{U}} \quad [\text{using (9)}] \\
&= \left( \begin{array}{cccc} \sqrt{\tilde{\mathbf{Y}}^{*'} \tilde{\mathbf{Y}}^*} & 0 & \dots & 0 \end{array} \right) \begin{bmatrix} b^{11} & b^{12} & \dots & b^{1p} \\ b^{21} & b^{22} & \dots & b^{2p} \\ \vdots & \vdots & & \vdots \\ b^{p1} & b^{p2} & \dots & b^{pp} \end{bmatrix} \begin{bmatrix} \sqrt{\tilde{\mathbf{Y}}^{*'} \tilde{\mathbf{Y}}^*} \\ 0 \\ \vdots \\ 0 \end{bmatrix} \\
&= b^{11} \tilde{\mathbf{Y}}^{*'} \tilde{\mathbf{Y}}^* \quad \rightarrow (11)
\end{aligned}$$

where ,  $b^{11}$  is first diagonal element of  $\mathbf{B}^{-1}$ .  
But we know that ,

$$b^{11} = \frac{1}{b_{11} - \underline{\mathbf{b}}_1' \mathbf{B}_{22}^{-1} \underline{\mathbf{b}}_1}, \quad \text{where } \mathbf{B} = \begin{bmatrix} b_{11} & \underline{\mathbf{b}}_1' \\ \underline{\mathbf{b}}_1 & b_{22} \end{bmatrix}$$

Thus, from eq (11),

$$\frac{T^2}{m} = \frac{\underline{\mathbf{Y}}^{*'} \underline{\mathbf{Y}}^*}{b_{11.2\dots p}} \rightarrow (12)$$

where,  $b_{11.2\dots p} = b_{11} - \underline{\mathbf{b}}_1' \mathbf{B}_{22}^{-1} \underline{\mathbf{b}}_1$ .

Let us suppose that  $\mathbf{\Omega}$  is fixed (given). Then, just as we show  $\mathbf{A}^*$  is distributed as  $\sum_{\alpha=1}^m \underline{\mathbf{Z}}_{\alpha}^* \underline{\mathbf{Z}}_{\alpha}^{*'}$ , we can show that  $\mathbf{\Omega} \mathbf{A}^* \mathbf{\Omega}'$  is distributed as

$$\sum_{\alpha=1}^m \underline{\mathbf{V}}_{\alpha} \underline{\mathbf{V}}_{\alpha}', \quad \text{when } \underline{\mathbf{V}}_{\alpha} = \mathbf{\Omega} \underline{\mathbf{Z}}_{\alpha}^* \text{ and } \underline{\mathbf{V}}_{\alpha}'\text{'s are i.i.d } \sim N_p(\mathbf{0}, \mathbf{I}).$$

Now, with little difficult, we may show that  $b_{11.2\dots p} = b_{11} - \underline{\mathbf{b}}_1' \mathbf{B}_{22}^{-1} \underline{\mathbf{b}}_1$  is

$$\text{conditionally distributed as } \sum_{\alpha=1}^{m-(p-1)} \omega_{\alpha}^2,$$

where *each*  $\omega_{\alpha}$  is *i.i.d*  $\sim N(0,1)$ .

$$\text{Therefore, } \sum_{\alpha=1}^{m-(p-1)} \omega_{\alpha}^2 \sim \chi_{m-(p-1)}^2.$$

More over, the conditional distribution of  $b_{11.2\dots p}$  does not depend on  $\mathbf{\Omega}$ , we have  $b_{11.2\dots p}$  is unconditionally distributed as  $\chi_{m-(p-1)}^2$ .

$$\text{Also, since } \underline{\mathbf{Y}}^* \sim N_p(\underline{\mathbf{y}}^*, \mathbf{I}), \quad \underline{\mathbf{Y}}^{*'} \underline{\mathbf{Y}}^* = \sum_{i=1}^p Y_i^{*2}$$

where,  $Y_i^* \sim N(v_i^*, 1)$  and  $Y_i^*$ 's are independent.

$$\begin{aligned} \text{Thus, } \underline{\mathbf{Y}}^{*'} \underline{\mathbf{Y}}^* & \text{ has non-central } \chi^2 \text{-distribution with non-centrality} \\ & = \sum_{i=1}^p v_i^{*2} = \underline{\mathbf{y}}^{*'} \underline{\mathbf{y}}^* = \underline{\mathbf{y}}' \mathbf{C}' \mathbf{C} \underline{\mathbf{y}} \quad [\text{from(5)}] \\ & = \underline{\mathbf{y}}' \mathbf{\Sigma}^{-1} \underline{\mathbf{y}} \quad [\text{from (3)}] \end{aligned}$$

Thus,  $\frac{T^2}{m}$  is distributed as the ratio of non-central  $\chi_p^2$  and an independent central  $\chi_{m-(p-1)}^2$ .

Thus,  $\frac{T^2}{m} \left[ \frac{m - (p - 1)}{p} \right] \sim F_{(p, m-p+1)}$  (non-central) and non-centrality parameter  $\mathbf{y}'\Sigma^{-1}\mathbf{y}$ .

If  $\mathbf{y} = \mathbf{0}$  then,  $\mathbf{y}'\Sigma^{-1}\mathbf{y} = \mathbf{0}$  and therefore in this case, the distribution is central.  $F_{(p, m-p+1)}$ , the distribution of  $T^2$  is called  $T^2$ -distribution with 'm' degrees of freedom.

## COMPARISON OF SEVERAL MULTIVARIATE POPULATION MEANS

### MANOVA for one way classification:

Suppose we have 'g' populations, each is distributed multivariate normal with mean vectors  $\underline{\mu}_1, \underline{\mu}_2, \dots, \underline{\mu}_g$  respectively. Let us suppose that all populations have the same covariance matrix  $\Sigma$ . Thus, we have the 'g' populations.

$$\Pi_1 \sim N_p(\underline{\mu}_1, \Sigma)$$

$$\Pi_2 \sim N_p(\underline{\mu}_2, \Sigma)$$

⋮

$$\Pi_g \sim N_p(\underline{\mu}_g, \Sigma)$$

Now, we have a sample of size 'n<sub>i</sub>' from  $i^{th}$  population  $\Pi_i$ . Thus, we have 'g' samples from the 'g' populations as follows:

$$\text{Population } \Pi_1 : \underline{X}_{11}, \underline{X}_{12}, \dots, \underline{X}_{1n_1}$$

$$\text{Population } \Pi_2 : \underline{X}_{21}, \underline{X}_{22}, \dots, \underline{X}_{2n_2}$$

⋮

$$\text{Population } \Pi_g : \underline{X}_{g1}, \underline{X}_{g2}, \dots, \underline{X}_{gn_g}$$

Using the above random samples, MANOVA is used to investigate whether the population mean vectors are same and if not, which mean components differ significantly. Thus, the null hypothesis is

$$H_0 : \underline{\mu}_1 = \underline{\mu}_2 = \dots = \underline{\mu}_g \rightarrow (1)$$

#### ASSUMPTIONS CONCERNING THE STRUCTURE OF THE DATA :-

- (1)  $\underline{X}_{i1}, \underline{X}_{i2}, \dots, \underline{X}_{in_i}$  is a random sample of size  $n_i$  from a population with mean  $\underline{\mu}_i, i=1, 2, \dots, g$ . The random samples from different populations are independent.
- (2) All populations have a common covariance matrix  $\Sigma$ .



(3) Each population is a multivariate normal. Condition (3) can be relaxed by appealing to the central limit theorem, when the sample sizes  $n_i$ 's are large. If the mean vector of  $i^{th}$  population is written as

$$\underline{\mu}_i = \underline{\mu} + \underline{\tau}_i \quad \rightarrow (2)$$

Here,  $\underline{\mu}$  is the overall mean vector of all population and  $\underline{\tau}_i$  is a component due to the specific population, then the null hypothesis (1) can be written as

$$H_0 : \underline{\tau}_1 = \underline{\tau}_2 = \dots = \underline{\tau}_g = \mathbf{0} \quad \rightarrow (3)$$

The response  $\underline{X}_{ij}$ , distributed as  $N_p(\underline{\mu} + \underline{\tau}_i, \underline{\Sigma})$ , can be expressed in the suggestive form,

$$\underline{X}_{ij} = \underline{\mu} + \underline{\tau}_i + \underline{\varepsilon}_{ij} \quad \rightarrow (4)$$

$$\left( \begin{array}{c} \text{overall} \\ \text{mean} \end{array} \right) \left( \begin{array}{c} \text{treatment} \\ \text{effect} \end{array} \right) \left( \begin{array}{c} \text{random} \\ \text{error} \end{array} \right)$$

$i = 1, 2, \dots, g$  &  $j = 1, 2, \dots, n_i$ .

where,  $\underline{\varepsilon}_{ij} \sim N(\mathbf{0}, \underline{\Sigma})$  are independent random variables. (4) is called as MANOVA model for comparing of population mean vectors. Here  $\underline{\mu}$  is overall mean vector and  $\underline{\tau}_i$  represents the  $i^{th}$  treatment effect with

$$\sum_{i=1}^g n_i \underline{\tau}_i = \mathbf{0} \quad \rightarrow (5)$$

A vector of observations may be decomposed as suggested by model (4). Thus,

$$\underline{x}_{ij} = \bar{\underline{x}} + (\underline{x}_i - \bar{\underline{x}}) + (\underline{x}_{ij} - \underline{x}_i) \quad \rightarrow (6)$$

$$\left( \begin{array}{c} \text{observation} \end{array} \right) \left( \begin{array}{c} \text{overall} \\ \text{sample} \\ \text{mean } \underline{\mu} \end{array} \right) \left( \begin{array}{c} \text{estimated} \\ \text{treatment} \\ \text{effect } \underline{\tau}_i \end{array} \right) \left( \begin{array}{c} \text{residual} \\ \hat{\underline{\varepsilon}}_{ij} \end{array} \right)$$

When  $\bar{\underline{x}}_i =$  mean of  $i^{th}$  sample  $\underline{x}_{i1}, \underline{x}_{i2}, \dots, \underline{x}_{in_i}$

$$\bar{\underline{x}} = \frac{1}{g} (\bar{\underline{x}}_1 + \bar{\underline{x}}_2 + \dots + \bar{\underline{x}}_g) \quad (\text{general mean})$$

From (6), we may write the cross product,

$$\begin{aligned} (\underline{x}_{ij} - \bar{\underline{x}})(\underline{x}_{ij} - \bar{\underline{x}})' &= ((\underline{x}_{ij} - \underline{x}_i) + (\underline{x}_i - \bar{\underline{x}}))((\underline{x}_{ij} - \underline{x}_i) + (\underline{x}_i - \bar{\underline{x}}))' \\ &= (\underline{x}_{ij} - \underline{x}_i)(\underline{x}_{ij} - \underline{x}_i)' + (\underline{x}_{ij} - \underline{x}_i)(\underline{x}_i - \bar{\underline{x}})' \\ &\quad + (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})' + (\underline{x}_i - \bar{\underline{x}})(\underline{x}_{ij} - \underline{x}_i)' \end{aligned}$$

Summing the cross product over  $i$  and  $j$ , we get

$$\begin{aligned}
\sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}})(\mathbf{x}_{ij} - \bar{\mathbf{x}})' &= \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)' \\
&+ \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})' \\
&+ \sum_{i=1}^g \sum_{j=1}^{n_i} (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})' \\
&+ \sum_{i=1}^g \sum_{j=1}^{n_i} (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})' \quad \rightarrow (7)
\end{aligned}$$

But, since  $\sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i) = \mathbf{0}$ , Eq (7) becomes ,

$$\begin{aligned}
\sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}})(\mathbf{x}_{ij} - \bar{\mathbf{x}})' &= \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)' \\
&+ \sum_{i=1}^g n_i (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})' \quad \rightarrow (8)
\end{aligned}$$

$$\Rightarrow \begin{pmatrix} \text{total(corrected)} \\ \text{sum of square} \\ \text{\& cross products} \end{pmatrix} = \begin{pmatrix} \text{residual(within)} \\ \text{sum of squares} \\ \text{\& cross products} \end{pmatrix} + \begin{pmatrix} \text{treatment(between)} \\ \text{sum of squares} \\ \text{\& cross products} \end{pmatrix}$$

That is (8) may be written as

$$\sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}})(\mathbf{x}_{ij} - \bar{\mathbf{x}})' = W + B \quad \rightarrow (9)$$

$$\begin{aligned}
\text{where , } W &= \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)' \\
&= (n_1 - 1)S_1 + (n_2 - 1)S_2 + \dots + (n_g - 1)S_g
\end{aligned}$$

where ,  $S_i$  is sample covariance matrix of  $i^{th}$  sample .

$$\text{and } B = \sum_{i=1}^g n_i (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})' .$$

Now, we summarise the calculations leading to the test statistic in a MANOVA table .  
MANOVA table for comparing population mean vectors :-

Source of variation	Matrix of sum of squares & cross product	Degrees of freedom
---------------------	--	--------------------

Treatments	$B = \sum_{i=1}^g n_i (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})'$	g-1
Residual(error)	$W = \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)'$	n-g
Total (correlated for the mean )	$B + W = \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}})(\mathbf{x}_{ij} - \bar{\mathbf{x}})'$	n-1

Now one of the test statistic for testing (3) involves generalized variances and is given by

$$\Lambda^* = \frac{|W|}{|B+W|} \rightarrow (10)$$

The quantity  $\Lambda^*$  is called Wilk's lamda and related to likelihood ratio criterion. The exact distribution of  $\Lambda^*$  can be derived for the special cases listed in the following table .

Distribution of Wilk's lamda ,  $\Lambda^*$  :-

No . of variables	No .of groups	Sampling distribution for multivariate normal data
p - 1	$g \geq 2$	$\left[ \frac{\sum_{i=1}^g n_i - g}{g-1} \left[ \frac{1-\Lambda^*}{\Lambda^*} \right] \right] \sim F_{\left( g-1, \sum_{i=1}^g n_i - g \right)}$
p = 2	$g \geq 2$	$\left[ \frac{\sum_{i=1}^g n_i - g - 1}{g-1} \left[ \frac{1-\sqrt{\Lambda^*}}{\sqrt{\Lambda^*}} \right] \right] \sim F_{\left( 2(g-1), 2 \sum_{i=1}^g n_i - g - 1 \right)}$
p $\geq$ 1	$g = 2$	$\left[ \frac{\sum_{i=1}^g n_i - p - 1}{p} \left[ \frac{1-\Lambda^*}{\Lambda^*} \right] \right] \sim F_{\left( p, \sum_{i=1}^g n_i - p - 1 \right)}$

$p \geq 1$	$g = 3$	$\left[ \frac{\sum_{i=1}^g n_i - p - 2}{p} \right] \left[ \frac{1 - \sqrt{\Lambda^*}}{\sqrt{\Lambda^*}} \right] \sim F \left( 2p, 2 \left( \sum_{i=1}^g n_i - p - 2 \right) \right)$
------------	---------	--

Bartlett has shown that if  $H_0$  is true and  $\sum_{i=1}^g n_i = n$  is large,

$$-\left[ n - 1 - \frac{(p+g)}{2} \right] \ln \Lambda^* = -\left[ n - 1 - \frac{(p+g)}{2} \right] \ln \frac{|W|}{|B+W|}$$

has approximately a  $\chi^2$ -distribution with  $p(g-1)$  d.f. consequently.

# Principle Component Analysis

## **Introduction :**

Suppose  $X_1, X_2, \dots, X_p$  are the given random variables. Then, principle component analysis (P.C.A) is concerned with explaining the variance-covariance structure of the variables through a few standardized linear combinations (SLC) of the original variables (we call a linear combination  $l_1X_1 + l_2X_2 + \dots + l_pX_p$  as an SLC if  $\sum_i l_i^2 = 1$  ).

Algebraically principal components (PCs) are particular standard linear combinations (SLCs) of the components of the original pattern and geometrically, these LCs represent the selection of new coordinate system obtained by rotating the original system with  $X_1, X_2, \dots, X_p$  as the coordinate axes. The new axes represent the directions with maximum variability and provide a simpler and more parsimonious (avoiding of excess) description of the covariance structure . As we shall see , principle components depend solely on the covariance matrix (or the correlation matrix) of the random variables  $X_1, X_2, \dots, X_p$ . Their development does not require a multivariate normal assumption.

The general objections of P.C.A are

- (i) Data-reduction and
- (ii) Interpretation.

Although p components required to reproduce the total system variability, often much of this variability can be accounted by a small number 'k'(<p) of the principal components. If there is almost as much information in the k components as there is in the original 'p' variables, then the 'k' principal components replace the original 'p' components of the pattern. And the original data set consisting of n measurements on p-component pattern is reduced to one consisting of n measurements on k-principal component pattern. In other

words, PCA reduces the dimensionality of the given data, losing as little information as possible. This technique was developed by Hotelling(1933).

An analysis of principle components often reveals relationships that were not previously suspected and there by allows interpretations that would not ordinarily result. In other words, the key problem is the interpretation of the principle components.

PCs may be inputs to a multiple regression analysis or cluster analysis. Moreover , (scaled) principle components are one factoring of the covariance matrix for the factor analysis model.

An example:-

Suppose we consider a sample of  $n$  students and they are asked to write five papers mechanics ( $X_1$ ), vectors ( $X_2$ ), algebra ( $X_3$ ), analysis( $X_4$ ) and statistics ( $X_5$ ). The examination in the first two papers is conducted in the closed book system , where as in the remaining three papers in the open book system.

Thus , we have totally '5n' observations so that  $n$  observations on each paper . One question which can be asked concerning this data is how the results on the five different papers should be combined to produce an overall score various answers are possible. One obvious answer would be to use the overall mean that is the linear combination  $(X_1 + X_2 + X_3 + X_4 + X_5)/5$ . But , can one do better than this? . This is one of the questions that principle component analysis seeks to answer .

**Definition Of Principle Component :-**

**Definition:-**

If  $\underline{X}$  is a random vector with mean  $\underline{\mu}$  and variance – covariance matrix  $\Sigma$ , then the principle component transformation is the transformation.

$$\underline{X} \rightarrow \underline{Y} = \Omega' (\underline{X} - \underline{\mu}) \rightarrow(1)$$

where ,  $\Omega$  is orthogonal matrix , such that

$$\Omega' \Sigma \Omega = \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p), \quad \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0$$

The strict positivity of the eigen values  $\lambda_i$  is guaranteed if,  $\Sigma$  is positive definite. The  $i^{\text{th}}$  principle component of  $\underline{X}$  may be defined as the  $i^{\text{th}}$  element of the vector  $\underline{Y}$ , namely as  $Y_i = \omega_i' (\underline{x} - \underline{\mu})$ , where  $\omega_i$  is the  $i^{\text{th}}$  column of  $\Omega$  and may be called the  $i^{\text{th}}$  vector of principle components leadings.

**An Alternative Definitions:-**

If  $\underline{X}$  is a  $p \times 1$  random vector with mean  $\underline{\mu}$  and variance-covariance matrix  $\Sigma$ , then the first principle component of  $\underline{X}$  is defined by the Standardized Linear Combination (SLC) of  $\underline{X}$ .

$$Y_1 = \beta_1' \underline{X} \quad (\text{where, } \beta_1' \beta_1 = 1)$$

such that  $V(Y_1)$  is larger than the variance of any other SLC  $\alpha' \underline{X}$ .

i.e.,  $V(Y_1) \geq V(\alpha' \underline{X})$  for an arbitrary ' $\alpha$ ' ( $\alpha' \alpha = 1$ ). In otherwords  $Y_1$  has the largest variance among all SLC's of  $\underline{X}$ .

**Principle Component Definition :-**

In general the  $k^{\text{th}}$  principle component of  $\underline{X}$  is defined by the SLC of  $\underline{X}$ ,  $Y_k = \beta_k' \underline{X}$  (where,  $\beta_k' \beta_k = 1$ ), which is uncorrelated with first  $k-1$  principle component and  $V(Y_k) \geq V(Y_i)$ , for  $i=1,2,\dots,k-1$ .

**Derivation Of The Principle Components:-**

Suppose  $\underline{X}$  is  $p \times 1$  random vector with mean vector  $\underline{\mu}$  and covariance-matrix  $\Sigma$  i.e.,  $\underline{X} \sim (\underline{\mu}, \Sigma)$ , then by definition, the first principle component is the SLC of  $\underline{X}$  which has largest variance among all SLC's of  $\underline{X}$ . Thus we should seek a LC of  $\underline{X}$  viz.,

$$Y = \omega' \underline{X} \quad \rightarrow (1)$$

with largest variance,

$$V(Y) = \omega' V(X) \omega = \omega' \Sigma \omega \quad \rightarrow (2)$$

such that  $\omega' \omega = 1$  .

Thus , we have to maximize (2) subject to the condition

$$\omega' \omega = 1 \quad \rightarrow (3)$$

which is equivalent to maximizing the function ,

$$\phi(\omega, \lambda) = \omega' \Sigma \omega - \lambda(\omega' \omega - 1) \quad \rightarrow (4)$$

w.r.t  $\omega$  and  $\lambda$  , where ‘ $\lambda$ ’ is a Lagrange multiplier. This implies to solve

$$\text{the equations , } \frac{\partial \phi}{\partial \omega} = 0 \quad \Rightarrow \quad \Sigma \omega = \lambda \omega$$

$$\text{i.e., } (\Sigma - \lambda I) \omega = 0 \quad \rightarrow (5)$$

$$\frac{\partial \phi}{\partial \lambda} = 0 \quad \Rightarrow \quad \omega' \omega = 1 \quad \rightarrow (6)$$

$$\text{Using (5) \& (6) , from (2) , we get , } V(Y) = \lambda \omega' \omega = \lambda \quad \rightarrow (7)$$

From (5), to have a non-zero solution for  $\omega$  , we must have ,

$$|\Sigma - \lambda I| = 0 \quad \rightarrow (8)$$

We know that (8) is a characteristic equation , and ‘ $\lambda$ ’ is a latent root and

From (5),  $\omega$  is the corresponding latent vector of the equation . But, we

know that , solving (8) for ‘ $\lambda$ ’ gives p-latent roots (positive) ,

$$\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_p \geq 0 \quad \rightarrow (9)$$

With the corresponding latent vectors ,  $\omega_1, \omega_2, \dots, \omega_p$  respectively ,

$$\text{i.e., we have from (5) , } \Sigma \omega_i = \lambda_i \omega_i, \quad i = 1, 2, \dots, p \quad \rightarrow (9.a)$$

Since  $\lambda_1$  is the largest latent root among all latent roots and  $\omega_1$  is the corresponding latent vector .

From (1) and (7) ,  $Y_1 = \omega_1' X$ , is the first principle component with variance,

$$V(Y_1) = \lambda_1 .$$

$$\text{Let us denote the first principle component by } Y_1 . \text{ Now , } Y_1 = \omega_1' X \quad \rightarrow (10)$$

$$V(Y_1) = \lambda_1 \quad \rightarrow (11)$$

$$\text{Now , let us show that for } 2 \leq k \leq p , \quad Y_k = \omega_k' X \quad \rightarrow (12)$$

$$\text{is the } K^{\text{th}} \text{ principle component with variance , } V(Y_k) = \lambda_k \quad \rightarrow (13)$$

By definition ,  $Y_k$  should uncorrelated with  $Y_1, Y_2, \dots, Y_{k-1}$  , which can be



easily verified as follows (for  $J = 1, 2, \dots, k-1$ ).

$$\begin{aligned}
 \text{Cov}(Y_k, Y_i) &= \text{Cov}(\omega_k' \mathbf{X}, \omega_i' \mathbf{X}) \\
 &= \omega_k' \text{Cov}(\mathbf{X}, \mathbf{X}') \omega_i \\
 &= \omega_k' \Sigma \omega_i \\
 &= \lambda_i \omega_k' \omega_i && \text{[From(5)]} \\
 &= 0 \quad (\because \omega_k \text{ \& } \omega_j \text{ are orthogonal vectors})
 \end{aligned}$$

Also by definition, the  $k^{\text{th}}$  PC  $Y_k$  has largest variance than  $Y_{k+1}, \dots, Y_p$  which can also be verifying from (13).

$$\begin{aligned}
 \lambda_k &\geq \lambda_{k+1} \geq \dots \geq \lambda_p \geq 0 \\
 \Rightarrow V(Y_k) &\geq V(Y_{k+1}) \geq \dots \geq V(Y_p) \geq 0.
 \end{aligned}$$

Hence the proof.

**Remark :-**

The above result may be asked as no standard linear combination (SLC) of  $\underline{X}$  has a variance larger than  $\lambda_1$ , the variance of first principle combination.

From the above result, we may say that construction (derivation) of principle components of a given random vector  $\underline{X}$  is equivalent to the problem of the construction (derivation) of the latent roots and latent vectors of the variance-covariance matrix  $\Sigma$  of  $\underline{X}$  in case of known  $\Sigma$ .

**Note:-**

- (1) If  $\Sigma$  is not known, we may construct the principle component of the random vector  $\underline{X}$  based on the sample variance – covariance matrix or sample correlation matrix.
- (2) If the population correlation matrix ‘ $\rho$ ’ is given, we may use it in place of  $\Sigma$  to construct the principle components.
- (3) The principle components of the random vector  $\underline{X}$  derived from population (sample) correlation matrix are different from the principle components derived from population (sample) covariance matrix.

**Properties Of Principle Components :-**

(1) Sum of the variances of all p.c's equal to the trace of  $\Sigma$ .

**Proof:-** Let  $Y_1, Y_2, \dots, Y_p$  are the p.c's obtained from random variable  $\underline{X}$ .

Let us denote,  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)_{p \times p}$

$$\Omega = (\omega_1, \omega_2, \dots, \omega_p)_{p \times p}$$

where,  $\lambda_i$ 's are the latent roots and  $\omega_i$ 's are the latent vectors of the covariance matrix  $\Sigma$  of the random variable  $\underline{X}$ . Then, we have,

$$\Omega' \Sigma \Omega = \Lambda$$

$$\Rightarrow \text{Tr}(\Lambda) = \text{Tr}(\Omega' \Sigma \Omega)$$

$$\Rightarrow \lambda_1 + \lambda_2 + \dots + \lambda_p = \text{Tr}(\Sigma \Omega \Omega')$$

$$\Rightarrow V(Y_1) + V(Y_2) + \dots + V(Y_p) = \text{Tr}(\Sigma I)$$

$$= \text{Tr}(\Sigma)$$

$$= V(X_1) + V(X_2) + \dots + V(X_p)$$

$$(\lambda_i = V(Y_i)).$$

(2) Product of the variance of p.c's is equal to the determinant of  $\Sigma$  i.e.,  $|\Sigma|$  (or generalized variance).

**Proof:-** We have  $\Omega' \Sigma \Omega = \Lambda$

$$\Rightarrow |\Lambda| = |\Omega' \Sigma \Omega|$$

$$\Rightarrow \lambda_1 \lambda_2 \dots \lambda_p = |\Sigma \Omega \Omega'|$$

$$\Rightarrow V(Y_1) V(Y_2) \dots V(Y_p) = |\Sigma I|$$

$$= |\Sigma|$$

Hence the proof.

(3) The sum of the first K eigen values divided by the sum of all eigen values,

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\lambda_1 + \lambda_2 + \dots + \lambda_p} = \frac{\lambda_1 + \lambda_2 + \dots + \lambda_k}{\text{Tr}(\Sigma)}$$

represents the 'Proportion of total variation' explained by the first K principle components.

(4) The principle components of a random vector are not scale invariant. It is one disadvantage of principle component analysis.

**Theorem :-**

An orthogonal transformation  $\underline{Y} = C\underline{X}$  of a random vector  $\underline{X}$  leaves invariant the generalized variance and the sum of the variance of the components .

**Proof**:- We have given  $\underline{X}$  is the original random vector and  $\underline{Y}$  is the transformed random variable using the orthogonal matrix C . Now, we have to show that ,  $|\text{cov}(\underline{X},\underline{X}')| = |\text{cov}(\underline{Y},\underline{Y}')|$  , where  $||$

is determinant and  $\sum_{i=1}^p V(X_i) = \sum_{i=1}^p V(Y_i)$

since , 'C' is orthogonal we have,  $C'C = CC' = I \rightarrow (1)$

Now ,  $\text{cov} (Y,Y') = \text{cov} (C\underline{X},(C\underline{X})')$   
 $= C \text{cov} (\underline{X},\underline{X}')C'$   
 $= C\Sigma C'$

$\Rightarrow |\text{cov}(\underline{Y},\underline{Y}')| = |C\Sigma C'|$   
 $= |C||\Sigma||C'|$   
 $= |\Sigma||CC'|$   
 $= |\Sigma||I| \quad (\because \text{from (1)})$   
 $= |\Sigma|$   
 $= |\text{cov}(\underline{X},\underline{X}')|$

$\Rightarrow$  generalized variance of  $\underline{Y} =$  generalized variance of  $\underline{X}$  .

We have ,  $\sum_{i=1}^p V(X_i) = \text{Tr} (\Sigma)$   
 $= \text{Tr} (\Sigma I)$   
 $= \text{Tr} (\Sigma CC')$  ( $\because$  from (1))  
 $= \text{Tr} (C\Sigma C')$   
 $= \sum_{i=1}^p V(Y_i) \quad (\because C\Sigma C' \text{ is covariance matrix of } \underline{Y}) \Rightarrow$

Sum of the variances of original variables (total population variance)  
 $=$  Sum of variances of principle components .  
 $= \lambda_1 + \lambda_2 + \dots + \lambda_p$  .

**Note**:-

The above theorem may be stated as follows . The generalized variance of the vector of principle components is the generalized variance of the original vector and the sum of the variances of the principle components is the sum of the variances of the original variates .

**Results :-**

If  $\underline{X}$  is a random vector with covariance matrix  $\Sigma$  and  $Y_i = \omega_i' \underline{X}$  is the  $i^{\text{th}}$  principle component of the random vector  $\underline{X} = (X_1, X_2, \dots, X_p)'$ , then the correlation coefficient between  $i^{\text{th}}$  principle component and  $J^{\text{th}}$  original variable (that correlation coefficient between  $Y_i$  and  $X_J$ ) is given by

$$\rho_{Y_i, X_J} = \frac{\sqrt{\lambda_i}}{\sqrt{\sigma_{JJ}}} \omega_{ij} \quad i, J = 1, 2, \dots, p, \quad \text{where } \sigma_{JJ} = V(X_J)$$

$\lambda_i$  is  $i^{\text{th}}$  largest root of  $\Sigma$  and  $\omega_{ij}$  is  $J^{\text{th}}$  component of  $\omega_i$ , when  $\omega_i$  is the latent vector of  $\Sigma$  corresponding to  $\lambda_i$ .

**Proof :-** Denote  $l_j = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ J \\ 0 \\ 0 \end{pmatrix} \rightarrow J^{\text{th}} \text{ position}$

Now,  $X_J = l_j' \underline{X}$ . Also we have given,  $Y_i = \omega_i' \underline{X}$  → (1)

$$\begin{aligned} \rho_{Y_i, X_J} &= \frac{\text{cov}(\omega_i' \underline{X}, l_j' \underline{X})}{\sqrt{\lambda_i} \sigma_{JJ}} = \frac{\text{cov}(Y_i, X_J)}{\sqrt{V(Y_i)V(X_J)}} \\ &= \frac{\text{cov}(\omega_i' \underline{X}, l_j' \underline{X})}{\sqrt{\lambda_i} \sigma_{JJ}} \quad (\because V(Y_i) = \lambda_i, \text{ the } i^{\text{th}} \text{ larger latent root} \end{aligned}$$

of  $\Sigma$ , using (1) &  $\sigma_{JJ}$  is  $J^{\text{th}}$  diagonal element of  $\Sigma$ ).

$$\begin{aligned} &= \frac{\omega_i' \text{cov}(\underline{X}, \underline{X}') l_j}{\sqrt{\lambda_i} \sqrt{\sigma_{JJ}}} \\ &= \frac{\omega_i' \Sigma l_j}{\sqrt{\lambda_i} \sqrt{\sigma_{JJ}}} \quad \rightarrow (2) \end{aligned}$$

( $\because \Sigma$  is covariance matrix of  $\underline{X}$ )

Since  $\omega_i$  is the latent vector of  $\Sigma$  corresponding to latent root  $\lambda_i$ , we have

$$\Sigma \omega_i = \lambda_i \omega_i$$

$$\Rightarrow \omega_i' \Sigma = \lambda_i \omega_i' \quad (\text{Taking transpose \& } \Sigma = \Sigma') \quad \rightarrow (3)$$

Using (3) in (2), we get ,

$$\begin{aligned} \rho_{Y_i, X_j} &= \frac{\lambda_i \omega_i' l_j}{\sqrt{\lambda_i} \sqrt{\sigma_{jj}}} \\ &= \frac{\sqrt{\lambda_i}}{\sqrt{\sigma_{jj}}} (\omega_{i1} 0 + \omega_{i2} 0 + \dots + \omega_{ij} 1 + \omega_{i,j+1} 0 + \dots + 0) \\ &= \frac{\sqrt{\lambda_i}}{\sqrt{\sigma_{jj}}} \omega_{ij} \quad \text{for } i, j = 1, 2, \dots, p \end{aligned}$$

Hence the proof .

**Principle Components Obtained From Standardized Variables :-**

Suppose  $\underline{X}$  is a random vector with mean,  $\underline{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix}$

and covariance matrix,  $\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{pmatrix}$

Now, the standardized random vector of  $\underline{X}$  is given by

$$\begin{aligned} \underline{Z} &= \begin{pmatrix} \frac{X_1 - \mu_1}{\sqrt{\sigma_{11}}} \\ \frac{X_2 - \mu_2}{\sqrt{\sigma_{22}}} \\ \vdots \\ \frac{X_p - \mu_p}{\sqrt{\sigma_{pp}}} \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{\sigma_{11}}} & 0 & \dots & 0 \\ 0 & \frac{1}{\sqrt{\sigma_{22}}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{1}{\sqrt{\sigma_{pp}}} \end{pmatrix} \begin{pmatrix} X_1 - \mu_1 \\ X_2 - \mu_2 \\ \vdots \\ X_p - \mu_p \end{pmatrix} \\ &= D (\underline{X} - \underline{\mu}) \quad \rightarrow (1) \end{aligned}$$

Where,  $D = \text{diag} \left( \frac{1}{\sqrt{\sigma_{11}}}, \frac{1}{\sqrt{\sigma_{22}}}, \dots, \frac{1}{\sqrt{\sigma_{pp}}} \right)$

$$\begin{aligned} \text{Now, } \text{cov}(\underline{Z}, \underline{Z}') &= D \text{cov}(\underline{X} - \underline{\mu}, (\underline{X} - \underline{\mu})') D' \\ &= D \Sigma D \quad (\text{since } D = D' \text{ \& } \text{cov}(\underline{X}, \underline{X}') = \Sigma) \quad \rightarrow (2) \end{aligned}$$

But we verify that,

$$D\Sigma D = \begin{pmatrix} \frac{1}{\sqrt{\sigma_{11}}} & 0 & \dots & \dots & 0 \\ 0 & \frac{1}{\sqrt{\sigma_{22}}} & \dots & \dots & 0 \\ \vdots & \vdots & & & \\ \vdots & \vdots & & & \\ 0 & 0 & \dots & \dots & \frac{1}{\sqrt{\sigma_{pp}}} \end{pmatrix} \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \dots & \sigma_{2p} \\ \vdots & \vdots & & & \\ \vdots & \vdots & & & \\ \sigma_{p1} & \sigma_{p2} & \dots & \dots & \sigma_{pp} \end{pmatrix}$$

$$\begin{pmatrix} \frac{1}{\sqrt{\sigma_{11}}} & 0 & \dots & \dots & 0 \\ 0 & \frac{1}{\sqrt{\sigma_{22}}} & \dots & \dots & 0 \\ \vdots & \vdots & & & \\ \vdots & \vdots & & & \\ 0 & 0 & \dots & \dots & \frac{1}{\sqrt{\sigma_{pp}}} \end{pmatrix}$$

$$= \begin{pmatrix} 1 & \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}} & \dots & \dots & \frac{\sigma_{1p}}{\sqrt{\sigma_{11}\sigma_{pp}}} \\ \frac{\sigma_{21}}{\sqrt{\sigma_{11}\sigma_{22}}} & 1 & \dots & \dots & \frac{\sigma_{2p}}{\sqrt{\sigma_{22}\sigma_{pp}}} \\ \vdots & \vdots & & & \\ \vdots & \vdots & & & \\ \frac{\sigma_{p1}}{\sqrt{\sigma_{11}\sigma_{pp}}} & \dots & \dots & \dots & 1 \end{pmatrix}$$

=  $\rho$  , the correlation matrix  $\underline{X}$  . → (3)

From (2) & (3) ,  $\text{cov}(\underline{Z}, \underline{Z}') = \rho$

Thus , the covariance matrix of standardized random vector  $\underline{Z}$  is nothing but the correlation matrix of the original random vector  $\underline{X}$  .

Now , the principle component of  $\underline{Z}$  may be obtained from the eigen vectors of the correlation matrix  $\rho$  of  $\underline{X}$  . The  $i^{\text{th}}$  principle component of the standardized variables  $Z_1, Z_2, \dots, Z_p$  with  $\text{cov}(\underline{Z}, \underline{Z}') = \rho = \text{cov}(\underline{X}, \underline{X}')$

is given by  $Y_i = \omega_i' \underline{z} = \omega_i' D (\underline{X} - \underline{\mu})$  ,  $i=1, 2, \dots, p$  (  $\because$  from (1))

$$\begin{aligned} \text{Moreover, } \sum_{i=1}^p \text{var}(Y_i) &= \sum_{i=1}^p \lambda_i = \text{Tr}(\rho) \\ &= \sum_{i=1}^p V(Z_i) \\ &= p \quad (\because V(Z_i) = 1) \end{aligned}$$

and  $\text{cov}(Y_i, Z_j) = \omega_{ij} \sqrt{\lambda_i}$  where,  $i, j = 1, 2, \dots, p$ .

Here,  $(\lambda_i, \omega_i)$ ,  $i = 1, 2, \dots, p$  are eigen value – eigen vector pairs of ‘ $\rho$ ’ with  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$ .

**Note :-**

The principle components discussed so far are called as population principle components . The computation of principle components is possible only when  $\Sigma$  or  $\rho$  is known , which is a rare practical situation . Therefore , we have to estimate the population principle components from the given sample . The estimation of the population principle components is equivalent to finding the eigen values and the eigen vectors of the sample covariance matrix  $S$  obtained from the sample . Therefore estimated population principle components are also called as Sample Principle Components. Sometimes , we may obtain the sample principle components from the sample correlation matrix  $R$ . But, the sample principle components obtained from  $S$  are different from those obtained from  $R$  .

**Result :-**

If  $S = (s_{ij})_{p \times p}$  is the sample covariance matrix of the given sample  $x_1, x_2, \dots, x_n$  from a multivariate population and  $(\hat{\lambda}_1, \hat{\omega}_1)$   $(\hat{\lambda}_2, \hat{\omega}_2)$  .....  $(\hat{\lambda}_p, \hat{\omega}_p)$  are the eigen value – eigen vectors pairs of  $S$  , then the  $i^{\text{th}}$  sample principle component is given by

Consider a pattern (random vector)

$$\mathbf{X} = \begin{pmatrix} X_1 \\ X_2 \\ \cdot \\ \cdot \\ X_p \end{pmatrix} \in R^p \tag{1}$$

with Covariance matrix  $\mathbf{\Sigma}$  given by

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \cdot & \cdot & \cdot & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \cdot & \cdot & \cdot & \sigma_{2p} \\ \cdot & \cdot & & & & \cdot \\ \cdot & \cdot & & & & \cdot \\ \cdot & \cdot & & & & \cdot \\ \sigma_{p1} & \sigma_{p2} & \cdot & \cdot & \cdot & \sigma_{pp} \end{bmatrix} \quad (2)$$

where  $\sigma_{ij} = \text{Cov}(X_i, X_j)$



## **FACTOR ANALYSIS**

### **INTRODUCTION:-**

Factor analysis is a mathematical model which attempts to explain the correlation between a large set of variables in terms of a small number of underlying unobservable factors. In other words, the essential purpose of factor analysis is to describe, if possible, the covariance relationships among many variables in terms of a few underlying but unobservable, random quantities called factors. Basically, the factor model is motivated by the following argument. Suppose variables can be grouped by their correlations. That is all variables within a particular group are highly correlated among themselves but have relatively small correlations with variables in a different group. It is conceivable that each group of variables represents a single underlying construct, or factor, that is responsible for the observed correlations. Factor analysis was originally developed by psychologists interested in psychometric measurement.

Arguments over the psychological interpretations of several early studies and the lack of powerful computing facilities impelled its development as a statistical method. The advent of high speed computers has generated a renewed interest in the theoretical and computational aspects of factor analysis. Most of the original techniques have been abandoned and early controversies resolved in the wake of recent developments. It is still true that each application of the technique must be examined on its own merits to determine its success.

Factor analysis can be considered as an extension of principal component analysis. Both can be viewed as attempts to approximate the covariance matrix  $\Sigma$ . However, the approximation based

on the factor analysis model is more elaborate . The primary question is factor analysis is whether the data are consistent with a prescribed structure.

In order to get a feel for the subject we first describe a simple example .

Example 1 (Spearman,1904): In children examinations performance in classics ( $x_1$ ),French ( $x_2$ ) and English ( $x_3$ ) .It is found that the correlation matrix is given by

$$\begin{pmatrix} 1 & 0.83 & 0.78 \\ & 1 & 0.67 \\ & & 1 \end{pmatrix}$$

Although this matrix has full rank its dimantionality can be effectively reduced from  $p=3$  to  $p=1$ by expressing the three variables as follows

$$\left. \begin{aligned} x_1 &= \lambda_1 f + u_1 \\ x_2 &= \lambda_2 f + u_2 \\ x_3 &= \lambda_3 f + u_3 \end{aligned} \right\} \quad \text{--- (1)}$$

In these equations  $f$  is an underlying ‘common factor’ and  $\lambda_1, \lambda_2$  and  $\lambda_3$  are known as factor loadings. The terms  $u_1, u_2$  and  $u_3$  represent random disturbance terms. The common factor may be interpreted as ‘general ability’ (or ‘intelligence’) and  $u_i$  will have small variance is  $x_i$  is closely related to general ability. The variation is  $u_i$  consist of two parts which we shall not try to disentangle in practice. First, this variance represent the extent to which an individuals ability at classics , say, differs from his general ability and second it represents the fact that the examination is only an approximate measure of his ability in the subject.

The model defined in (1) can be generalized to include  $k > 1$  common factors.

### **ORTHOGONAL FACTOR MODEL:**

The observable random vector  $\underline{x}$  with  $p$  component has mean  $\underline{\mu}$  and covariance matrix  $\Sigma$  . The factor model postulates that  $\underline{x}$  is linearly dependent upon a few unobservable random variables  $F_1, \dots, F_k$  called common factors and  $p$  additional sources of variations  $u_1, u_2, \dots, u_p$  called random disturbances or error or specific factors . In particular, the factor analysis model is

$$\left. \begin{aligned} X_1 &= \mu_1 + \lambda_{11}F_1 + \lambda_{12}F_2 + \cdots + \lambda_{1k}F_k + u_1 \\ X_2 &= \mu_2 + \lambda_{21}F_1 + \lambda_{22}F_2 + \cdots + \lambda_{2k}F_k + u_2 \\ &\vdots \\ X_p &= \mu_p + \lambda_{p1}F_1 + \lambda_{p2}F_2 + \cdots + \lambda_{pk}F_k + u_p \end{aligned} \right\} \dots (1)$$

(or) in matrix notation.

$$\underline{\mathbf{X}}_{(p \times 1)} = \underline{\boldsymbol{\mu}}_{(p \times 1)} + \underline{\boldsymbol{\Lambda}}_{(p \times k)} \underline{\mathbf{F}}_{(k \times 1)} + \underline{\mathbf{u}}_{(p \times 1)} \quad (2)$$

$$\text{Where } \underline{\mathbf{X}} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix}, \underline{\boldsymbol{\mu}} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{pmatrix}, \underline{\mathbf{F}} = \begin{pmatrix} F_1 \\ F_2 \\ \vdots \\ F_k \end{pmatrix}, \underline{\mathbf{u}} = \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_p \end{pmatrix}, \underline{\boldsymbol{\Lambda}} = \begin{pmatrix} \lambda_{11} & \lambda_{12} & \cdots & \lambda_{1k} \\ \lambda_{21} & \lambda_{22} & \cdots & \lambda_{2k} \\ \vdots & \vdots & & \vdots \\ \lambda_{p1} & \lambda_{p2} & \cdots & \lambda_{pk} \end{pmatrix}$$

The matrix  $\Lambda$  is called the matrix of factor loadings, where  $\lambda_{ij}$  is the loading of  $i^{\text{th}}$  variable ( $X_i$ ) on  $j^{\text{th}}$  factor ( $F_j$ ). Note that the  $i^{\text{th}}$  specific factor  $u_i$  is associated only with the  $i^{\text{th}}$  response  $X_i$ .

The  $p$  deviations  $X_1 - \mu_1, X_2 - \mu_2, \dots, X_p - \mu_p$  are expressed in terms of  $k+p$  random variables  $F_1, F_2, \dots, F_k, \mu_1, \mu_2, \dots, \mu_p$  are unobservable.

From (1), it may be noted that each equation looks like a multiple regression equation but for one exception. The common factor in (1)  $F_1, F_2, \dots, F_k$  are unobservable whereas in multiple regression equation the independent variables can be observed. This distinguishes the factor model from the multivariate regression model. With so many unobservable quantities ( $k+p$ ) a direct verification of the factor model (1) from observations on  $X_1, \dots, X_p$  is hopeless. However, with some additional assumptions about the random vectors  $\underline{\mathbf{F}}$  and  $\underline{\mathbf{u}}$ , the model in (2) implies certain covariance relationships, which can be checked.

We assume that

$$\left. \begin{aligned} E(\underline{\mathbf{F}}) &= \underline{\mathbf{0}}, V(\underline{\mathbf{F}}) = E(\underline{\mathbf{F}}\underline{\mathbf{F}}') = \mathbf{I}_{k \times k} \\ E(\underline{\mathbf{u}}) &= \underline{\mathbf{0}}, V(\underline{\mathbf{u}}) = E(\underline{\mathbf{u}}\underline{\mathbf{u}}') \\ &= \boldsymbol{\Psi} = \text{diag}(\psi_1, \psi_2, \dots, \psi_p) \\ \text{and } \text{cov}(\underline{\mathbf{u}}, \underline{\mathbf{F}}) &= E(\underline{\mathbf{u}}\underline{\mathbf{F}}') = \mathbf{0}_{p \times k} \end{aligned} \right\} \text{---- (3)}$$

The model (2) with the assumptions (3) is called the 'Orthogonal Factor model'

The assumption (3) implies the following implicit assumptions.

- All common factors are standardized to have variance 1 and uncorrelated with one another ( $V(\underline{\mathbf{F}}) = \mathbf{I}$ )
- All specific factors (random disturbances) are have zero means and uncorrelated

$$V(\underline{\mathbf{u}}) = \Psi = \text{diag}(\psi_1, \psi_2, \dots, \psi_p)$$

- Common factor and specific factor are uncorrelated ( $\text{cov}(\underline{\mathbf{u}}, \underline{\mathbf{F}}) = 0$ ).

The Orthogonal model with  $k$  common factors

$$\underline{\mathbf{X}}_{(p \times 1)} = \underline{\boldsymbol{\mu}}_{(p \times 1)} + \Lambda_{(p \times k)} \underline{\mathbf{F}}_{(k \times 1)} + \underline{\mathbf{u}}_{(p \times 1)} \quad \text{---- (4)}$$

Where  $\mathbf{X}_i = i^{\text{th}}$  response variable

$$\boldsymbol{\mu}_i = \text{mean of } \mathbf{X}_i$$

$$\lambda_{ij} = \text{loading of } \mathbf{X}_i \text{ on } \mathbf{F}_j$$

$$\mathbf{F}_j = j^{\text{th}} \text{ common factor}$$

$$\mathbf{u}_i = i^{\text{th}} \text{ specific factor.}$$

The unobservable random vectors  $\underline{\mathbf{F}}$  and  $\underline{\mathbf{u}}$  satisfy  $\underline{\mathbf{F}}$  and  $\underline{\mathbf{u}}$  are independent

$$E(\underline{\mathbf{u}}) = 0, V(\underline{\mathbf{u}}) = \Psi = \text{diag}(\psi_1, \psi_2, \dots, \psi_p)$$

The orthogonal factor model implies a covariance structure for  $\underline{\mathbf{X}}$ . From the model in (4), we have

$$\begin{aligned} (\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})' &= (\Lambda \underline{\mathbf{F}} + \underline{\mathbf{u}})(\Lambda \underline{\mathbf{F}} + \underline{\mathbf{u}})' \\ &= \Lambda \underline{\mathbf{F}} \underline{\mathbf{F}}' \Lambda' + \Lambda \underline{\mathbf{F}} \underline{\mathbf{u}}' + \underline{\mathbf{u}} \underline{\mathbf{F}}' \Lambda' + \underline{\mathbf{u}} \underline{\mathbf{u}}' \end{aligned}$$

so that

$$\begin{aligned} \Sigma = V(\underline{\mathbf{X}}) &= E\left((\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})(\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}})'\right) \\ &= \Lambda E(\underline{\mathbf{F}} \underline{\mathbf{F}}') \Lambda' + \Lambda E(\underline{\mathbf{F}} \underline{\mathbf{u}}') + E(\underline{\mathbf{u}} \underline{\mathbf{F}}') \Lambda' + E(\underline{\mathbf{u}} \underline{\mathbf{u}}') \\ &= \Lambda \Lambda' + \Psi \quad \quad \quad (\text{From(3)}) \quad \text{---- (5)} \end{aligned}$$

Also from the model (4), we have

$$\begin{aligned} (\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}}) \underline{\mathbf{F}}' &= (\Lambda \underline{\mathbf{F}} + \underline{\mathbf{u}}) \underline{\mathbf{F}}' = \Lambda \underline{\mathbf{F}} \underline{\mathbf{F}}' + \underline{\mathbf{u}} \underline{\mathbf{F}}' \\ \text{cov}(\underline{\mathbf{X}}, \underline{\mathbf{F}}) &= E\left((\underline{\mathbf{X}} - \underline{\boldsymbol{\mu}}) \underline{\mathbf{F}}'\right) = \Lambda E(\underline{\mathbf{F}} \underline{\mathbf{F}}') + E(\underline{\mathbf{u}} \underline{\mathbf{F}}') = \Lambda \quad (\text{From(3)}) \quad \text{---- (6)} \end{aligned}$$

From the model (1), we have

$$X_i = \mu_i + \sum_{j=1}^k \lambda_{ij} F_j + u_i, \quad i=1, 2, \dots, p$$

Covariance-structure for the orthogonal factor model

$$\left. \begin{aligned} 1. V(\underline{\mathbf{X}}) &= \Lambda \Lambda' + \Psi \\ \text{or } X_i &= \mu_i + \sum_{j=1}^k \lambda_{ij} F_j + u_i \\ V(X_i) &= \sum_{j=1}^k \lambda_{ij}^2 + \psi_i \text{ and } \text{cov}(X_i, X_j) = \sum_{j=1}^k \lambda_{ij} \lambda_{jj} \\ 2. \text{cov}(\underline{\mathbf{X}}, \underline{\mathbf{F}}) &= \Lambda \text{ or } \text{cov}(X_i, F_j) = \lambda_{ij} \end{aligned} \right\} \quad \text{---- (7)}$$

From the above ,thus  $V(X_i)$  can be split into two parts.

First  $h_i^2 = \sum_{j=1}^k \lambda_{ij}^2$  is called the *communality* and represents the variance of  $X_i$  which is shared with the other variables via the common factors.

In particular  $\lambda_{ij}^2 = [\text{cov}(X_i, F_j)]^2$  represents the extent to which  $X_i$  depends on the  $j^{\text{th}}$  common factor . On the other hand  $\psi_i$  is called specific or unique variance and is due to the specific factor  $u_i$  it explains the variability in  $X_i$  not shared with other variables.

Thus from (7)

$$\frac{\sigma_{ii}}{V(X_i)} = \underbrace{\lambda_{i1}^2 + \lambda_{i2}^2 + \dots + \lambda_{ik}^2}_{\text{communality}} + \underbrace{\psi_i}_{\text{specific variance}} \quad \text{---- (8)}$$

$$\frac{\sigma_{ii}}{V(X_i)} = \underbrace{h_i^2}_{\text{communality}} + \underbrace{\psi_i}_{\text{specific variance}} \quad \text{---- (9)}$$

so that the  $i^{\text{th}}$  communality is the sum of squares of the loadings of the  $i^{\text{th}}$  variable on  $k$  common factors.

**NOTE:**The validity of the k-factor model can be expressed in terms of a simple condition on  $\Sigma$  From (5) we have

$$\Sigma = \Lambda \Lambda' + \psi \quad \text{---- (10)}$$

The converse also holds . If  $\Sigma$  can be decomposed into the form (10), then the k-factor model holds For  $\underline{X}$ . However ,  $\underline{F}$  and  $\underline{u}$  are not uniquely determined by  $\underline{X}$ .

Factor analysis is invariant of scaling of variables

$$\text{Suppose } \underline{X} = \underline{\mu} + \Lambda_x \underline{F} + \underline{u} \quad \text{---- (1)}$$

is the factor model.

Now rescaling the variables of  $\underline{X}$  is equivalent to set

$$\underline{Y} = C\underline{X}, \text{ where } C = \text{diag}(c_1, c_2, \dots, c_p)$$

$$= \begin{pmatrix} c_1 & 0 & 0 & \dots & 0 \\ 0 & c_1 & 0 & \dots & 0 \\ 0 & 0 & c_1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & c_1 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ \vdots \\ X_p \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ \vdots \\ Y_p \end{pmatrix} = \begin{pmatrix} c_1 X_1 \\ c_2 X_2 \\ \vdots \\ \vdots \\ c_p X_p \end{pmatrix} \quad \text{---- (2)}$$

Premultiplying (1) with C we get

$$\underline{Y} = C\underline{\mu} + C\Lambda \underline{F} + C\underline{u}$$

$$\text{and } V(\underline{Y}) = C\Lambda_x \Lambda_x' C' + C\psi_x C'$$

$$\text{i.e., } \Sigma_y = \Lambda_y \Lambda_y' + \psi_y \quad \text{-----(3)}$$

when  $\Lambda_y = C\Lambda_x$

$$\begin{aligned}\psi_y &= C\psi_x C \quad (\because C = C') \\ &= \text{diag}(c_1^2\psi_1, c_2^2\psi_2, \dots, c_p^2\psi_p)\end{aligned}$$

From (1)

$$\begin{aligned}V(\underline{\mathbf{X}}) &= \Lambda_x \Lambda_x' + \psi_x \\ \Rightarrow \Sigma_x &= \Lambda_x \Lambda_x' + \psi_x\end{aligned}\quad \text{-----(4)}$$

But we have

$$\begin{aligned}\Sigma_y &= C\Sigma_x C \quad (\because C = C') \\ &= C\Lambda_x \Lambda_x' C' + C\psi_x C \\ &= \Lambda_y \Lambda_y' + \psi_y\end{aligned}$$

Which is nothing but (3).

Thus the factor loading matrix  $\Lambda_y$  for the scaled random vector  $\underline{\mathbf{Y}}$  is obtained by scaling the factor loading matrix  $\Lambda_x$  of the original random vector  $\underline{\mathbf{X}}$ . Similarly the specific variance matrix  $\psi_y$  for the scaled random vector  $\underline{\mathbf{Y}}$  is obtained by premultiplying and postmultiplying the specific variance matrix  $\psi_x$  of the original r.v.  $\underline{\mathbf{X}}$  by  $C$ . In other words, factor analysis (unlike principal component analysis) is unaffected by a rescaling of the variables.

### Non-uniqueness of factor loadings (Rotated Factors)

Let  $T$  is any  $k \times k$  orthogonal matrix. So that,  $TT' = T'T = I$ . Then the factor model

$$\underline{\mathbf{X}} = \underline{\boldsymbol{\mu}} + \Lambda \underline{\mathbf{F}} + \underline{\mathbf{u}} \quad \text{---- (1)}$$

$$\begin{aligned}\underline{\mathbf{X}} &= \underline{\boldsymbol{\mu}} + \Lambda TT' \underline{\mathbf{F}} + \underline{\mathbf{u}} \\ &= \underline{\boldsymbol{\mu}} + (\Lambda T)(T' \underline{\mathbf{F}}) + \underline{\mathbf{u}}\end{aligned}$$

Can be written as

$$= \underline{\boldsymbol{\mu}} + \Lambda^* \underline{\mathbf{F}}^* + \underline{\mathbf{u}} \quad \text{----(2)}$$

$$\text{where, } \Lambda^* = \Lambda T \text{ and } \underline{\mathbf{F}}^* = T' \underline{\mathbf{F}}$$

Since,  $E(\underline{\mathbf{F}}^*) = T'E(\underline{\mathbf{F}}) = \underline{\mathbf{0}}$  and  $V(\underline{\mathbf{F}}^*) = T'V(\underline{\mathbf{F}})T = T'T = I$ .

It is impossible, on the basis of observations on  $\underline{\mathbf{X}}$  to distinguish the loadings  $\Lambda$  from those of  $\Lambda^*$ . That is the factor  $\underline{\mathbf{F}}$  and  $\underline{\mathbf{F}}^* = T' \underline{\mathbf{F}}$  have the same statistical properties and even though the loadings  $\Lambda^*$  are in general different from the loadings  $\Lambda$ , they both generate the same covariance matrix. That is

$$\Sigma = \Lambda\Lambda' + \psi \quad \text{----(3)}$$

$$= \Lambda T T' \Lambda' + \psi$$

$$= \Lambda^* \Lambda^{*'} + \psi \quad \text{----(4)}$$

Thus the variance-covariance matrix  $\Sigma$  can be decomposed as either (3) or (4). And if  $\Lambda$  is the factor loadings, then  $\Lambda^* = \Lambda T$  (for any orthogonal matrix T), is also the factor loadings.

However, the communalities given by the diagonal elements of  $\Lambda\Lambda' = \Lambda^* \Lambda^{*'}$  are unaffected by the choice of T.

This indeterminacy in the definition of factor loadings is usually resolved by rotating (multiplying by an orthogonal matrix). The factor loadings  $\Lambda$  to satisfy an arbitrary constant such as  $\Lambda' \psi^{-1} \Lambda$  is diagonal or  $\Lambda' D^{-1} \Lambda$  is diagonal,  $D = \text{diag}(\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp})$ . Where in either case the diagonal elements are written in decreasing order. Once the loadings and specific variances are obtained, factors are identified and estimated values for the factors themselves (called factor scores) are frequently constructed.

#### METHODS OF ESTIMATION:-

Given observations  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  on p generally correlated variables, factor analysis seeks to the question, 'Does the factor model

$$\mathbf{X} = \boldsymbol{\mu} + \Lambda \mathbf{F} + \mathbf{u} \quad \text{----(1)}$$

With a small number of factors, adequately represent the data?

In essence, we tackle this statistical model building problem by trying to verify the covariance relationship

$$\Sigma = \Lambda\Lambda' + \psi \quad \text{----(2)}$$

The sample covariance matrix S is an estimator of the unknown  $\Sigma$ . If the off-diagonal elements are small or those of the sample correlation matrix R are essentially zero, the variables are not related and factor analysis will not prove useful. In these circumstances, the specific factors play the dominant role, whereas the major aim of the factor analysis is to determine a few important common factors.

If  $\Sigma$  appears to deviate significantly from diagonal matrix then a factor model can be entertained and the initial problem is one of the estimating the factor loadings  $\lambda_{ij}$ 's and specific variances  $\psi_i$ 's. We shall consider two of the most popular methods of parameter estimation.

1. Principal factor method (Analysis).
2. Maximum likelihood method (factor analysis)

The solution from either method can be rotated in order to simplify the interpretation of factors

Principal factor analysis:

We have the factor model (k-factor)

$$\underline{\mathbf{X}}_{(p \times 1)} = \underline{\boldsymbol{\mu}}_{(p \times 1)} + \underline{\boldsymbol{\Lambda}}_{(p \times k)} \underline{\mathbf{F}}_{(k \times 1)} + \underline{\mathbf{u}}_{(p \times 1)} \quad \text{----(1)}$$

Where  $\underline{\mathbf{X}}$  = p-component random vector

$\underline{\boldsymbol{\mu}}$  = mean of  $\underline{\mathbf{X}}$

$\underline{\boldsymbol{\Lambda}}$  = matrix of factor loadings

$\underline{\mathbf{F}}$  = vector of common factors

$\underline{\mathbf{u}}$  = p-component random vector

with covariance matrix of  $\underline{\mathbf{X}}$

$$\underline{\boldsymbol{\Sigma}} = \underline{\boldsymbol{\Lambda}} \underline{\boldsymbol{\Lambda}}' + \underline{\boldsymbol{\Psi}} \quad \text{----(2)}$$

Where  $\underline{\boldsymbol{\Psi}} = V(\underline{\mathbf{u}}) = \text{diag}(\psi_1, \psi_2, \dots, \psi_k)$ .

In practical situation, since  $\underline{\boldsymbol{\Sigma}}$  is not known,  $\underline{\boldsymbol{\Sigma}}$  is replaced by its estimate the sample covariance matrix  $\underline{\mathbf{S}}$  which is obtained from the observations  $X_1, \dots, X_n$ . Since, factor analysis is invariant of the scaling of the variables the correlaton matrix  $\underline{\mathbf{R}}$ , computed from the observations  $X_1, \dots, X_n$  on p-variable random vector  $\underline{\mathbf{X}}$ , may also be used in place of  $\underline{\mathbf{S}}$ .

Let us suppose the data is summerised by the correlation matrix  $\underline{\mathbf{R}}$  so that an estimate of  $\underline{\boldsymbol{\Lambda}}$  and  $\underline{\boldsymbol{\Psi}}$  is rought for the standardised variables.

Now our problem is to obtain the estimates of  $\underline{\boldsymbol{\Lambda}}$  and  $\underline{\boldsymbol{\Psi}}$  from equation (2), replacing the unknown  $\underline{\boldsymbol{\Sigma}}$  with known  $\underline{\mathbf{R}}$  (when the variables standardised  $\underline{\boldsymbol{\Sigma}}$  is equivalent to the population correlation matrix  $\underline{\boldsymbol{\rho}}$ ). Then we have

$$\underline{\mathbf{R}} = \hat{\underline{\boldsymbol{\Lambda}}} \hat{\underline{\boldsymbol{\Lambda}}} + \hat{\underline{\boldsymbol{\Psi}}} \quad \text{----(3)}$$

Comparing the diagonal elements on both sides, we get

$$1 = \hat{h}_i^2 + \hat{\psi}_i \quad \text{for } i=1,2,\dots,p$$

$$\text{where } \hat{h}_i^2 = \sum_{j=1}^k \hat{\lambda}_{ij}^2$$



Is the preliminary estimate of the  $i^{\text{th}}$  communality  $h_i^2$  and may be obtained either of the following two ways:

- 1) The square of the multiple correlation coefficient of the  $i^{\text{th}}$  variable  $X_i$  on the remaining  $p-1$  variables.
- 2) The largest absolute correlation coefficient between  $X_i$  and one of the remaining  $p-1$  variables.

i.e.,  $\max_{j \neq i} |r_{ij}|$

Note that the estimated communality  $h_i^2$  is higher when  $X_i$  is highly correlated with the other as we would expect. Now  $\hat{\Psi} = \text{diag}(\hat{\psi}_i) = \text{diag}(1 - \hat{h}_i^2)$  has to be subtracted from  $\mathbf{R}$  to obtain the matrix

$$\mathbf{R} - \hat{\Psi} = \begin{bmatrix} \hat{h}_1^2 & r_{12} & \cdots & r_{1p} \\ r_{12} & \hat{h}_2^2 & \cdots & r_{2p} \\ \vdots & \vdots & & \vdots \\ r_{1p} & r_{2p} & \cdots & \hat{h}_p^2 \end{bmatrix} \quad \text{--- (4)}$$

Which is called the reduced correlation matrix because the 1's on the diagonal have been replaced by the estimated *communalities*  $\hat{h}_i^2$ .

Suppose  $a_1 \geq a_2 \geq \cdots \geq a_p$  are eigen values of  $\mathbf{R} - \hat{\Psi}$  and  $\omega_1, \omega_2, \dots, \omega_p$  are the corresponding eigen vectors, then we may decompose  $\mathbf{R} - \hat{\Psi}$  as

$$\mathbf{R} - \hat{\Psi} = \sum_{i=1}^p a_i \omega_i \omega_i' \quad \text{--- (5)}$$

Suppose the first  $K$  eigen values  $a_1, a_2, \dots, a_k$  are positive then

$$\mathbf{R} - \hat{\Psi} = \sum_{i=1}^k a_i \omega_i \omega_i' = \hat{\Lambda} \hat{\Lambda}' \quad \text{--- (6)}$$

Where,  $\hat{\Lambda} = \begin{bmatrix} \sqrt{a_1} \omega_1 & \sqrt{a_2} \omega_2 & \cdots & \sqrt{a_k} \omega_k \end{bmatrix}_{p \times k}$

$$= \mathbf{\Omega} \mathbf{A}^{\frac{1}{2}} \quad \text{--- (7)}$$

$\mathbf{\Omega} = (\omega_1 \ \omega_2 \ \cdots \ \omega_k)$  and  $\mathbf{A} = \text{diag}(a_1, a_2, \dots, a_k)$  is the estimate of the factor loading matrix  $\hat{\Lambda}$ . Since,  $\mathbf{\Omega}$  is orthogonal matrix, we may see that

$$\hat{\Lambda}' \hat{\Lambda} = \mathbf{A}^{1/2} \mathbf{\Omega}' \mathbf{\Omega} \mathbf{A}^{1/2} = \mathbf{A}^{1/2} \mathbf{I} \mathbf{A}^{1/2} = \mathbf{A} \quad \text{--- (8)}$$

Finally, the revised estimates of the specific variances are given in terms of  $\hat{\Lambda}$  by

$$\hat{\psi}_i = 1 - \sum_{j=1}^k \hat{\lambda}_{ij}^2, i = 1, 2, \dots, p \quad \text{--- (9)}$$

Where  $\hat{\lambda}_{ij}$  is the  $(i, j)^{th}$  element of the estimated factor loading matrix  $\hat{\Lambda}$  given by (7). Then the principal factor solution is permissible if all the  $\hat{\psi}_i$  are non-negative.

Thus for the k factor model (1) the principal factor estimates of the factors loading matrix  $\Lambda$  is given by (7) and the estimates of communalities  $h_i^2$  are given by the diagonal elements of  $\hat{\Lambda}\hat{\Lambda}'$ .

$$\text{i.e. } h_i^2 = \sum_{j=1}^k \hat{\lambda}_{ij}^2 \quad \text{--- (10)}$$

The estimates of the specific variables  $\hat{\psi}_i$ 's are given by (9)

NOTE:

- The principal factor analysis can be performed iteratively with the communality estimates given by (10)
- becoming the initial estimates for the next stage.
- If we are given the sample covariance matrix  $S$ , it may be converted into  $R$  and then above analysis can be performed.

Example:

Consider the open/closed book data of the following table with correlation matrix.

$$\begin{bmatrix} 1 & 0.553 & 0.547 & 0.410 & 0.389 \\ & 1 & 0.610 & 0.485 & 0.437 \\ & & 1 & 0.711 & 0.665 \\ & & & 1 & 0.607 \\ & & & & 1 \end{bmatrix}$$

If  $k > 2$  then  $S < 0$  and the factor model is not well defined. The principal factor solutions for  $k=1$  and  $k=2$ , where we estimate the  $i^{th}$  communality  $h_i^2$  by  $\max_j |r_{ij}|$ , are given in the table. The

eigen values of the reduced correlation matrix are 2.84, 0.38, 0.08, 0.02 and -0.05, suggesting that the two-factor solution fits the data well.

In the above table principal factor solutions for the open/closed book data with  $k=1$  and  $k=2$  factors.

variable	$\overbrace{\hat{h}_i^2 \quad \lambda(1)}^{k=1}$		$\overbrace{\hat{h}_i^2 \quad \lambda(1) \quad \lambda(2)}^{k=2}$		
	1	0.417	0.646	0.543	0.646
2	0.506	0.711	0.597	0.711	0.303
3	0.746	0.864	0.749	0.864	-0.051
4	0.618	0.786	0.680	0.786	-0.249
5	0.551	0.742	0.627	0.742	-0.276

The first factor represents overall performance and for k=2, the second factor, which is much less important

( $a_2 = 0.38 \ll 2.84 = a_1$ ), represents a contrast across the range  $h_i^2 \ll 1$  for all i, and therefore a fair proportion of the variance of each variable is left unexplained by the common factor.

**PRINCIPAL COMPONENT METHOD (PRINCIPAL COMPONENT SOLUTION OF THE FACTOR MODEL)**

Suppose  $\underline{x}_1, \dots, \underline{x}_n$  are observations on  $\mathbf{p}$  generally correlated variables and the data is summarised either into the sample correlation matrix  $\mathbf{R}$ .

Let the orthogonal factor model with k common factors

$$\underline{\tilde{X}}_{(p \times 1)} = \underline{\tilde{\mu}}_{(p \times 1)} + \Lambda_{(p \times k)} \underline{\tilde{F}}_{(k \times 1)} + \underline{\tilde{u}}_{(p \times 1)} \quad \text{----(1)}$$

Where  $\underline{\tilde{X}}$  = p-component random vector

$\underline{\tilde{\mu}}$  = mean of  $\underline{\tilde{X}}$

$\Lambda$  = matrix of factor loadings

$\underline{\tilde{F}}$  = vector of common factors

$\underline{\tilde{u}}$  = vector of random disturbances

with  $V(\underline{\tilde{X}}) = \Sigma = \Lambda \Lambda' + \Psi \quad \text{----(1.a)}$

$$\boldsymbol{\psi} = \mathbf{V}(\mathbf{u})$$

Now the principal component method is to obtain the estimates of  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\psi}$  using the sample covariance matrix  $\mathbf{S}$  or sample correlated matrix  $\mathbf{R}$ .

Suppose  $a_1 \geq a_2 \geq \dots \geq a_p$  are the latent roots of  $\mathbf{S}$  (or  $\mathbf{R}$ ) and let us consider the first 'k' roots

i.e.  $a_1, a_2, \dots, a_p$

Let  $\boldsymbol{\omega}_1, \boldsymbol{\omega}_2, \dots, \boldsymbol{\omega}_k$  be the corresponding latent vectors. Then the estimated matrix of factor loadings is given by

$$\hat{\boldsymbol{\Lambda}} = \begin{bmatrix} \sqrt{a_1} \boldsymbol{\omega}_1 & \sqrt{a_2} \boldsymbol{\omega}_2 & \dots & \sqrt{a_k} \boldsymbol{\omega}_k \end{bmatrix}_{p \times k} \quad \text{-----(2)}$$

and the estimated specific variances are provided by the diagonal elements of the matrix

$$\mathbf{S} - \hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Lambda}}' (\mathbf{R} - \hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Lambda}}') \quad \text{-----(3)}$$

so that  $\hat{\boldsymbol{\psi}} = \text{diag}(\hat{\psi}_1, \hat{\psi}_2, \dots, \hat{\psi}_k)$  with  $\hat{\psi}_i = s_{ii} - \sum_{j=1}^k \hat{\lambda}_{ij}^2$   $\left( \hat{\psi}_i = r_{ii} - \sum_{j=1}^k \hat{\lambda}_{ij}^2 \right)$  estimates of communalities

are given by the diagonal elements of  $\hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Lambda}}'$

i.e.

$$h_i^2 = \sum_{j=1}^k \hat{\lambda}_{ij}^2$$

**NOTE:**

- Consider the residual matrix  $\mathbf{S} - (\hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Lambda}}' + \hat{\boldsymbol{\psi}})$  resulting from the approximation of  $\mathbf{S}$  by the principal component solution. The diagonal elements are zero and if the other elements are also small, we may subjectively take the 'k' factor model to be appropriate.
- The contribution to the total sample variance =  $\text{Tr}(\mathbf{S})$  from the  $j^{\text{th}}$  common factor is given by

$$\begin{aligned} \sum_{i=1}^p \hat{\lambda}_{ij}^2 &= (\sqrt{a_j} \boldsymbol{\omega}_j)' (\sqrt{a_j} \boldsymbol{\omega}_j) \\ &= a_j \quad (\because \boldsymbol{\omega}_j' \boldsymbol{\omega}_j = 1) \\ &= j^{\text{th}} \text{ latent root of } \mathbf{S} \quad (\text{Where } a_j \text{ is the } j^{\text{th}} \text{ latent root of } \mathbf{S}) \end{aligned}$$

Thus, proportion of total sample variance due to  $j^{\text{th}}$  factor

$$= \begin{cases} \frac{a_j}{\text{Tr}(S)} & \text{for factor analysis of 'S'} \\ \frac{a_j}{p} & \text{for factor analysis of 'R'} \end{cases}$$

- (For worked out examples see page nos: 388-391 of Applied Multivariate Analysis by Johnson&Wichern)

### MAXIMUM LIKELYHOOD FACTOR ANALYSIS:

Suppose  $\underline{x}_1, \dots, \underline{x}_n$  are 'n' observations drawn on  $\underline{X}$  which follows population  $Np(\underline{\mu}, \Sigma)$  and  $\underline{X}$  is having the following k factor model

$$\underline{X}_{(p \times 1)} = \underline{\mu}_{(p \times 1)} + \Lambda_{(p \times k)} \underline{F}_{(k \times 1)} + \underline{u}_{(p \times 1)}$$

$\underline{\mu}$  = mean of  $\underline{X}$

$\Lambda$  = matrix of factor loadings

$\underline{F}$  = vector of common factors

$\underline{u}$  = vector of random disturbances

with the assumptions

$$E(\underline{F})=0=E(\underline{u})$$

$$V(\underline{F})=I, \quad V(\underline{u})=\psi=\text{diag}(\psi_1, \psi_2, \dots, \psi_p)$$

$$\text{cov}(\underline{F}, \underline{u})=0$$

These assumptions implicitly imposes the restriction on  $\Sigma$  as follows

$$\Sigma = \Lambda \Lambda' + \psi \quad \text{-----(2)}$$

Since  $\underline{X} \sim Np(\underline{\mu}, \Sigma)$ , its log -likelihood is given by

$$\log L = \frac{-n}{2} \log |2\pi \Sigma| - \frac{1}{2} \sum_{i=1}^n (\underline{x}_i - \underline{\mu})' \Sigma^{-1} (\underline{x}_i - \underline{\mu})$$

if we with its MLE  $\bar{\underline{X}}$ , then  $\log L$  becomes

$$\begin{aligned} l = \log L &= \frac{-n}{2} \log |2\pi \Sigma| - \frac{1}{2} \sum_{i=1}^n (\underline{x}_i - \underline{\mu})' \Sigma^{-1} (\underline{x}_i - \underline{\mu}) \\ &= - \left( \frac{n}{2} \log |2\pi \Sigma| + \frac{n}{2} \text{Tr}(\Sigma^{-1} S_n) \right) \end{aligned} \quad \text{-----(3)}$$

Where

$$S_n = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})' (\mathbf{x}_i - \boldsymbol{\mu})$$

$\boldsymbol{\Sigma}$  is as given by (2)

Maximizing (3) is equivalent to minimizing the following function w.r.t.  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\psi}$

$$F(\boldsymbol{\Lambda}, \boldsymbol{\psi}) = \log |\boldsymbol{\Sigma}| + \text{Tr}(\boldsymbol{\Sigma}^{-1} S_n) - \log |S_n| - p \quad \text{---- (4)}$$

( $\because |S_n|$  and  $p$  are constants)

Since from (2),  $\boldsymbol{\Lambda}$  is not uniquely determined. We have minimize (4) subject to the following uniqueness condition

$$\boldsymbol{\Lambda}' \boldsymbol{\psi}^{-1} \boldsymbol{\Lambda} = \Delta, \text{ a diagonal matrix} \quad \text{----(5)}$$

The MLEs  $\hat{\boldsymbol{\Lambda}}$  and  $\hat{\boldsymbol{\psi}}$  obtained by minimizing (4) subject to (5) satisfy

$$\left( \hat{\boldsymbol{\psi}}^{-\frac{1}{2}} S_n \hat{\boldsymbol{\psi}}^{-\frac{1}{2}} \right) \left( \hat{\boldsymbol{\psi}}^{-\frac{1}{2}} \hat{\boldsymbol{\Lambda}} \right) = \left( \hat{\boldsymbol{\psi}}^{-\frac{1}{2}} \hat{\boldsymbol{\Lambda}} \right) (1 + \hat{\boldsymbol{\Lambda}}) \quad \text{----(6)}$$

so that the  $j^{\text{th}}$  column of  $\hat{\boldsymbol{\psi}}^{-\frac{1}{2}} \hat{\boldsymbol{\Lambda}}$  is the (non-normalised) eigen vector of  $\hat{\boldsymbol{\psi}}^{-\frac{1}{2}} S_n \hat{\boldsymbol{\psi}}^{-\frac{1}{2}}$  corresponding to eigen value  $1 + \hat{\boldsymbol{\Lambda}}_i$

where  $\hat{\boldsymbol{\Lambda}}_1 \geq \hat{\boldsymbol{\Lambda}}_2 \geq \dots \geq \hat{\boldsymbol{\Lambda}}_k$

clearly, for the above, the MLE of  $\hat{\boldsymbol{\Lambda}}$  can be obtained only for a given  $\hat{\boldsymbol{\psi}}$ , whose initial value can be taken as

$$\hat{\boldsymbol{\psi}}^{(0)} = \text{diag}(\hat{\boldsymbol{\psi}}_1^{(0)}, \hat{\boldsymbol{\psi}}_2^{(0)}, \dots, \hat{\boldsymbol{\psi}}_k^{(0)})$$

where  $\hat{\boldsymbol{\psi}}_i = \left(1 - \frac{1}{2} \frac{k}{p}\right) \left(\frac{1}{s^{ii}}\right)$

where  $s^{ii}$  is the  $i^{\text{th}}$  diagonal element of  $\mathbf{S} = \frac{n S_n}{(n-1)}$

The next modified value of  $\hat{\boldsymbol{\psi}}$  is given by

$$\hat{\boldsymbol{\psi}}^{(1)} = \text{diag}(\hat{\boldsymbol{\psi}}_1^{(1)}, \hat{\boldsymbol{\psi}}_2^{(1)}, \dots, \hat{\boldsymbol{\psi}}_k^{(1)})$$

Where  $\hat{\boldsymbol{\psi}}_i^{(1)}$  is the  $i^{\text{th}}$  diagonal element of the computed matrix  $S_n - \hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Lambda}}'$

Using this  $\hat{\boldsymbol{\psi}}^{(1)}$ , we can obtain the revised value of  $\hat{\boldsymbol{\Lambda}}$  using (6).

This procedure is to be continued until the latest estimates  $\hat{\boldsymbol{\Lambda}}$  and  $\hat{\boldsymbol{\psi}}$  satisfy the relation (5).

**NOTE:**

Ordinarily the observations are standardised and a sample correlation matrix is factor analysed. Of the data is

summerised into a sample correlation matrix  $\mathbf{R}$ , then the above method of maximum likelihood factor analysis may be

carried out replacing  $\mathbf{S}_n$  or  $\mathbf{S}$  by  $\mathbf{R}$  to get the some estimates of  $\mathbf{\Lambda}$  and  $\mathbf{\Psi}$ . This is due to the fact that the MLEs are scale invariant.

An worked out example is given in page no:394 of Applied Multivariate Statistical analysis by Richard A.Jhon and Wichern.

**GOODNES OF FIT TEST:**

One of the main advantages of the maximum likelihood technique is that it provides a test of the hypothesis

$H_0$  : k comman factors are sufficient to against describe the data .

$H_1$ :  $\Sigma$  has no constraints.

The likelihood ratio statistic  $\lambda$  is given by

$$-2 \log \lambda = np(\hat{a} - \log \hat{g} - 1) \tag{1}$$

where  $\hat{a}$  and  $\hat{g}$  are the arithmatic and geomatric means of the eigen values of  $\Sigma^{-1}\mathbf{S}_n$  .

But we have

$$\text{Tr}(\Sigma^{-1}\mathbf{S}) - \log|\hat{\Sigma}^{-1}\mathbf{S}| - p = p(\hat{a} - \log \hat{g} - 1) \tag{2}$$

When  $\hat{\Sigma} = \hat{\mathbf{\Lambda}}\hat{\mathbf{\Lambda}}' + \hat{\mathbf{\Psi}}$

comparing (1) & (2), we have

$$-2 \log \lambda = n(\text{Tr}(\Sigma^{-1}\mathbf{S}) - \log|\hat{\Sigma}^{-1}\mathbf{S}| - p)$$

$$\text{when } \hat{\Sigma} = \hat{\mathbf{\Lambda}}\hat{\mathbf{\Lambda}}' + \hat{\mathbf{\Psi}}$$

Now under  $H_0$  the statistic  $-2 \log \lambda$  has an asymptotic  $\chi_m^2$  distribution ,

where

$$m = \frac{1}{2}(p - k)^2 - \frac{1}{2}(p - k)$$

According to Bartlet, an improved test statistic is given by (3) replacing n by

$$n' = (n - 1 - (2p + 4k + 5) / 6)$$

**NOTE:** See for example 9.7 page no.399 of AMSA by Johnson.

**FACTOR ROTATION:**

We have

$$\begin{aligned} \Sigma &= \Lambda\Lambda' + \psi \\ &= \Lambda TT'\Lambda' + \psi \\ &= (\Lambda T)(\Lambda T)' + \psi \\ &= \Lambda^* (\Lambda^*)' + \psi \end{aligned}$$

where  $\Lambda^* = \Lambda T$

T is an orthogonal matrix

Thus if  $\Lambda$  is a factor loadings matrix which reproduce  $\Sigma$ , then any other factor loadings matrix  $\Lambda^*$  obtained from  $\Lambda$  by an orthogonal transformation (T) have the same ability to reproduce the covariance matrix (or correlation matrix). From matrix algebra, we know that an orthogonal transformation corresponds to a rigid rotation of the coordinate axes. For this reason an orthogonal transformation of the factor loadings and the implied orthogonal transformation of the factor is called "factor rotation".

Let  $\hat{\Lambda}$  be the  $p \times k$  matrix of estimated factor loadings obtained by any method, then

$$\hat{\Lambda}^* = \hat{\Lambda}T \text{ where } TT' = T'T = I \tag{1}$$

is a  $p \times k$  matrix of rotated loadings. Moreover, the estimated covariance (or correlation) matrix remains unchanged,

$$\text{since } \hat{\Lambda}\hat{\Lambda}' + \hat{\psi} = \hat{\Lambda}^* (\hat{\Lambda}^*)' + \hat{\psi} \tag{2}$$

since, the original loadings may not be readily interpretable, it is usual practice to rotate them until a "sample structure" is achieved. Ideally we should like to see a pattern of loadings such that each variable loadings highly on a single factor and has

small-to-moderate loadings on the remaining factors. Of course, it is not always possible to get this simple structure. A convenient analytical choice of rotation is given by the "varimax method" described below: The varimax method of orthogonal rotation was provided by **kaiser**(1958). Its rationale is to provide axes with a few large loadings and as many near zero loadings as possible. This is accomplished by an iterative maximization of a quadratic function of the loadings.



Devote the matrix of rotated loadings as

$$\hat{\Lambda} = \hat{\Lambda}T$$

Now the  $(i, j)^{th}$  element of  $\hat{\Lambda}$  viz;  $\delta_{ij}$  represents the loadings of the  $i^{th}$  variable on the  $J^{th}$  factor .

The function  $\phi$  that the variance criterion maximizes is the sum of the variances of the squared loadings within each

column of the loadings is normalised by its communality, that is

$$\phi = \sum_{i=1}^k \sum_{j=1}^p (d_{ij}^2 - \bar{d}_i)^2 = \sum_{i=1}^k \sum_{j=1}^p d_{ij}^4 - p \sum_{i=1}^k \bar{d}_i^2$$

$$\text{Where } d_{ij} = \frac{\delta_{ij}}{h_i} \text{ and } \bar{d} = \frac{1}{p} \sum_{j=1}^p d_{ij}^2$$

$h_i^2$  is the  $i^{th}$  communality is the  $i^{th}$  diagonal element of  $\hat{\Lambda} \hat{\Lambda}'$

The varimax criterion  $\phi$  is a function of  $T$ , and the iterative algorithm proposed by Kaiser finds the orthogonal matrix

$G$  which maximizes  $\phi$ .

In the case where  $k=2$ , the calculations simplify. For then  $T$  is given by

$$T = \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}$$

and represents a rotation of the coordinate axis clockwise by an angle  $\theta$ . The value of  $\theta$  can be determined by

the relation  $T'T=I$

In the case where  $k>2$ , an iterative solution for the rotation is used.

See example 9.8, 9.9, 9.10, 9.11 in the pages 401-408 of AMVA by Richard Johanson & Wichern.

# DISCRIMINATION AND CLASSIFICATION

## INTRODUCTION:

Discriminant analysis and classification are multivariate techniques concerned with separating distinct sets of objects (or observations) and with allocating new objects (observations) to previously defined groups. Discriminant analysis is rather exploratory in nature. As a separatory procedure, it is often employed on a onetime basis in order to investigate observed differences when causal relationships are not well understood. Classification procedure are less exploratory in the sense that they lead to well defined rules, which can be used for assigning new objects. Classification ordinarily requires more problem structure than discrimination. Thus, the immediate goals of discrimination and classification, respectively, are as follows:

Goal 1: To describe either graphically (in three or fewer dimensions) or algebraically, the differential features of objects (observations) from several known collections (populations). We try to find “discriminants” whose numerical values are such that the collection are separated as much as possible.

Goal 2: To sort objects (observations) into two or more labeled classes. The emphasis is on deriving a rule that can be used to optimally assign a new object to the labeled classes.

We shall follow convention and use the term discrimination to refer to ‘Goal 1’. This terminology was introduced by R.A. Fisher in the first modern treatment of separatory problems. A more descriptive term for this goal, however, is separation. We shall refer to the second goal as classification, or allocation. A function that separates may some times serve as an allocator, and, conversely, an allocator rule may suggest a discriminatory procedure. Thus, in practice Goal 1 & Goal 2 frequently overlap and distinction between separation and allocation is not clear.

The problem of classification arises when an investigator makes a number of measurements on an individual and wishes to classify the individual into one of several categories on the basis of these measurements. The investigator cannot identify the individual with a category directly but must use these measurements. In many cases it can be assumed that these are a finite number of categories or populations from which the individual may have come and each population is characterised by a probability distribution of the measurements. Thus, an individual is considered as a random observation from this population. The question is : Given an individual with certain measurements, from which population did it arise?

In some, instances, the categories are specified before hand in the sense that the probability distributions of the measurements are completely known. In other cases, the form of each distribution may be known, but the parameters of the distribution must be estimated from a sample from that population .In some other cases, the form of the distribution of the populations may not be known.

Let us give an example of a problem of discrimination and classification. Prospective students applying for admission into college are given a battery of tests; the vector of scores is a set of measurements  $\mathbf{x}$ . The prospective students may be a member of one population consisting of these students who will successfully complete college training or, rather, have potentialities for successfully completing training, or he/she may be member of the other population, those who will not complete the course successfully. The problem is to classify a student applying for admission on the basis of these scores on the entrance examination. Before that we have to describe or explore the differential scores between the two categories of the students from the past information. Also, we have to prepare a discriminant function that separates the two categories of students clearly as much as possible. This problem is called discrimination.

#### **SEPERATION AND CLASSIFICATION FOR TWO POPULATION:**

To fix ideas, we list below situations where one may be interested in

(1). Separating or discriminating two classes of objects.

Or (2). Assigning a new object to one of the two classes .

Or both (1)&(2).

It is convenient to label the classes  $\pi_1$  &  $\pi_2$ . The objects are ordinarily separated or classified on the basis of measurements on, for instance, P associated random variables.  $\mathbf{X}' = (X_1, X_2, \dots, X_p)$ . The observed values of  $\mathbf{X}$  differ to some extent from one class to the other (of the values of  $\mathbf{X}$  were not very different for objects in  $\pi_1$  &  $\pi_2$ , there would be no problem; i.e., the would be indistinguishable and new objects could be assigned to either class indiscriminately). We can think of the totality of values from the first class as being the population of  $\mathbf{x}$  values for  $\pi_1$  and those from the second class as the population of  $\mathbf{x}$  values for  $\pi_2$ . These two populations can then be described by probability density functions

$f_1(\mathbf{x})$  &  $f_2(\mathbf{x})$ , and consequently, we can talk of assigning observations to populations (or objects to classes).

The following are some more examples:

- (1). Separation of two species of chickweed based on the measurements sepal and petal lengths, petal left depth, bract length, scarious tip length and pollen diameter.
- (2). Discrimination of successful and unsuccessful college students based on the entrance examination scores, high school grade point average and number of high school activities.
- (3). Classification of purchasers of a new product and laggards (those slow to purchase) based on particulars of education, income, family size and amount of previous brand switching.
- (4). Discriminating male-skulls and female-skulls based on the anthropological measurements like circumference and volume on ancient skulls.
- (5). Separating good and poor credit risks based on the particulars of income, age, member of credit cards and family size.

From the above examples, it is clear that allocation or classification rules are usually developed from learning samples. Measured characteristics of randomly selected objects known to come from each of the two populations are examined for differences. Essentially, the set of possible sample outcomes is divided into two regions  $R_1$  &  $R_2$ , such that if a new observation falls in  $R_1$ , it is allocated to population  $\pi_1$  and if it falls in  $R_2$ , we allocate it to population  $\pi_2$ . Thus one set of observed values favours  $\pi_1$ , the other set of values favours  $\pi_2$ . Here, it may be noted that classification rules cannot usually provide an error-free method of assignment. This is because there may not be a clear distinction between the measured characteristics of the populations; i.e. the groups may overlap. It is then possible, for example, to incorrectly classify a  $\pi_2$  object as belonging to  $\pi_1$  or a  $\pi_1$  object as belonging to  $\pi_2$ .

A good classification procedure should result in a few misclassifications. In other words, the chances or probabilities of misclassification should be small. As we shall see, there are additional features that an “optimal” classification rule should be possessed.

### STANDARDS OF GOOD CLASSIFICATION:

In constructing a procedure of classification, it is desired to minimize the probability of misclassification or more specifically it is desired to minimize on the average the bad effects of misclassification.

Suppose an individual is an observation from either population  $\pi_1$  or population  $\pi_2$ . The classification of an observation depends on the vector of measurements

$$\underline{x} = (x_1, x_2, \dots, x_p)'_{p \times 1}$$

on that individual. We set up a rule that if an individual is characterized by certain sets of values of  $x_1, x_2, \dots, x_p$  it will be classified as from  $\pi_1$ ; if it has other values it is classified as from  $\pi_2$ .

We can think of an observation  $\underline{x}$  as a point in a P-dimensional space. We divide this space into two regions  $R_1$  &  $R_2$  if the observation falls in  $R_1$ , we classify it as coming from  $\pi_1$  and if it falls in  $R_2$  we classify it as coming from  $\pi_2$ .

Usually, the statistician can make two kinds of errors in classification. If the individual is actually from  $\pi_1$  and is misclassified into  $\pi_2$ ; or if it is actually from  $\pi_2$  and is misclassified into  $\pi_1$ . We need to know the relative undesirability of these two kinds of misclassification.

Let  $f_1(\underline{x})$  &  $f_2(\underline{x})$  be the p.d.f.'s associated with the  $p \times 1$  random vector  $\underline{x}$  for populations  $\pi_1$  &  $\pi_2$  respectively. An object, with associated measurements  $\underline{x}$ , must be assigned to either  $\pi_1$  (or)  $\pi_2$ . Let  $\Omega$  be the sample space that is the collection of all possible observations  $\underline{x}$ . Let  $R_1$  be that set of  $\underline{x}$  values for which we classify objects as  $\pi_1$  and  $R_2 = \Omega - R_1$  be the remaining  $\underline{x}$  values for which we classify objects as  $\pi_2$ . Since every object must be assigned to one and only one of the two populations, the sets  $R_1$  &  $R_2$  be mutually exclusive and exhaustive.

### EXPECTED (OR AVERAGE) COST OF MISCLASSIFICATION (ECM):

In order to obtain ECM we consider the following conditional probabilities:  
P (correctly classifying an observation (object) that actually is drawn from  $\pi_1$ )

$$= P(\underline{X} \in R_1 / \pi_1) = \int_{R_1} f_1(\underline{x}) d\underline{x} = P(1/1) \text{ (say)} \quad \rightarrow (1)$$

P(correctly classifying an observation that actually is drawn from  $\pi_2$ )

$$= P(\underline{X} \in R_2 / \pi_2) = \int_{R_2 = \Omega - R_1} f_2(\underline{x}) d\underline{x} = P(2/2) \text{ (say)} \quad \rightarrow (2)$$

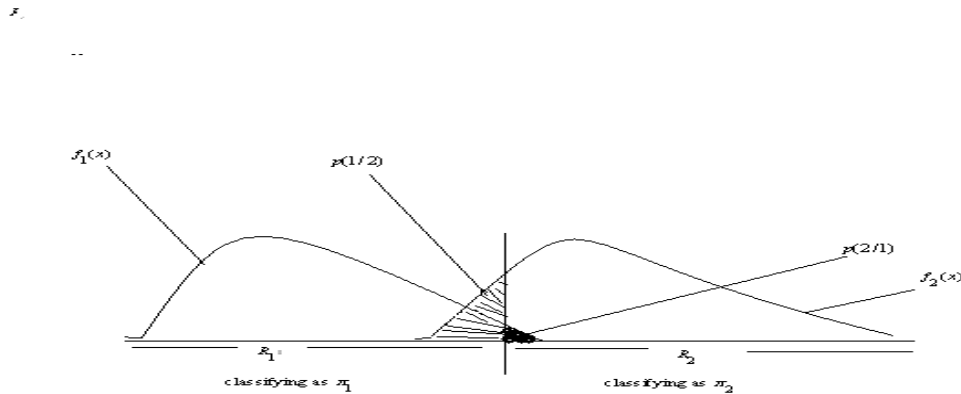
P(misclassifying an observation that is drawn from  $\pi_1$ )

$$= P(\underline{X} \in R_2 / \pi_1) = \int_{R_2} f_1(\underline{x}) d\underline{x} = P(2/1) \text{ (say)} \quad \rightarrow (3)$$

P(misclassifying an observation that is drawn from  $\pi_2$ )

$$= P(\underline{X} \in R_1 / \pi_2) = \int_{R_1} f_2(\underline{x}) d\underline{x} = P(1/2) \text{ (say)} \quad \rightarrow (4)$$

Misclassification probabilities when  $p=1$ :



Let

$P_1$  = prior probability of  $\pi_1$

$$= P(\text{drawing an observation from } \pi_1) = P(\pi_1) \quad \rightarrow (5)$$

and  $P_2$  = prior probability of  $\pi_2$

$$= P(\text{drawing an observation from } \pi_2) = P(\pi_2) \quad \rightarrow (6)$$

Now the overall probabilities of correctly or incorrectly classifying objects can be derived as the product of the prior and conditional classification probabilities. Thus we get

P(correctly classified as  $\pi_1$ ) = P(observations comes from  $\pi_1$  and is correctly Classified as  $\pi_1$ )

$$= P(\underline{X} \in R_1 / \pi_1).P(\pi_1) = P(1/1).P_1 \quad (\text{from (1)\&(5)}) \rightarrow (7)$$

similarly

$$P(\text{correctly classified as } \pi_2) = P(2/2).P_2 \quad (\text{from (2)\&(6)}) \rightarrow (8)$$

$$\begin{aligned} P(\text{misclassified as } \pi_1) &= P(\text{observations comes from } \pi_2 \text{ and is misclassified as } \pi_1) \\ &= P(\underline{X} \in R_1 / \pi_2).P(\pi_2) = P(1/2).P_2 \quad (\text{from (4)\&(6)}) \rightarrow (9) \end{aligned}$$

$$\begin{aligned} P(\text{misclassified as } \pi_2) &= P(\text{observation comes from } \pi_1 \text{ and is misclassified as } \pi_2) \\ &= P(\underline{X} \in R_2 / \pi_1).P(\pi_1) = P(2/1).P_1 \quad (\text{from (3)\&(5)}) \rightarrow (10) \end{aligned}$$

A good classification rule must take into account the misclassification costs. Although the statistician may not know these costs in each case, he will often have at least a rough idea of them. The costs of misclassification can be defined by a cost matrix C:

True population | Classified as

	$\pi_1$	$\pi_2$
$\pi_1$	0	C(2/1)
$\pi_2$	C(1/2)	0

$\rightarrow(11)$

The costs are

- (1). Zero for correct classification .
- (2). C(1/2) is cost involved when an observation drawn from  $\pi_2$  is incorrectly classified into  $\pi_1$  .
- (3). C(2/1) is cost involved when an observation actually drawn from  $\pi_1$  is incorrectly classified as  $\pi_2$  .

Clearly, a good classification procedure is one which minimize in some sense or the cost of misclassification. Now , the expected cost of misclassification(ECM) is obtained by multiplying the off-diagonal entries in (11) by their probabilities of occurrence. Consequently a reasonable classification rule should has an ECM as small as possible. From the above the ECM may be defined as follows :

$$\begin{aligned} \text{ECM} &= C(1/2) .P(\text{misclassification into } \pi_1) + C(2/1) .P(\text{misclassification into } \pi_2) \\ &= C(1/2) .P(1/2).P_2 + C(2/1) . P(2/1) . P_1 \quad \rightarrow (12) \end{aligned}$$

**DEFINITION:**

Expected (or average) cost of misclassification (ECM) is the sum of the products of costs of each misclassification multiplied by the probability of its occurrence. Its formula is given by (12).

**THEOREM (Optional ECM regions or Bayes regions):**

The regions  $R_1$  &  $R_2$  that minimize ECM are defined by the values of  $\underline{x}$  for which the following inequalities hold.

$$R_1 : \frac{f_1(\underline{x})}{f_2(\underline{x})} \geq \left( \frac{C(1/2)}{C(2/1)} \right) / \left( \frac{P_1}{P_2} \right) \rightarrow (1)$$

(density ratio)  $\geq$  (cost ratio)/(prior probability ratio)

$$R_2 : \frac{f_1(\underline{x})}{f_2(\underline{x})} < \left( \frac{C(1/2)}{C(2/1)} \right) / \left( \frac{P_1}{P_2} \right) \rightarrow (2)$$

**PROOF:**

We have the expected cost of misclassification (ECM) as

$$ECM = C(1/2) \cdot P(1/2) + C(2/1) \cdot P(2/1) \rightarrow (3)$$

where  $C(1/2) = \dots\dots\dots$

$C(2/1) = \dots\dots\dots$

$P_1 = \dots\dots\dots$

$P_2 = \dots\dots\dots$

$P(1/2) = \dots\dots\dots$

$P(2/1) = \dots\dots\dots$

But , we have

$$P(1/2) = \int_{R_1} f_2(\underline{x}) d\underline{x}$$
$$P(2/1) = \int_{R_2} f_1(\underline{x}) d\underline{x} \rightarrow (4)$$

using (4) in (3) we get

$$ECM = C(1/2) \int_{R_1} f_2(\underline{x}) d\underline{x} + C(2/1) \int_{R_2} f_1(\underline{x}) d\underline{x} \rightarrow (5)$$

Noting that  $\Omega = R_1 \cup R_2$  so that the total probability

$$1 = \int_{\Omega} f_1(\underline{x}) d\underline{x} = \int_{R_1} f_1(\underline{x}) d\underline{x} + \int_{R_2} f_1(\underline{x}) d\underline{x} \quad (\because R_1 \text{ \& } R_2 \text{ are disjoint}) \rightarrow (6)$$



using (6) in (5), we get

$$\begin{aligned} \text{ECM} &= C(1/2) P_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x} + C(2/1) P_1 \left[ 1 - \int_{R_1} f_1(\mathbf{x}) d\mathbf{x} \right] \\ &= \int_{R_1} [C(1/2)p_2 f_2(\mathbf{x}) - C(2/1)p_1 f_1(\mathbf{x})] d\mathbf{x} + C(2/1)p_1 \rightarrow (7) \end{aligned}$$

Now  $P_1, P_2, C(1/2)$  and  $C(2/1)$  are non-negative. In addition  $f_1(\mathbf{x})$  &  $f_2(\mathbf{x})$  are Non-negative for all  $\mathbf{x}$  and are the only quantities in ECM that depend on  $\mathbf{x}$ . Therefore, minimization of ECM is equivalent to minimize the function

$$\int_{R_1} [C(1/2)p_2 f_2(\mathbf{x}) - C(2/1)p_1 f_1(\mathbf{x})] d\mathbf{x} \rightarrow (8)$$

But, from the theory of integration (8) will be minimized is  $R_1$  includes there values of  $\mathbf{x}$  for which the integrand

$$C(1/2) P_2 f_2(\mathbf{x}) - C(2/1) p_1 f_1(\mathbf{x}) \leq 0 \rightarrow (9)$$

and for all  $\mathbf{x}$  those not included in  $R_1$  or equivalently for all  $\mathbf{x}$  those included in  $R_2$

$$C(1/2) p_2 f_2(\mathbf{x}) - C(2/1) p_1 f_1(\mathbf{x}) > 0 \rightarrow (10)$$

Thus from (9),

$$\begin{aligned} R_1 &= \{\mathbf{x} / C(1/2) p_2 f_2(\mathbf{x}) - C(2/1) p_1 f_1(\mathbf{x}) \leq 0\} \\ &= \{\mathbf{x} / C(2/1) p_1 f_1(\mathbf{x}) \geq C(1/2) p_2 f_2(\mathbf{x})\} \\ &= \left\{ \mathbf{x} / \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \left( \frac{C(1/2)}{C(2/1)} \right) / \left( \frac{p_1}{p_2} \right) \right\} \rightarrow (11) \end{aligned}$$

( $\because$  all  $f_1, f_2, p_1, p_2, C(1/2)$  &  $C(2/1)$  are all positive)

similarly from (10),

$$R_2 = \left\{ \mathbf{x} / \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \left( \frac{C(1/2)}{C(2/1)} \right) / \left( \frac{p_1}{p_2} \right) \right\} \rightarrow (12)$$

where (12) gives (2).

**REMARK:**

It is clear from (1) & (2) that the implementation of the minimum ECM

rules requires

(1). The ratio of p.d.f.'s is  $f_1 / f_2$  is to be evaluated at a new observation  $\mathbf{x}_0$ .

(2). The cost ratio  $\frac{C(1/2)}{C(2/1)}$

(3). The prior probability ratio  $\frac{p_1}{p_2}$

The appearance of ratios in the definition of the optimal classification regions has significance as often it is much easier to specify the ratios than their component parts.

### SPECIAL CASES OF MINIMUM "ECM" REGIONS

#### CASE (1):

(Equal prior probabilities i.e.  $p_1 = p_2$  or  $\frac{p_1}{p_2} = 1$ )

In this case (1) & (2) become

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{C(1/2)}{C(2/1)};$$

$$R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{C(1/2)}{C(2/1)}$$

#### CASE (2):

(Equal misclassification costs that is  $C(1/2)=C(2/1)$ )

In this case  $\frac{C(1/2)}{C(2/1)} = 1$  and therefore (1) & (2) become

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{p_2}{p_1};$$

$$R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1}$$

#### CASE (3):

$p_1 = p_2$  &  $C(1/2) = C(2/1)$

In this case  $\frac{p_1}{p_2} = 1 = \frac{C(1/2)}{C(2/1)}$  and therefore (1) & (2) become

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq 1;$$

$$R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < 1$$

#### NOTE:

(1). When the prior probabilities are not known, they are often taken to be

equal.

(2). Similarly when the misclassification costs are unknown, they are often taken to be equal.

(3). If  $\frac{C(1/2)}{C(2/1)} = \frac{p_1}{p_2}$  then  $C(1/2)p_2 = C(2/1)p_1$  and hence

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq 1;$$

$$R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < 1$$

### OPTIMAL TOTAL PROBABILITY OF MISCLASSIFICATION (TPM) REGIONS :

Criteria other than the ECM can be used to derive “optimal” classification procedures.

For example, one might ignore the costs of misclassification and

Choose  $R_1$  &  $R_2$  to minimize the total probability of misclassification (TPM).

TPM=P(misclassifying as  $\pi_1$  observation or misclassifying a  $\pi_2$  observation)

= P( $\mathbf{x}$  comes from  $\pi_1$  and is misclassified )+

P( $\mathbf{x}$  comes from  $\pi_2$  and is misclassified)

$$\Rightarrow TPM = P(\mathbf{X} \in R_2 / \pi_1).P(\pi_1) + P(\mathbf{X} \in R_1 / \pi_2).P(\pi_2)$$

$$= P_1 \int_{R_2} f_1(\mathbf{x})d\mathbf{x} + P_2 \int_{R_1} f_2(\mathbf{x})d\mathbf{x}$$

$$= p_1P(2/1) + p_2P(1/2) \quad \rightarrow (1)$$

But, when  $C(1/2)=C(2/1)$  (i.e. when misclassification costs are equal)

we get from equation (12) of page 14,

$$ECM = C(1/2)[p_1P(2/1) + p_2P(1/2)] \quad \rightarrow (2)$$

Now, from (1) & (2), it can be easily seen that minimizing (1) is equivalent

to minimizing (2). In other words, minimizing TPM is equivalent to

minimizing ECM with equal misclassification costs. Thus the optional TPM

regions  $R_1$  &  $R_2$  are same as those given in case(2) of page 20. Thus

$$\begin{aligned}
R_1 &= \left\{ \mathbf{x} / \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{p_2}{p_1} \right\} \\
R_2 &= \left\{ \mathbf{x} / \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1} \right\}
\end{aligned}
\rightarrow (3)$$

**ALLOCATING AN NEW OBSERVATION  $\mathbf{x}_0$  BASED ON BAYE'S  
POSTERIOR PROBABILITIES**

We can also allocate a new observation  $\mathbf{x}_0$  to the population with the largest posterior probability  $P(\pi_i / \mathbf{x}_0)$ . By Baye's rule, the "posterior" probabilities are

$$\begin{aligned}
P(\pi_1 / \mathbf{x}_0) &= P(\pi_1 \text{ occurs and observe } \mathbf{x}_0) / P(\text{observe } \mathbf{x}_0) \\
&= P(\text{observe } \mathbf{x}_0 / \pi_1) \cdot P(\pi_1) / \\
&\quad P(\text{observe } \mathbf{x}_0 / \pi_1) \cdot P(\pi_1) + P(\text{observe } \mathbf{x}_0 / \pi_2) \cdot P(\pi_2) \\
&= \frac{f_1(\mathbf{x}_0) \cdot p_1}{f_1(\mathbf{x}_0) \cdot p_1 + f_2(\mathbf{x}_0) \cdot p_2} \\
&= \frac{p_1 f_1(\mathbf{x}_0)}{p_1 f_1(\mathbf{x}_0) + p_2 f_2(\mathbf{x}_0)}
\end{aligned}
\rightarrow (1)$$

$$\begin{aligned}
P(\pi_2 / \mathbf{x}_0) &= 1 - P(\pi_1 / \mathbf{x}_0) \\
&= 1 - \frac{p_1 f_1(\mathbf{x}_0)}{p_1 f_1(\mathbf{x}_0) + p_2 f_2(\mathbf{x}_0)} \\
&= \frac{p_2 f_2(\mathbf{x}_0)}{p_1 f_1(\mathbf{x}_0) + p_2 f_2(\mathbf{x}_0)}
\end{aligned}
\rightarrow (2)$$

Now classify an observation  $\mathbf{x}_0$  into  $\pi_1$  when

$$\begin{aligned}
P(\pi_1 / \mathbf{x}_0) &> P(\pi_2 / \mathbf{x}_0) \\
\Rightarrow p_1 f_1(\mathbf{x}_0) &> p_2 f_2(\mathbf{x}_0) \\
&(\because \text{Numerators of (1) \& (2) are equal}) \\
\Rightarrow \frac{f_1(\mathbf{x}_0)}{f_2(\mathbf{x}_0)} &> \frac{p_2}{p_1}
\end{aligned}
\rightarrow (3)$$

Now from(3), it can be seen that allocating a new observation to a population based on Baye's posterior probabilities is same as optional TPM rule.

**NOTE:**

The above method is also equivalent to classify a new observation using optional ECM (Baye's method ) rule when misclassification costs are equal.

**CLASSIFICATION INTO ONE OF TWO KNOWN MULTIVARIATE NORMAL POPULATIONS (With common covariance matrix  $\Sigma$  )**

Classification procedures based on normal populations predominate in statistical practice because of their simplicity and reasonably high efficiency across a wide variety of population models. We assume  $f_1(\mathbf{x})$  &  $f_2(\mathbf{x})$  are multivariate normal densities; the first with mean vector  $\underline{\mu}_1$ , and the second with mean vector  $\underline{\mu}_2$  and both with common matrix  $\Sigma$ . Now the p.d.f. of the two populations  $\pi_1$  &  $\pi_2$  is given by

$$f_i(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\underline{\mu}_i)' \Sigma^{-1} (\mathbf{x}-\underline{\mu}_i)} \quad \text{for } i=1,2,\dots \rightarrow (1)$$

The ratio of densities after simplification is

$$\begin{aligned} \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} &= e^{-\frac{1}{2}(\mathbf{x}-\underline{\mu}_1)' \Sigma^{-1} (\mathbf{x}-\underline{\mu}_1) + \frac{1}{2}(\mathbf{x}-\underline{\mu}_2)' \Sigma^{-1} (\mathbf{x}-\underline{\mu}_2)} \\ &= e^{-\frac{1}{2}\mathbf{x}' \Sigma^{-1} \mathbf{x} + \frac{1}{2}\mathbf{x}' \Sigma^{-1} \underline{\mu}_1 + \frac{1}{2}\underline{\mu}_1' \Sigma^{-1} \mathbf{x} - \frac{1}{2}\underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 + \frac{1}{2}\mathbf{x}' \Sigma^{-1} \mathbf{x} - \frac{1}{2}\mathbf{x}' \Sigma^{-1} \underline{\mu}_2 - \frac{1}{2}\underline{\mu}_2' \Sigma^{-1} \mathbf{x} + \frac{1}{2}\underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2} \\ &\Rightarrow \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} = e^{\underline{\mu}_1' \Sigma^{-1} \mathbf{x} - \underline{\mu}_2' \Sigma^{-1} \mathbf{x} - \frac{1}{2}(\underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 - \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2)} \\ &\quad (\because \underline{\mu}_i' \Sigma^{-1} \mathbf{x} = \mathbf{x}' \Sigma^{-1} \underline{\mu}_i \text{ for } i=1,2) \\ &= e^{[(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \mathbf{x} - \frac{1}{2}(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2)]} \rightarrow (2) \end{aligned}$$

$$\begin{aligned} \text{for } (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) &= \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 + \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_2 - \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_1 - \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2 \\ &= \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_1 - \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_2 \quad (\because \underline{\mu}_1' \Sigma^{-1} \underline{\mu}_2 = \underline{\mu}_2' \Sigma^{-1} \underline{\mu}_1) \end{aligned}$$

By minimum ECM classification rule, we have

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \left( \frac{C(1/2)}{C(2/1)} \right) \left( \frac{p_2}{p_1} \right)$$

$$R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \left( \frac{C(1/2)}{C(2/1)} \right) \left( \frac{p_2}{p_1} \right)$$

From (2) we have, after taking logarithms on both sides

$$R_1 : (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \mathbf{x} - \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \geq \log K$$

$$R_2 : (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \mathbf{x} - \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) < \log K \quad \rightarrow (3)$$

where  $K = \frac{C(1/2) \cdot p_2}{C(2/1) p_1}$

The regions  $R_1$  &  $R_2$  given by (3) are called as minimum ECM regions for two normal populations.

#### NOTES:

(1). The first term of (3) viz

$$(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \mathbf{x} = \mathbf{l}' \mathbf{x}, \quad \text{where } \mathbf{l} = \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) \quad \rightarrow (4)$$

is the well known Fisher (linear) discriminant function, which is actually obtained by Fisher with entirely different argument which we will discuss later.

(2). The minimum ECM classification rule for two normal populations is given by allocating  $\mathbf{x}_0$  to  $\pi_1$  is

$$(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \mathbf{x}_0 - \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \geq \log K \quad \rightarrow (5)$$

$$\text{where } K = \frac{C(1/2) p_2}{C(2/1) p_1}$$

otherwise, allocate  $\mathbf{x}_0$  to  $\pi_2$

(3). The regions  $R_1$  &  $R_2$  given by equation (3) are called as best regions of classification.

(4). The regions  $R_1$  &  $R_2$  given by equation (3) are called as Baye's regions of classification.

(5). If the misclassification costs are equal i.e.  $C(1/2)=C(2/1)$ , then the regions

$R_1$  &  $R_2$  given by (3) with  $K = \frac{p_2}{p_1}$  are called as minimum(optional) TPM

regions for two normal populations.

(6). The minimum or optional TPM classification rule for two normal populations is given by allocate  $\underline{x}_0$  to  $\pi_1$  if

$$(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x}_0 - \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \geq \log K \quad \rightarrow (6)$$

$$\text{where } K = \frac{p_2}{p_1}$$

(7). If misclassification costs are equal and prior probabilities are equal i.e.

$C(1/2)=C(2/1)$  and  $p_1 = p_2$ , then  $K=1$  and consequently  $\log K=0$  and

those from (3), we have

$$R_1 : (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} \geq \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \quad \rightarrow (7)$$

$$\text{and } R_2 : (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x} < \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2)$$

In this case, the classification rule is allocate  $\underline{x}_0$  to  $\pi_1$  if

$$(\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{x}_0 \geq \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \quad \rightarrow (8)$$

allocate  $\underline{x}_0$  to  $\pi_2$  otherwise.

This case may be used when both misclassification costs as well as prior probabilities are unknown.

(8). Let

$$\begin{aligned} U &= (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \underline{X} - \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 + \underline{\mu}_2) \\ &= (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \left[ \underline{X} - \frac{1}{2} (\underline{\mu}_1 + \underline{\mu}_2) \right] \end{aligned}$$

Now let us obtain the distribution of U. since U is a linear combination of multivariate normal vector  $\underline{X}$ , U is distributed as univariate normal.

If  $\underline{X} \sim N_p(\underline{\mu}_1, \Sigma)$

$$\begin{aligned} E(U) &= \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) \\ &= \frac{\alpha}{2}, \text{ where } \alpha = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) \end{aligned}$$

$$\begin{aligned}
V(U) &= (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \text{cov}(\underline{X}, \underline{X}') \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) \\
&= (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} \Sigma \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) \\
&= (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2) \\
&= \alpha
\end{aligned}$$

Thus  $U \sim N\left(\frac{\alpha}{2}, \alpha\right)$ .

If  $\underline{X} \sim N_p(\underline{\mu}_2, \Sigma)$

In this case, we may show that  $U \sim N\left(\frac{-\alpha}{2}, \alpha\right)$ .

Here  $\alpha$  is the mahalanobis squared distance  $D^2$  between  $\underline{\mu}_1$  &  $\underline{\mu}_2$ .

### CLASSIFICATION INTO ONE OF TWO MULTIVARIATE NORMAL POPULATIONS WHEN THE PARAMETREERS ARE UNKNOWN

Suppose  $\underline{X}_{11}, \underline{X}_{12}, \dots, \underline{X}_{1n_1}$ , be a random sample of size ' $n_1$ ', from population  $\pi_1 : N(\underline{\mu}_1, \Sigma)$  and let  $\underline{X}_{21}, \underline{X}_{22}, \dots, \underline{X}_{2n_2}$  be a random sample of size ' $n_2$ ' from population  $\pi_2 : N(\underline{\mu}_2, \Sigma)$ . Since  $\underline{\mu}_1, \underline{\mu}_2$  &  $\Sigma$  are unknown we replace them with their unbiased estimators viz.,

$$\bar{\underline{X}}_1 = \frac{1}{n_1} \sum_{\alpha=1}^{n_1} \underline{X}_{1\alpha}, \bar{\underline{X}}_2 = \frac{1}{n_2} \sum_{\alpha=1}^{n_2} \underline{X}_{2\alpha} \quad \rightarrow (1)$$

and 
$$S = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2} \quad \rightarrow (2)$$

where

$$\begin{aligned}
S_1 &= \frac{1}{n_1 - 1} \sum_{\alpha=1}^{n_1} (\underline{X}_{1\alpha} - \bar{\underline{X}}_1)(\underline{X}_{1\alpha} - \bar{\underline{X}}_1)' \\
S_2 &= \frac{1}{n_2 - 1} \sum_{\alpha=1}^{n_2} (\underline{X}_{2\alpha} - \bar{\underline{X}}_2)(\underline{X}_{2\alpha} - \bar{\underline{X}}_2)'
\end{aligned} \quad \rightarrow (3)$$

Now, the estimated (or sample) minimum ECM regions can be obtained from the above method replacing  $\underline{\mu}_1, \underline{\mu}_2$  &  $\Sigma$  with their unbiased estimators  $\bar{\underline{X}}_1, \bar{\underline{X}}_2$  &  $S$  (given by (1) & (2)) respectively. They are form equations as follows :



$$\hat{R}_1 : (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)' S^{-1} \mathbf{X} - \frac{1}{2} (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)' S^{-1} (\bar{\mathbf{X}}_1 + \bar{\mathbf{X}}_2) \geq \log [c(1/2)p_2/c(2/1)p_1] \rightarrow (4)$$

$$\hat{R}_2 : (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)' S^{-1} \mathbf{X} - \frac{1}{2} (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)' S^{-1} (\bar{\mathbf{X}}_1 + \bar{\mathbf{X}}_2) < \log [c(1/2)p_2/c(2/1)p_1] \rightarrow (5)$$

from (4)&(5), the estimated sample minimum classification ECM rule for two normal populations is given by

Allocate  $\mathbf{X}_0$  to  $\pi_1$  if

$$(\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)' S^{-1} \mathbf{X} - \frac{1}{2} (\bar{\mathbf{X}}_1 - \bar{\mathbf{X}}_2)' S^{-1} (\bar{\mathbf{X}}_1 + \bar{\mathbf{X}}_2) \geq \log K \rightarrow (6)$$

$$\text{where } K = [c(1/2)p_2/c(2/1)p_1]$$

Allocate  $\mathbf{X}_0$  to  $\pi_2$  otherwise.

**NOTE:**

- (1) The estimated or sample minimum TPM rule for two normal populations with unknown parameters can be obtained from (6) replacing  $K$  with  $(p_2 / p_1)$ .
- (2) When  $p_1 = p_2$  &  $c(1/2) = c(2/1)$ , the estimated or sample minimum ECM rule is equivalent to sample ML rule and is given by

Allocate  $\mathbf{x}_0$  to  $\pi_1$  if

$$(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' S^{-1} \mathbf{x}_0 \geq \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' S^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \rightarrow (6)$$

Allocate  $\mathbf{x}_0$  to  $\pi_2$  otherwise.

- (3) The estimated minimum ECM rule or sample ML rule amounts to comparing the scalar variable (univariate normal variable)

$$y = \hat{l}' \mathbf{X} \quad , \text{ where } \hat{l} = S^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \rightarrow (7)$$

evaluated at  $\mathbf{x}_0$  is

$$y_0 = \hat{l}' \mathbf{x}_0$$

with the number

$$\begin{aligned} \hat{m} &= \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' S^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \\ &= \frac{1}{2} (\bar{y}_1 + \bar{y}_2) \end{aligned} \rightarrow (8)$$

where

$$\bar{y}_1 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^{-1} S^{-1} \bar{\mathbf{x}}_1 = \hat{l}' \bar{\mathbf{x}}_1$$

$$\bar{y}_2 = \hat{l}' \bar{\mathbf{x}}_2$$

Thus allocate  $\mathbf{x}_0$  to  $\pi_1$  if

$$y_0 \geq \hat{m} \quad \rightarrow (9)$$

otherwise allocate  $\mathbf{x}_0$  to  $\pi_2$

That is, the estimated minimum ECM rule for two normal populations is

to creating two univariate normal populations for the y values by taking an appropriate linear combination of the observations from populations

$\pi_1$  and  $\pi_2$  and then assigning a new new observation  $\mathbf{x}_0$  to  $\pi_1$  or  $\pi_2$  depending upon whether  $y_0 = \hat{l}' \mathbf{x}_0$  falls to the right or left of the midpoint  $\hat{m}$ , between the two normal means  $\bar{y}_1$  and  $\bar{y}_2$ .

- (4) The linear function (7) is known as Fisher linear discriminant function. Which is obtained by Fisher with a different argument for separating two populations.

## CLASSIFICATION OF NORMAL POPULATION WHEN $\Sigma_1 \neq \Sigma_2$

Here we have  $\pi_1 : N(\boldsymbol{\mu}_1, \Sigma_1)$  and  $\pi_2 : N(\boldsymbol{\mu}_2, \Sigma_2)$  when  $\Sigma_1 \neq \Sigma_2$ .

Let  $f_1(\mathbf{x})$  be the p.d.f. of  $\pi_1$  and  $f_2(\mathbf{x})$  be the p.d.f. of  $\pi_2$ . Then on simplification,

$$\log \left[ \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \right] = \log f_1(\mathbf{x}) - \log f_2(\mathbf{x})$$

$$= 1/2 \mathbf{x}' (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1' \Sigma_1^{-1} - \boldsymbol{\mu}_2' \Sigma_2^{-1}) \mathbf{x} - \lambda \quad \rightarrow (1)$$

$$\text{where } \lambda = 1/2 \log \left( \frac{|\Sigma_1|}{|\Sigma_2|} \right) + 1/2 (\boldsymbol{\mu}_1' \Sigma_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2' \Sigma_2^{-1} \boldsymbol{\mu}_2) \quad \rightarrow (2)$$

we have general formula for minimum ECM region and is given by

$$R_1 : \log [f_1(\mathbf{x}) / f_2(\mathbf{x})] \geq \log k, \text{ where } K = c(1/2)p_2 / c(2/1)p_1$$

$$R_2 : \log [f_1(\mathbf{x}) / f_2(\mathbf{x})] < \log k \quad \rightarrow (3)$$

Now, the minimum ECM regions for classification of two normal populations when  $\Sigma_1 \neq \Sigma_2$  is given by:

$$R_1 : -1/2 \mathbf{x}' (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1' \Sigma_1^{-1} - \boldsymbol{\mu}_2' \Sigma_2^{-1}) \mathbf{x} - \lambda \geq \log k$$

$$R_2 : -1/2 \mathbf{x}' (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1' \Sigma_1^{-1} - \boldsymbol{\mu}_2' \Sigma_2^{-1}) \mathbf{x} - \lambda < \log k$$

where  $\lambda$  &  $k$  are given as (2) & (3)  $\rightarrow$  (4)

The allocation rule that minimizes the ECM is given by :

Allocate  $\mathbf{x}_0$  to  $\pi_1$  if

$$-1/2\mathbf{x}'_0(\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\Sigma}_2^{-1})\mathbf{x}_0 + (\boldsymbol{\mu}'_1\boldsymbol{\Sigma}_1^{-1} - \boldsymbol{\mu}'_2\boldsymbol{\Sigma}_2^{-1})\mathbf{x}_0 - \lambda \geq \log k \rightarrow (5)$$

Allocate  $\mathbf{x}_0$  to  $\pi_2$  otherwise.

In practice, the classification rule in (5) is implemented by substituting the sample quantities  $\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2, S_1$  and  $S_2$  for  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1$  and  $\boldsymbol{\Sigma}_2$  respectively.

### QUADRATIC CLASSIFICATION RULE (NORMAL POPULATIONS WITH $\boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2$ )

Allocate  $\mathbf{x}_0$  to  $\pi_1$  if

$$1/2\mathbf{x}'_0(S_1^{-1} - S_2^{-1})\mathbf{x}_0 + (\bar{\mathbf{x}}'_1S_1^{-1} - \bar{\mathbf{x}}'_2S_2^{-1})\mathbf{x}_0 - \hat{\lambda} \geq \log k \rightarrow (6)$$

allocate  $\mathbf{x}_0$  to  $\pi_2$  otherwise.

$$\text{Where } \hat{\lambda} = 1/2 \log \left( \frac{|S_1|}{|S_2|} \right) + 1/2(\bar{\mathbf{x}}'_1S_1^{-1}\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}'_2S_2^{-1}\bar{\mathbf{x}}_2) \rightarrow (7)$$

#### NOTE:

- (1). Minimum TPM rule or quadratic classification rule when  $\boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2$  is a special case of (6) when  $K = p_2 / p_1$ .
- (2). If the misclassification costs are equal and prior probabilities are equal (i.e.  $C(1/2) = C(2/1)$  &  $p_1 = p_2$ ). Then the MC rule or QCR is obtained by taking  $K=1$  or  $\log K=0$  in the rule (6).

### FISHERS DISCRIMINANT FUNCTION-SEPARATION OF TWO POPULATIONS (NOT NECESSARY MULTIVARIATE NORMAL)

Fishers idea was to transform the multivariate observations  $\mathbf{x}$ 's to univariate observation  $y$ 's such that the  $y$ 's derived from population  $\pi_1$  and  $\pi_2$  were separated as much as possible . Fisher suggested taking linear combination of  $\mathbf{x}$ 's to create  $y$ 's because they are simple function of  $\mathbf{x}$  and are easily handled mathematically .

Fisher's approach does not assume that the populations are normal.

If does, however, implicitly assume the population covariance matrices are equal because a pooled estimate of the common covariance matrix is used.

Let  $\mathbf{x}_{11}, \mathbf{x}_{12}, \dots, \mathbf{x}_{1n_1}$  be a random sample of size  $n_1$  from population  $\pi_1$  and let

$\mathbf{x}_{21}, \mathbf{x}_{22}, \dots, \mathbf{x}_{2n_2}$  be a random sample of size  $n_2$  from population  $\pi_2$ . Now

$\bar{\mathbf{x}}_1$  be the mean of 1<sup>st</sup> sample

$S_1$  be the sample covariance matrix of 1<sup>st</sup> sample

$\bar{\mathbf{x}}_2$  be the mean of 2<sup>nd</sup> sample

$S_2$  be the sample covariance matrix of 2<sup>nd</sup> sample

$$\text{Denote } S = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2} \rightarrow (1)$$

Which is a pooled sampled covariance matrix.

Now, Fisher's idea is as follows

Consider the linear combination

$$y = \mathbf{w}'\mathbf{x}, \text{ when } \mathbf{w} \text{ is } |\mathbf{x}| \text{ vector of real number } \rightarrow (2)$$

using the linear transformation, the multivariate observation of 1<sup>st</sup> sample will be transformed into univariate observations given by

$$y_{11}, y_{12}, \dots, y_{1n_1}$$

$$\text{when } y_{1i} = \mathbf{w}'\mathbf{x}_{1i}, i = 1, 2, \dots, n_1$$

similarly the second sample  $\mathbf{x}_{21}, \mathbf{x}_{22}, \dots, \mathbf{x}_{2n_2}$  will be transformed into

$$y_{21}, y_{22}, \dots, y_{2n_2}$$

$$\text{when } y_{2i} = \mathbf{w}'\mathbf{x}_{2i}, i = 1, 2, \dots, n_2$$

$$\text{Now } \bar{y}_1 = \mathbf{w}'\bar{\mathbf{x}}_1$$

$$\bar{y}_2 = \mathbf{w}'\bar{\mathbf{x}}_2$$

$$\text{and } s_y^2 = \frac{\sum_{i=1}^{n_1} (y_{1i} - \bar{y}_1)^2 + \sum_{i=1}^{n_2} (y_{2i} - \bar{y}_2)^2}{n_1 + n_2 - 2} \rightarrow (3)$$

consider

$$\text{separation} = \frac{|\bar{y}_1 - \bar{y}_2|}{s_y} \rightarrow (4)$$

Now, Fisher's idea is to select the linear combination  $\mathbf{w}$  such that the separation given in (4) is maximum. In other words, the objective is to select the linear combination of  $\mathbf{x}$  (i.e.  $\mathbf{w}'\mathbf{x}$ )

to achieve maximum separation between the sample means  $\bar{y}_1$  &  $\bar{y}_2$ . Equation (4) may be written as

$$\begin{aligned} \text{separation}^2 &= \frac{(\text{squared distance between sample mean } \bar{y}_1 \text{ and } \bar{y}_2)}{\text{pooled sample variance of } y} \\ &= \frac{(\bar{y}_1 - \bar{y}_2)^2}{s_y^2} \\ &= \frac{[\underline{\mathbf{w}}'(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)]^2}{s_y^2} \quad (\text{from (3)}) \end{aligned}$$

but  $s_y^2 = \underline{\mathbf{w}}'S\underline{\mathbf{w}}$  (from (3) & (1))

$$\therefore \frac{\text{squared distance between } \bar{y}_1 \text{ \& } \bar{y}_2}{\text{pooled variance of } y} = \frac{(\underline{\mathbf{w}}'\underline{\mathbf{d}})^2}{\underline{\mathbf{w}}'S\underline{\mathbf{w}}} = \phi \text{ say} \quad \rightarrow (5)$$

where  $\underline{\mathbf{d}} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$

Now, as per Fisher's idea, (5) has to be maximized w.r.t.  $\underline{\mathbf{w}}$ .

Which implies

$$\begin{aligned} \frac{\partial \phi}{\partial \underline{\mathbf{w}}} &= \frac{(\underline{\mathbf{w}}'S\underline{\mathbf{w}}) \frac{\partial}{\partial \underline{\mathbf{w}}} (\underline{\mathbf{w}}'\underline{\mathbf{d}})^2 - (\underline{\mathbf{w}}'\underline{\mathbf{d}})^2 \frac{\partial}{\partial \underline{\mathbf{w}}} (\underline{\mathbf{w}}'S\underline{\mathbf{w}})}{\underline{\mathbf{w}}'S\underline{\mathbf{w}}} = 0 \\ &\Rightarrow (\underline{\mathbf{w}}'S\underline{\mathbf{w}})2(\underline{\mathbf{w}}'\underline{\mathbf{d}})\underline{\mathbf{d}} - 2(\underline{\mathbf{w}}'\underline{\mathbf{d}})^2 S\underline{\mathbf{w}} = 0 \\ &\Rightarrow (\underline{\mathbf{w}}'S\underline{\mathbf{w}})\underline{\mathbf{d}} - (\underline{\mathbf{w}}'\underline{\mathbf{d}})S\underline{\mathbf{w}} = 0 \\ &\Rightarrow \underline{\mathbf{w}} = \frac{(\underline{\mathbf{w}}'S\underline{\mathbf{w}})}{(\underline{\mathbf{w}}'\underline{\mathbf{d}})} S^{-1}\underline{\mathbf{d}} \quad (\because S \text{ is positive defined matrix}) \\ &= CS^{-1}\underline{\mathbf{d}} \quad \rightarrow (6) \end{aligned}$$

where  $C = \frac{\underline{\mathbf{w}}'S\underline{\mathbf{w}}}{\underline{\mathbf{w}}'\underline{\mathbf{d}}}$  and C is ratio of two scalars thus  $\underline{\mathbf{w}}$  is a scalar multiplier of the vector

$S^{-1}\underline{\mathbf{d}}$ .

Using

$$\begin{aligned} \underline{\mathbf{w}} &= CS^{-1}\underline{\mathbf{d}} \text{ in (5) we get} \\ \phi &= \frac{C^2(\underline{\mathbf{d}}'S^{-1}\underline{\mathbf{d}})^2}{C^2\underline{\mathbf{d}}'S^{-1}SS^{-1}\underline{\mathbf{d}}} \\ &= \underline{\mathbf{d}}'S^{-1}\underline{\mathbf{d}} \quad \rightarrow (7) \end{aligned}$$

Now using  $\underline{\mathbf{w}} = CS^{-1}\underline{\mathbf{d}}$  in (5) we get

$$\phi = \underline{\mathbf{d}}'S^{-1}\underline{\mathbf{d}} \quad \rightarrow (8)$$

Thus, from (7) & (8), we can see that for either

$$\underline{\mathbf{w}} = CS^{-1}\underline{\mathbf{d}} \text{ or } \underline{\mathbf{w}} = S^{-1}\underline{\mathbf{d}}$$

the same ratio  $\phi$ , we are setting. Thus  $\phi$  will be maximized if we take

$$\mathbf{w} = S^{-1}\mathbf{d} = S^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \quad \rightarrow (9)$$

and the maximum value of  $\phi$  is

$$\begin{aligned} \phi_m &= \mathbf{d}'S^{-1}\mathbf{d} \\ &= (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'S^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \\ &= D^2 \quad (\text{say}) \end{aligned} \quad \rightarrow (10)$$

Now, the linear function

$$\begin{aligned} Y &= \mathbf{w}'\mathbf{X} \\ &= (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'S^{-1}\mathbf{X} \quad (\text{from (2) \& (9)}) \end{aligned} \quad \rightarrow (11)$$

is called as Fisher's linear discriminant function . and the maximum ratio  $D^2$ , where  $D^2$  given by (10), is called the sample squared distance or squared Mahalanobis distance between sample means  $\bar{\mathbf{x}}_1$  and  $\bar{\mathbf{x}}_2$ .

The linear discriminant function given by (11) converts the two multivariate samples into two univariate samples such that the corresponding univariate sample means are separated as much as possible to the relative to pooled sample variance .

We can employ (11) as a classification device as given below.

### **AN ALLOCATION RULE BASED ON FISHER'S DISCRIMINANT FUNCTION :**

We have the Fisher's linear discriminant function

$$y = \mathbf{w}'\mathbf{x}, \quad \text{where } \mathbf{w} = S^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \quad \rightarrow (1)$$

Let 'm' be the midpoint between  $\bar{y}_1$  and  $\bar{y}_2$  and is given by

$$\begin{aligned} m &= (\bar{y}_1 + \bar{y}_2)/2 \\ &= 1/2 (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'S^{-1}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \end{aligned} \quad \rightarrow (2)$$

Now , the allocation rule or classification rule based on Fisher's discriminant function is as follows:

Allocate  $\mathbf{x}_0$  to  $\pi_1$  ,if

$$y_0 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'S^{-1}\mathbf{x}_0 \geq m \text{ or } y_0 - m \geq 0$$

Allocate  $\mathbf{x}_0$  to  $\pi_2$  ,if

$$y_0 < m \text{ or } y_0 - m < 0 \quad \rightarrow (3)$$

### **NOTE:**

(1). If  $\pi_1 \sim \mu_1, \Sigma$  and  $\pi_2 \sim \mu_2, \Sigma$  then the Mahalanobis distance between  $\mu_1$  and  $\mu_2$  is

denoted by  $\Delta_{\underline{\mu}_1, \underline{\mu}_2}$  and is given by

$$\Delta_{\underline{\mu}_1, \underline{\mu}_2}^2 = (\underline{\mu}_1 - \underline{\mu}_2)' \Sigma^{-1} (\underline{\mu}_1 - \underline{\mu}_2)$$

(2).  $\Delta_{\underline{x}, \underline{\mu}}^2 = (\underline{x} - \underline{\mu})' \Sigma^{-1} (\underline{x} - \underline{\mu})$

(3). Mahalanobis  $D^2$  test statistic to test separation between  $\pi_1$  and  $\pi_2$  (or)

$$H_0: \underline{\mu}_1 = \underline{\mu}_2 \text{ vs } H_1: \underline{\mu}_1 \neq \underline{\mu}_2$$

suppose  $\pi_1: N_p(\underline{\mu}_1, \Sigma)$  and  $\pi_2: N_p(\underline{\mu}_2, \Sigma)$   $\bar{x}_1, S_1$  are the sample mean and sample covariance matrix of a sample drawn from  $\pi_1$  and  $\bar{x}_2, S_2$  are ...  $\pi_2$ .

Now Mahalanobis  $D^2$  test statistic is given by

$$D^2 = (\bar{x}_1 - \bar{x}_2)' S^{-1} (\bar{x}_1 - \bar{x}_2)$$

where  $S = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2}$

under  $H_0: \left( \frac{n_1 + n_2 - p - 1}{(n_1 + n_2 - 2)p} \right) \left( \frac{n_1 n_2}{n_1 + n_2} \right) D^2 \sim F_{p, n_1 + n_2 - p - 1}$

which can be used as for testing the significant difference  $\underline{\mu}_1 - \underline{\mu}_2$ . If  $H_0$  is rejected, we can conclude that the separation between the two populations  $\pi_1$  and  $\pi_2$  is significant.

(4). Two sample  $T^2$  and Mahalanobis  $D^2$  are closely associated as

$$T^2 = \left( \frac{n_1 n_2}{n_1 + n_2} \right) D^2$$

(5). In case of two normal populations with common covariance matrix, Fisher's method is corresponds to a particular case of minimum ECM rule with equal prior probabilities and equal costs of TPM rule with equal prior probabilities. Further, it is same as ML rule.

(6). The expression in minimum ECM rule for two multivariate normal populations  $w = (\bar{x}_1 - \bar{x}_2)' S^{-1} (\bar{x} - 1/2(\bar{x}_1 + \bar{x}_2))$  is frequently called Anderson's classification.

(7). Fisher's method is also a special case of allocation rule based on Bayesian posterior probabilities when the prior probabilities  $p_1$  and  $p_2$  are same for the case of two multivariate normal populations.

## CLASSIFICATION WITH SEVERAL POPULATIONS

### THE MINIMUM “TPM” METHOD:

Let  $\pi_1, \pi_2, \dots, \pi_g$  are  $g$  populations. Let  $f_i(\underline{x})$  be the density associated with population  $\pi_i, i=1, 2, \dots, g$ . Let

$$p_i = \text{the prior probability of population } \pi_i, i=1, 2, \dots, g \rightarrow (1)$$

Now, for a given new observation  $\underline{x}_0$ , we have to allocate it to one of the populations  $\pi_1, \pi_2, \dots, \pi_g$ . It is equivalent to, if the total sample space  $\Omega$  is divided into  $R_1, R_2, \dots, R_g$  disjoint regions, and if the new observation  $\underline{x}_0$  falls in region  $R_k$ , then  $\underline{x}_0$  will be allocated to the population  $\pi_k$ . Now, if  $\underline{x}_0$  is actually drawn from population  $\pi_i$  and if it is allocated to  $\pi_k$  ( $i \neq k$ ) then we say that this is misclassification.

Thus, if  $R_k$  be the set of  $\underline{x}$ 's classified as  $\pi_i$  and let

$$\begin{aligned} P(k/i) &= P(\text{classify } \underline{x} \text{ in to } \pi_k / \pi_i) \\ &= \int_{R_k} f_i(\underline{x}) d\underline{x} \end{aligned} \rightarrow (2)$$

The conditional probability of misclassifying an  $\underline{x}$  from  $\pi_1$  into  $\pi_2$  or  $\pi_3$  or... or  $\pi_g$  is given by

$$\begin{aligned} \text{CPM}(1) &= P(2/1) + P(3/1) + \dots + P(g/1) \\ &= \sum_{i=2}^g P(i/1) \end{aligned} \rightarrow (3)$$

In a similar manner, we can obtain the conditional probabilities of misclassification  $\text{CPM}(2), \text{CPM}(3), \dots, \text{CPM}(g)$ .

Multiplying each conditional CPM by its prior probability and summing gives the total probability of misclassification (TPM).

Thus total probability of misclassification can be obtained as follows:

$$\begin{aligned} \text{TPM} &= P(\text{misclassification an observation } \underline{x}) \\ &= P(\text{misclassification } \underline{x} \text{ from } \pi_1) \cdot P(\pi_1) + \dots \\ &\quad + P(\text{misclassification } \underline{x} \text{ from } \pi_g) \cdot P(\pi_g) \end{aligned}$$



$$\begin{aligned}
&= \sum_{i=1}^g CPM(i)P_i \quad \text{where } p_i = P(\pi_i) \\
&= \sum_{i=1}^g P_i \left( \sum_{k \neq i=1}^g P(k/i) \right) \quad \text{(from (3))} \\
&= \sum_{i=1}^g P_i \left( \sum_{k \neq i=1}^g \int_{R_k} f_i(\mathbf{x}) d\mathbf{x} \right) \quad \rightarrow (4)
\end{aligned}$$

**Determining optimal classification procedures based on TPM amounts choosing the mutually exclusive regions  $R_1, R_2, \dots, R_g$  such that (4) is minimum.**

**THE MINIMUM “ECM” METHOD:**

The expected cost of misclassification (ECM) can be obtained as follows:

ECM=Expected cost of misclassifying an observation  $\mathbf{x}$

$$\begin{aligned}
&= (\text{conditional expected cost of misclassifying an } \mathbf{x} \text{ from } \pi_1) P(\pi_1) \\
&\quad + \dots + (\text{conditional expected cost of misclassifying an } \mathbf{x} \text{ from } \pi_g). P(\pi_g).
\end{aligned}$$

$$= p_1.ECM(1) + p_2.ECM(2) + \dots + p_g.ECM(g)$$

$$= \sum_{i=1}^g p_i.ECM(i) \quad \rightarrow (5)$$

where  $ECM(i)$ =conditional expected cost of misclassifying an  $\mathbf{x}$  from  $\pi_i$

$$\begin{aligned}
&= \sum_{\substack{k=1 \\ k \neq i}}^g (\text{conditional expected cost of misclassifying an } \mathbf{x} \text{ from } \pi_i \text{ into } \pi_k) \\
&= \sum_{\substack{k=1 \\ k \neq i}}^g c(k/i) (\text{conditional probability of misclassifying an } \mathbf{x} \text{ from } \pi_i \text{ into } \pi_k)
\end{aligned}$$

$$ECM(i) = \sum_{\substack{k=1 \\ k \neq i}}^g c(k/i)p(k/i) \quad \rightarrow (6)$$

where  $c(k/i)$ =the cost of misclassifying  $\mathbf{x}$  to  $\pi_k$  when actually it belongs to  $\pi_i (i \neq k)$  and  $p(k/i)$  is as given by (2).

(Here it may be noted  $c(i/i) = 0$ , the misclassification cost for correct classification)

using (6) in (5), we get

$$\begin{aligned}
ECM &= \sum_{i=1}^g p_i \left[ \sum_{\substack{K=1 \\ K \neq i}}^g c(k/i) p(k/i) \right] \\
&= \sum_{i=1}^g p_i \left[ \sum_{\substack{K=1 \\ K \neq i}}^g c(k/i) \int_{R_k} f_i(\mathbf{x}) d\mathbf{x} \right] \quad \rightarrow (7)
\end{aligned}$$

Thus (7) gives the expected cost of misclassification , Now, an optional classification procedure obtained by minimizing ECM (given by (7)) for the choice of  $R_1, R_2, \dots, R_g$  is called minimum ECM method.

**The choice of mutually exclusive regions  $R_1, R_2, \dots, R_g$  for when (7) is minimum are called minimum ECM regions.**

### **THEOREM (MINIMUM “ECM” CLASSIFICATION RULE)**

The classification regions that minimize the

$$ECM = \sum_{i=1}^g p_i \left[ \sum_{\substack{K=1 \\ K \neq i}}^g c(k/i) \int_{R_k} f_i(\mathbf{x}) d\mathbf{x} \right]$$

are defined by allocating  $\mathbf{x}$  to that population  $\pi_k$ ,  $k=1,2,\dots,g$

for which

$$\sum_{\substack{i=1 \\ i \neq k}}^g p_i c(k/i) f_i(\mathbf{x}) \quad \rightarrow (8)$$

is smallest. If a tie occurs  $\mathbf{x}$  can be assigned to any of the tied populations.

**Minimum ECM classification rule with equal misclassification costs**

**(OR)**

**Minimum TPM classification rule:**

Suppose all the misclassification costs are equal , in which case the minimum ECM rule reduces to minimum TPM rule. Thus in this case minimum ECM classification rule becomes (from (8))

Allocate  $\mathbf{x}$  to  $\pi_k$  if

$$\sum_{\substack{i=1 \\ K \neq i}}^g p_i f_i(\mathbf{x}) \quad \rightarrow (9)$$

is smallest. If tie occurs  $\underline{x}$  can be assigned to any of the tied populations.

Now (9) will be smallest when the omitted term  $p_k f_k(\underline{x})$  is largest. Consequently, the minimum ECM rule with equal misclassification costs (OR) minimum TPM rule is as follows:

Allocate  $\underline{x}$  to  $\pi_k$  if

$$p_k f_k(\underline{x}) > p_i f_i(\underline{x}) \quad \forall i \neq k \quad \rightarrow (10)$$

or equivalently, Allocate  $\underline{x}$  to  $\pi_k$  if

$$\log(p_k f_k(\underline{x})) > \log(p_i f_i(\underline{x})) \quad \forall i \neq k \quad \rightarrow (11)$$

### NOTES:

- (1). It is interesting to note that the classification rule (10) is identical to the one that maximizes the posterior probability,

$$\begin{aligned} p(\pi_k / \underline{x}) &= p(\underline{x} \text{ comes from } \pi_k \text{ given that } \underline{x} \text{ was observed}) \\ &= \frac{p_k f_k(\underline{x})}{\sum_{i=1}^g p_i f_i(\underline{x})} = \frac{(\text{prior}) \times (\text{likelihood})}{\sum (\text{prior}) \times (\text{likelihood})} \end{aligned}$$

Thus rule (10) or rule (11) is identical with the rule obtained based on Bayes posterior Probabilities.

- (2). On case if prior probabilities are equal, rule (10) (or rule (11)) reduces to

Allocate  $\underline{x}$  to  $\pi_k$  if

$$f_k(\underline{x}) > f_i(\underline{x}) \quad \forall i \neq k \quad \rightarrow (12)$$

or equivalently Allocate  $\underline{x}$  to  $\pi_k$  if

$$\log f_k(\underline{x}) > \log f_i(\underline{x}) \quad \forall i \neq k \quad \rightarrow (13)$$

The rule (12) or (13) is called as ML rule which is a special case of minimum TPM rule as well as minimum ECM rule.

- (3). Generally , the minimum ECM rules have these components
- a) prior probabilities
  - b) misclassification costs
  - c) density functions

These components must be specified or estimated before the rules can be implemented.

### CLASSIFICATION WITH NORMAL POPULATIONS

We have 'g' multivariate normal populations

$$\begin{aligned}\boldsymbol{\pi}_1 &: N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \\ \boldsymbol{\pi}_2 &: N_p(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \\ &\vdots \\ \boldsymbol{\pi}_g &: N_p(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)\end{aligned}$$

In this case

$$\begin{aligned}f_i(\mathbf{x}) &= \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_i|^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_i)' \boldsymbol{\Sigma}_i^{-1}(\mathbf{x}-\boldsymbol{\mu}_i)} \quad \text{for } i=1,2,\dots,g \\ \log f_i(\mathbf{x}) &= -\left(\frac{p}{2}\right) \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \frac{1}{2} ((\mathbf{x}-\boldsymbol{\mu}_i)' \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}-\boldsymbol{\mu}_i)) \quad \rightarrow (14)\end{aligned}$$

### CASE (1): UNEQUAL $\boldsymbol{\Sigma}_i$

In this case the minimum TPM rule (or minimum ECM rule with equal misclassification costs) is as follows:

From (14) we have

$$\log(p_k f_k(\mathbf{x})) = \log p_k - \frac{p}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}-\boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}-\boldsymbol{\mu}_k)$$

Now the above rule becomes (from (11))

Allocate  $\mathbf{x}$  to  $\boldsymbol{\pi}_k$  if

$$\begin{aligned}\log(p_k f_k(\mathbf{x})) &= \max_i \log(p_i f_i(\mathbf{x})) \\ &= \log p_k - \frac{p}{2} \log(2\pi) - \frac{1}{2} \log |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x}-\boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}-\boldsymbol{\mu}_k) \quad \rightarrow (15)\end{aligned}$$

The constant  $p/2 \log(2\pi)$  can be ignored in (15) since it is same for all populations. We therefore define the quadratic discrimination score for  $i^{th}$  population is

$$d_i^Q(\mathbf{x}) = \log p_i - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \frac{1}{2} (\mathbf{x}-\boldsymbol{\mu}_i)' \boldsymbol{\Sigma}_i^{-1} (\mathbf{x}-\boldsymbol{\mu}_i) \quad \text{for } i=1,2,\dots,g \quad \rightarrow (16)$$

The quadratic score  $d_i^Q(\mathbf{x})$  is composed contributions from the generalized variance  $|\boldsymbol{\Sigma}_i|$ , the prior probability  $p_i$ , and Mahalanobis (or statistical) squared distance between  $\mathbf{x}$  and population mean  $\boldsymbol{\mu}_i$ .

Using discriminant scores the classification rule (15) becomes

Allocate  $\mathbf{x}$  to  $\boldsymbol{\pi}_k$

$$\text{The quadratic score } d_k^Q(\mathbf{x}) = \max_i \{d_i^Q(\mathbf{x})\} \quad \rightarrow (17)$$

Where  $d_i^Q(\mathbf{x})$  is given by (16).

In practice, the  $\mu_i$  and  $\Sigma_i$  are unknown and hence a training set of correctly classified observations is often available for the construction of estimates. The relevant sample quantities for population  $\pi_i$  are

$$\begin{aligned}\bar{\mathbf{x}}_i &= \text{sample mean vector} \\ S_i &= \text{sample covariance matrix and} \\ n_i &= \text{sample size}\end{aligned}$$

Using the above estimation in (16), we get the estimate of the quadratic discriminant score  $\hat{d}_i^o(\mathbf{x})$  as

$$\hat{d}_i^o(\mathbf{x}) = \log p_i - \frac{1}{2} \log |S_i| - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_i)' S_i^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i) \quad \rightarrow (18)$$

and the classification rule based on the sample is as follows:

### ESTIMATED MINIMUM “TPM” RULE FOR SEVERAL NORMAL POPULATIONS- UNEQUAL $\Sigma_i$

Allocate  $\mathbf{x}$  to  $\pi_k$  if the quadratic score

$$\hat{d}_i^o(\mathbf{x}) = \max_i \{ \hat{d}_i^o(\mathbf{x}) \} \quad \rightarrow (19)$$

where  $\hat{d}_i^o(\mathbf{x})$  is given by (18) for  $i = 1, 2, \dots, g$ .

### CASE-2:- $\Sigma_i$ 's ARE EQUAL

#### ESTIMATED MINIMUM “TPM” RULE (BAYES ALLOCATION RULE ) NORMAL POPULATION –EQUAL $\Sigma_i$ :

In the case we have  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma$  (say) and hence the discriminant score in (16) becomes

$$\hat{d}_i^o(\mathbf{x}) = -\frac{1}{2} \log |\Sigma| - \frac{1}{2} \mathbf{x}' \Sigma^{-1} \mathbf{x} + \mu_i' \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \log p_i$$

The first two terms are same for  $d_1^o(\mathbf{x}), d_2^o(\mathbf{x}), \dots, d_g^o(\mathbf{x})$  and consequently, they can be ignored for allocatory purposes. Define the linear discriminant score

$$d_i(\mathbf{x}) = \mu_i' \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \log p_i \quad \rightarrow (20)$$

An estimate of  $d_i(\mathbf{x})$  viz,  $\hat{d}_i(\mathbf{x})$  is based on the pooled estimate of  $\Sigma$ ,

$$S = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2 + \dots + (n_g - 1)S_g}{n_1 + n_2 + \dots + n_g - g} \quad \rightarrow (21)$$

and is given by

$$\hat{d}_i(\mathbf{x}) = \bar{\mathbf{x}}_i' S^{-1} \mathbf{x} - \frac{1}{2} \bar{\mathbf{x}}_i' S^{-1} \bar{\mathbf{x}}_i + \log p_i \quad \rightarrow (22)$$

consequently, the estimated minimum TPM rule for equal covariance normal populations is as follows:

$$\text{Allocate } \mathbf{x} \text{ to } \pi_k \text{ if the linear discriminant score } \hat{d}_k(\mathbf{x}) = \max_i \{\hat{d}_i(\mathbf{x})\} \quad \rightarrow (23)$$

Where  $\hat{d}_i(\mathbf{x})$  is given by (22)

**NOTE:-**

(1). In the above minimum TPM rules , for any case , if  $p_1 = p_2 = p_3 = \dots = p_g = 1/g$  ,

we may ignore those term  $\log p_i$  is discriminant scores , as it is same for all discriminant scores. In this case the minimum TPM rule is reduced to ML rule in which case the allocation rules are same as above except ignoring  $\log p_i$  .

(2). Expression (20) is a convenient linear function of  $\mathbf{x}$  . An equivalent classifier for equal covariance matrices case can be obtained from (16) by ignoring the term  $-\frac{1}{2} \log |\Sigma|$  and is given by

$$-\frac{1}{2} (\mathbf{x} - \underline{\mu}_i)' \Sigma^{-1} (\mathbf{x} - \underline{\mu}_i) + \log p_i$$

The classification rule with sample estimates instead for unknown populations quantities is given by Allocate  $\mathbf{x}$  to  $\pi_k$  ,if

$$-\frac{1}{2} D_k^2(\mathbf{x}) + \log p_k \text{ is largest for } k=1,2,\dots,g. \quad \rightarrow (24)$$

$$\text{where } D_k^2(\mathbf{x}) = (\mathbf{x} - \bar{\mathbf{x}}_k)' S^{-1} (\mathbf{x} - \bar{\mathbf{x}}_k)$$

is Mahalanobis squared distance between  $\mathbf{x}$  and the sample mean  $\bar{\mathbf{x}}_k$  .

Thus , we see the rule (24) or equivalently rule (23) assigns  $\mathbf{x}$  to the closest population (The distance is penalized by  $\log p_i$  ).

(3). In note(2) , if we assume  $p_1, p_2, p_3, \dots, p_g$  are equal and hence allocation rule may be significant as follows:

$$\text{Allocate } \mathbf{x} \text{ to } \pi_k, \text{ if } -\frac{1}{2} D_k^2(\mathbf{x}) \text{ is largest}$$

Or equivalently  $D_k^2(\mathbf{x})$  smallest  $\rightarrow (25)$

In other words, we are allocating  $\mathbf{x}$  to that population whose sample mean vector is closest to  $\mathbf{x}$ . This rule is also called as ML classification rule.

**FISHER'S METHOD FOR DISCRIMINATING AMONG SEVERAL POPULATIONS  
WHEN PARAMETERS ARE SPECIFIED**

Fisher also proposed a several population extension of his discriminant method, which was discussed for the case of two populations. The motivation behind the Fisher discriminant analysis is the need to obtain a reasonable representation of the population that involves only a few linear combinations of the observations, such as  $\mathbf{l}'_1\mathbf{x}, \mathbf{l}'_2\mathbf{x}$  and so on. His approach has several advantages and one is interested in separating several populations for

- 1) Visual inspection or
- 2) Graphical descriptive purposes.

It allows for the follows:-

1. Convenient representation of the  $g$  populations that reduce the dimension from a very large number of characteristics to a relatively few linear combinations. Of course, some information – needed for optimal classification- may be lost unless the population means lie completely in the lower dimensional space selected.
2. Plotting of the means of the first two or three linear combinations (discriminates). This helps display the relation ship and possible groupings of the populations.
3. Scatter plots of the sample values of the first two discriminates, which can indicate outliers or other abnormalities in the data.

The primary purpose of Fishers Discriminant analysis is to separate populations. However, it can also be used to classify a new observation into one of the populations. It is not necessary to assume that the  $g$  populations are multi variate normal. However we assume the population covariance matrices are equal and of full rank. That is  $\Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma$ . Thus, we have  $g$  populations with mean vectors  $\underline{\mu}_1, \underline{\mu}_2, \dots, \underline{\mu}_g$  and common covariance matrix  $\Sigma$ .

Let  $\bar{\mu} = \frac{1}{g} \sum_{i=1}^g \mu_i$

and  $B = \sum_{i=1}^g (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})'$  -----(1)

we consider the linear combination  $y = \bar{l}'X$

which has expected value

$$\begin{aligned} E(y) &= \bar{l}'E(X/\pi_i) = \bar{l}'\mu_i \text{ (for population } \pi_i) \\ &= \mu_{iy} \text{ (say)} \end{aligned}$$

and variance  $V(y) = \bar{l}'cov(X, X')\bar{l}$   
 $= \bar{l}'\Sigma\bar{l} = \sigma_y^2$  for all populations. -----(2)

we defuse the overall mean,

$$\begin{aligned} \bar{\mu}_y &= \frac{1}{g} \sum_{i=1}^g \mu_{iy} = \frac{1}{g} \sum_{i=1}^g \bar{l}'\mu_i \\ &= \bar{l}'\left(\frac{1}{g} \sum_{i=1}^g \mu_i\right) \\ &= \bar{l}'\bar{\mu} \quad \text{(From (1))} \end{aligned} \quad \text{-----(3) and form the ratio}$$

sum of squared distances from populations to over all mean of Y  
common population variance of Y

$$\begin{aligned} &= \frac{\sum_{i=1}^g (\mu_{iy} - \bar{\mu}_y)^2}{\sigma_y^2} \\ &= \frac{\sum_{i=1}^g (\bar{l}'\mu_i - \bar{l}'\bar{\mu})^2}{\bar{l}'\Sigma\bar{l}} \\ &= \frac{\bar{l}' \sum_{i=1}^g (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})' \bar{l}}{\bar{l}'\Sigma\bar{l}} \\ &= \frac{\bar{l}'B\bar{l}}{\bar{l}'\Sigma\bar{l}} \quad \text{(from (1))} \end{aligned}$$



$$\text{Thus } \frac{\sum_{i=1}^g (\mu_{iy} - \bar{\mu}_y)^2}{\sigma_y^2} = \frac{l' B l}{\tilde{l}' \tilde{\Sigma} \tilde{l}} \quad \text{-----(4)}$$

The ratio (4) measures the variability between the groups of Y- values relative to the common variability within the groups. We can then choose  $\tilde{l}$  to maximize the ratio (4) Thus if we write

$$\lambda = \frac{l' B l}{\tilde{l}' \tilde{\Sigma} \tilde{l}} \quad \text{-----(5)}$$

Then we have to maximize (5) with respect to  $\tilde{l}$  when implies

$$\begin{aligned} \frac{\partial \lambda}{\partial \tilde{l}} = 0 &\Rightarrow (l' \tilde{\Sigma} l) \frac{\partial l'}{\partial \tilde{l}} B l - (l' B l) \frac{\partial l'}{\partial \tilde{l}} \tilde{\Sigma} l = 0 \\ &\Rightarrow (l' \tilde{\Sigma} l) B l - (l' B l) \tilde{\Sigma} l = 0 \\ &\Rightarrow B l - \left( \frac{l' B l}{\tilde{l}' \tilde{\Sigma} \tilde{l}} \right) \tilde{\Sigma} l = 0 \\ &\Rightarrow \Sigma^{-1} B l - \left( \frac{l' B l}{\tilde{l}' \tilde{\Sigma} \tilde{l}} \right) l = 0 \\ &\Rightarrow (\Sigma^{-1} B - \lambda I) l = 0 \quad (\text{using (5)}) \quad \text{-----(6)} \end{aligned}$$

Thus  $\tilde{l}$  is the latent vector corresponding to a latent root  $\lambda$  of  $\Sigma^{-1} B$ . As, we are seeking for a  $\tilde{l}$  which maximizes  $\lambda$ , let  $\lambda_1$  be the non zero largest latent root of  $\Sigma^{-1} B$  and  $l_1$  be the corresponding latent vector. Now, the linear combination,  $Y_1 = l_1' X$  is called Fisher's first linear discriminant .

Similarly if  $\lambda_2$  is the next non Zero largest latent root of  $\Sigma^{-1} B$  and  $l_2$  correspondent latent vector then ,  $Y_2 = l_2' X$  is Fisher's second linear discriminant.

Let  $\lambda_1 > \lambda_2 > \dots > \lambda_s > 0$  denote the  $s \leq \min(g-1, p)$  non zero eigen values of  $\Sigma^{-1} B$  and let  $l_1, l_2, \dots, l_s$  be the corresponding latent vectors. Now, the linear combinations

$$Y_k = l_k' X \quad (k \leq s) \quad \text{.....(7)}$$

is Fisher's  $k^{\text{th}}$  linear discriminant.

**Fishers method for discriminating several populations when parameters are unknown**

**Fisher's sample linear discriminants:**

In general,  $\Sigma$  and the  $\mu_i$ 's are unknown, but we have a training set consisting of correctly classified observations. Suppose the training set consist of a random sample of size  $n_i$  from population  $\pi_i, i = 1, 2, 3, \dots, g$

Let  $\bar{x}_i$  be the mean vector and  $S_i$  be the covariance matrix of ith sample. Now denote the sample between groups matrix.

$$B_0 = \sum_{i=1}^g (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})'$$

Where, 
$$\bar{x} = \frac{1}{g} \sum_{i=1}^g \bar{x}_i \quad \dots\dots(8)$$

$B_0$  is an estimate of  $B$

Also, an estimate of  $\Sigma$  is based on the sample within groups matrix is

$$W = \sum_{i=1}^g (n_i - 1)S_i \quad \dots\dots(9)$$

Consequently, 
$$S_p = \frac{W}{(n - g)}, n = \sum_{i=1}^g n_i \quad \dots\dots(10)$$

is an estimate of  $\Sigma$ .

We consider the linear transformation,

$$y = l'x \quad \dots\dots(11)$$

Under the linear transformation, (11) the given multi variate samples can be transformed into univariate samples whose means and variances are given by

Means :  $\bar{y}_1, \bar{y}_2, \dots, \bar{y}_g$

Variances:  $s_{y_1}^2, s_{y_2}^2, \dots, s_{y_g}^2$

We denote the overall sample as

$$\bar{y} = \frac{1}{g} \sum_{i=1}^g \bar{y}_i \quad \dots\dots(12)$$

Now form the ratio,

$$\lambda = \frac{\text{sum of squared distances fro sample means to overall mean}}{\text{Total within samples variation}}$$

$$\begin{aligned} &= \frac{\sum_{i=1}^g (\bar{y}_i - \bar{y})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2} \\ &= \frac{\sum_{i=1}^g (l' \bar{x}_i - l' \bar{x})^2}{\sum_{i=1}^g \sum_{j=1}^{n_i} (l' x_{ij} - l' \bar{x}_i)^2} \quad (\text{from (11)}) \\ &= \frac{l' \sum_{i=1}^g (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})' l}{l' \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_j)(x_{ij} - \bar{x}_j)' l} \\ &= \frac{l' \sum_{i=1}^g (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x}_i)' l}{l' \sum_{i=1}^g (n_i - 1) S_i l} \end{aligned}$$

(By using the definition of sample covariance matrix)

$$= \frac{l' B_0 l}{l' W l} \quad (\text{from (8) \& (9) )} \quad \dots(13)$$

The ratio (13) measures the variability between the groups of g values relative to the total variability within the groups.

Now, Fisher suggested to choose  $l$  such that  $\lambda$  given by (13) is maximum, Maximization of  $\lambda$  with respect to  $l$  implies.

$$\frac{\partial \lambda}{\partial l} = 0 \Rightarrow (W^{-1} B_0 - \lambda I) l = 0 \quad \dots(14)$$

(See (6) of page 29 for derivation particulars)

Now, if we denote  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s > 0$  (where  $s = \min(g-1, p)$ ) are s eigen values of (14) and let  $l_1, l_2, \dots, l_s$  the corresponding eigen vectors, then

Fisher's K-th sample linear discriminant is given by

$$y_k = l'_k x(k \leq s)$$

Thus, Fisher's sample linear discriminants are eigen vectors of  $W^{-1}B_0$ ,

Where  $B_0$  and  $W$  are as assigned in (8) & (9).

**Note:**

- 1) If sample means  $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_g$  and sample covariance matrix  $S_1, S_2, \dots, S_g$  are given, then  $B_0$  and  $W$  can be completed using (8) and (9) respectively.
- 2) If raw samples from  $g$  populations are given, then  $B_0$  and  $W$  can be computed as follows:

First compute the individual sample covariance matrices  $S_1, S_2, \dots, S_g$  from the given samples and then use (9) to compute  $W$ . Now, compute the sample covariance matrix  $S$  from the combined samples of  $g$  samples given by

$$S = \frac{1}{n-1} \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x})(x_{ij} - \bar{x})', \text{ when } n = \sum_{i=1}^g n_i.$$

Now,  $B_0$  can be computed from the following relationship.

$$(n-1)S = W + B_0.$$

- 3) It may be noted that the pooled sample covariance matrix  $S_p$  and combined sample covariance matrix  $S$  are connected by the

$$(n-1)S = W + B_0$$

Thus, if the individual sample covariance matrix  $S_1, S_2, \dots, S_g$  and the combined sample covariance matrix  $S$  are given, then one can obtain  $W$  and  $B_0$  can be obtained as follows

$$W = (n_1 - 1)S_1 + (n_2 - 1)S_2 + \dots + (n_g - 1)S_g$$

$$B_0 = (n-1)S - W$$

Now, Fisher's discriminates can be constructed using the eigen vectors of  $W^{-1}B_0$ .

- 4) We know that  $W^{-1}B_0$  is not a symmetric matrix. Many computer

Packages can compute eigen values and eigen vectors only for the symmetric matrices. However, the eigen vectors of  $W^{-1}B_0$  can be computed as follows:

Suppose,  $\lambda$  is a ch root and  $\underline{l}$  is Ch. Vector of  $W^{-1}B_0$ , then we have

$$(W^{-1}B_0 - \lambda I)\underline{l} = 0$$

The above equations may be rewritten as

$$\begin{aligned} (W^{-1}B_0 - \lambda W^{-\frac{1}{2}}W^{-\frac{1}{2}})\underline{l} &= 0 \\ \Rightarrow (W^{-\frac{1}{2}}B_0 - \lambda W^{\frac{1}{2}})\underline{l} &= 0 \\ \Rightarrow (W^{-\frac{1}{2}}B_0W^{-\frac{1}{2}} - \lambda)W^{\frac{1}{2}}\underline{l} &= 0 \\ \Rightarrow (W^{-\frac{1}{2}}B_0W^{-\frac{1}{2}} - \lambda)\underline{w} &= 0 \end{aligned}$$

( $W^{\frac{1}{2}}$  is square root of  $W$  and  $W^{-\frac{1}{2}}$  is a inverse of  $W^{\frac{1}{2}}$ ).

Where,  $\underline{w} = W^{\frac{1}{2}}\underline{l}$  or  $\underline{l} = W^{-\frac{1}{2}}\underline{w}$ .

Thus if  $\underline{w}$  is latent vector of the matrix,  $W^{-\frac{1}{2}}B_0W^{-\frac{1}{2}}$ . Corresponding to the latent root of  $\lambda$ , then latest vector  $\underline{l}$  of  $W^{-\frac{1}{2}}B_0$  corresponding root  $\lambda$  may be obtained as  $\underline{l} = W^{-\frac{1}{2}}\underline{w}$

For all practical purposes, for the construction of Fisher's discriminant functions we use the above method.

### **Classification of a new observation among several populations using Fisher's discriminants:-**

Mainly, Fisher's discriminates were derived for the purpose of obtaining a low dimensional representation of the data that separate the populations as much as possible. Although they were derived from separatory considerations, the discriminates also provide the basis for a classification rule.

## CLUSTER ANALYSIS

### CONCEPT:

Multivariate methods deal with the analysis of data of more than two variables recorded from  $n$  sample objects selected from a specified population. Since the sample objects are selected from a specified population, the units are assumed to be homogeneous in respect of some characteristics. However, the values of different variables recorded from sample objects are not strictly uniform, though there should not be any systematic difference in the objects. In general, we expect some variations in the values of the variables, even if the sample objects are uniform in respect of some characters. For example, the income or the expenditure of middle class of people in a country are not exactly uniform, though they belong to the same class.

Again, the people of a country can be classified as rich, upper middle class, lower middle class and poor. For each class of people there may be common variable which influences the economic condition. For example, the income of a person depends on his education. This is true for every class of people. But their income or expenditure are not uniform. Therefore, there may be some systematic difference in values of the variables recorded from sample objects, there may be some similarities in the recorded observations of sample objects. Those sample objects which are similar in their recorded information may form a group. Dissimilar objects fall in different groups. In general, the objects that share similar characteristics are found together. In statistics, the search for relatively homogeneous objects is called cluster analysis.

The cluster analysis has wide applications in biology, medicine, agriculture, marketing, etc. The numerical taxonomy in the field of biology is used to classify the animals into class, order and families. Different species of plants have different characteristics. Therefore, plant specimens can be classify into homogeneous groups. In agriculture, the land fertility of a particular region may not be homogeneous for any type of crop. Then the pieces of land sharing similar fertility for a particular may be grouped together. The milk production of cows, even of the same type, may vary due to lactation period. Then the cows of the same lactation period may be grouped together. In

economics, the people of a city center may be grouped according to their socio-economic condition. In marketing, people can be grouped according to the similar buying habits. In medicine, the patients having similar disease may be clustered together.

Since similar objects form a cluster, all the sample points in any cluster will provide similar information about the population characteristics. Thus, for further analysis one may include one object from each cluster analysis is a data reduction technique in rows of the data matrix.

Let  $X(n \times p)$  be a data matrix from a specified population. Let the values of the  $p$  variables observed from  $n$  sample objects be denoted by  $X_1, X_2, \dots, X_n$ . The objective of the cluster analysis is to group these  $n$  vector of values into  $n_1$  ( $n_1 < n$ ) vectors so that the elements in a group are homogeneous. Here the method of clustering is on the basis of one-sample observations. Let  $X_{ij}[i = 1, 2, \dots, n_j; j = 1, 2, \dots, m]$  be the vector of values of  $p$  variables of  $i$ -th object in  $j$ -th sample. Here the objective of clustering is to form  $m_1$  groups ( $m_1 < m$ ) of sample observations in different groups are heterogeneous.

From above discussion it is clear that the CA reduces the sample observations in size. It has similar property of other data reduction technique. Namely, PCA. This analysis has a similarity with DA in respect of classification of observations. But DA derives a rule for allocating an object to its proper properties based on some prior information of the group membership of the object. Whereas, the CA identifies homogeneous groups or clusters.

There is no unified approach on what actually constitute a cluster. As per the definition what we have discussed above, a cluster constitutes with a similar object. Then, we need to decide on a measure of inter-object similarities. Also, a decision is needed to specify a procedure for forming the clusters, based on the chosen measure of similarity. The criterion of similarity in observations varies from researcher to researcher. However, the basic criterion is that the objects in a cluster should be closer to each other than to objects in other clusters. As a preliminary technique to identify the similarity of objects, one can use the diagram of sample objects. Let us consider that from each of 'n' sample object values of  $p$  variables are recorded. These values can be

represented a p-dimensional diagram. The values of each variable is plotted in each separate axis. If n sets of values are plotted in p-axes, a diagram will be formed. The cluster can be formed with those objects, which lie nearer in an area of the diagram but are dispersed from another area. The cluster can also be formed mathematically calculating distances among sample objects.

#### **BASIC STEPS OF CA:**

In CA, the sample objects are clustered on the basis of some characteristics. Therefore, to start with the analysis, a number of decisions must be made regarding the characteristics to be considered, the variables to be included in the analysis, the measurement of distance between objects and the criterion to group the objects.

The selection of variables for any CA is important, since the exclusion of important variables will be poor or misleading findings. For eg., if any marketing research the consumers are needed to be clustered, their tastes and habits and their economic capacities must be considered. Otherwise, the clustering of consumers will not be fruitful. The initial choice of variables determines the characteristics that can be used to identify subgroups.

After the selection of variables, the next important point to be considered is to measure the distance and similarity between objects. Two objects will be included in two separate groups, if their distance is maximum and they will be included in one group if they are close to each other. Therefore, one of the important steps in cluster analysis is to measure the distance among objects.

#### **SIMILARITY MEASURES:**

The measurement of similarity or distance is divided into two main parts. One of this is (a) distance – type measure, and another is (b) matching – type measure.



**Distance type similarity measures:**

**Euclidean distance:**

Let  $\underline{x} = [x_1, x_2, \dots, x_p]'$  and  $\underline{y} = [y_1, y_2, \dots, y_p]'$  are two p-dimensional observations or items or objects, then Euclidean distance between these two points is defined as

$$\begin{aligned} d(\underline{x}, \underline{y}) &= \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \dots + (x_p - y_p)^2} \\ &= \sqrt{(\underline{x} - \underline{y})' (\underline{x} - \underline{y})} \end{aligned}$$

**Statistical distance (Mahalonibis distance):**

Let  $\underline{x} = [x_1, x_2, \dots, x_p]'$  and  $\underline{y} = [y_1, y_2, \dots, y_p]'$  are two p-dimensional observations or items or objects, then Statistical distance between these two points is defined as

$$D(\underline{x}, \underline{y}) = \sqrt{(\underline{x} - \underline{y})' \mathbf{A} (\underline{x} - \underline{y})}$$

Ordinarily,  $\mathbf{A} = \mathbf{S}^{-1}$ , where  $\mathbf{S}$  contains the sample variances and covariances. However, without prior knowledge of the distinct groups, these sample quantities cannot be computed. For this reason Euclidean distance is often preferred for clustering.

**Minkowski distance:**

Let  $\underline{x} = [x_1, x_2, \dots, x_p]'$  and  $\underline{y} = [y_1, y_2, \dots, y_p]'$  are two p-dimensional observations or items or objects, then Minkowski distance between these two points is defined as

$$d(\underline{x}, \underline{y}) = \left( \sum_{i=1}^p |x_i - y_i|^m \right)^{\frac{1}{m}}$$

**Note 1:** When  $m=1$  Minkowski distance is called as city block distance

**Note 2:** When  $m=2$  Minkowski distance is reduced to Euclidean distance thus Euclidean distance is special case of Minkowski distance.

## Hierarchical clustering methods

### Agglomerative hierarchical clustering algorithm:

The following are the steps in the agglomerative hierarchical clustering algorithm for grouping  $N$  objects (items or variables).

1. Starts with  $N$  clusters, each containing a single entity and an  $N \times N$  symmetric matrix of distance (or similarities)  $D = \{d_{ik}\}$ .
2. Search the distance matrix for the nearest (most similar) pair of clusters. Let the distance between “most similar” clusters  $U$  and  $V$  be  $d_{uv}$ .
3. Merge cluster  $U$  and  $V$ . Label the newly formed cluster ( $UV$ ). Update the entries in the distance matrix by (a) deleting the rows and columns corresponding to cluster  $U$  and  $V$  and (b) adding a row and column giving the distances between cluster ( $UV$ ) and the remaining clusters.
4. Repeat steps 2 and 3 a total of  $N-1$  times. (All objects will be in a single cluster at termination of the algorithm.) Record the identity of cluster that are merged and the levels (distances or similarities) at which the mergers take place.

### Single Linkage (nearest neighbor) Clustering method:

In this method two objects with lowest distance are merged into a cluster. This cluster is the first one. At the second step, either a third object is merged to the first formed cluster, or the two closest unclustered objects are joined to form a second cluster. The decision rests on whether the distance from one of the unclustered objects to the first formed cluster is shorter than the distances between the two closest unclustered objects. The process continues until all objects belong to a single cluster. At any step, if a new cluster is formed and if the distance between the new cluster and the old one is shortest, then both the clusters are combined. **Since an object of a cluster cannot be split from the cluster, the distance between two clusters is equal to the distance between nearest objects.** After clustering, the objects in the clusters can be represented by a diagram. This diagram is known as Dendrogram.

The following are the steps in the single linkage clustering algorithm for grouping  $N$  objects (items or variables) .

1. Starts with  $N$  clusters, each containing a single entity and an  $N \times N$  symmetric matrix of distance (or similarities)  $D = \{d_{ik}\}$ .
2. Search the distance matrix for the smallest distance (other than zero) which gives the distance between nearest (most similar) pair of clusters. Let the distance between “most similar” clusters  $U$  and  $V$  be  $d_{UV}$ .
3. Merge clusters  $U$  and  $V$ . Label the newly formed cluster  $(UV)$ . Update the entries in the distance matrix by

(a) deleting the rows and columns corresponding to cluster  $U$  and  $V$

and (b) adding a row and column giving the distances between the cluster  $(UV)$  and the remaining clusters. Here, the distance between the cluster  $(UV)$  and any other cluster  $W$  is computed by

$$d_{(UV)W} = \min \{d_{UW}, d_{VW}\}$$

In general, if we take two clusters, one is with  $m$  objects and another is with  $n$  objects, then we will get ‘ $mn$ ’ pairs of objects, one is from the first cluster and another is from the second cluster. Thus, we get totally ‘ $mn$ ’ distances. **Now the distance between these two clusters is defined as the minimum (smallest) of these ‘ $mn$ ’ distances.**

4. Repeat steps 2 and 3 a total of  $N-1$  times.(All objects will be in a single cluster at termination of the algorithm,.) Record the identity of cluster that are merged and the levels(distances or similarities) at which the mergers take place.

### **Complete Linkage (farthest neighbor) Clustering method:**

In this method two objects with lowest distance are merged into a cluster. This cluster is the first one. At the second step, either a third object is merged to the first

formed cluster, or the two closest unclustered objects are joined to form a second cluster. The decision rests on whether the distance from one of the unclustered objects to the first formed cluster is shorter than the distances between the two closest unclustered objects. The process continues until all objects belong to a single cluster. At any step, if a new cluster is formed and if the distance between the new cluster and the old one is shortest, then both the clusters are combined. **The distance between two clusters is defined as the distance between most distant pair of objects.** After clustering, the objects in the clusters can be represented by a diagram. This diagram is known as Dendrogram.

The following are the steps in the complete linkage clustering algorithm for grouping  $N$  objects (items or variables) .

1. Starts with  $N$  clusters, each containing a single entity and an  $N \times N$  symmetric matrix of distance (or similarities)  $D = \{d_{ik}\}$ .
2. Search the distance matrix for the smallest distance (other than zero) which gives the distance between nearest (most similar) pair of clusters. Let the distance between “most similar” clusters  $U$  and  $V$  be  $d_{UV}$ .
3. Merge clusters  $U$  and  $V$ . Label the newly formed cluster  $(UV)$ . Update the entries in the distance matrix by
  - (a) deleting the rows and columns corresponding to cluster  $U$  and  $V$
  - and (b) adding a row and column giving the distances between the cluster  $(UV)$  and the remaining clusters. Here, the distance between the cluster  $(UV)$  and any other cluster  $W$  is computed by

$$d_{(UV)W} = \max \{d_{UW}, d_{VW}\}$$

In general, if we take two clusters, one is with  $m$  objects and another is with  $n$  objects, then we will get ‘ $mn$ ’ pairs of objects, one is from the first cluster and another is from the second cluster. Thus, we get totally ‘ $mn$ ’ distances. **Now the distance between these two clusters is defined as the maximum (largest) of these ‘ $mn$ ’ distances.**

4. Repeat steps 2 and 3 a total of  $N-1$  times.(All objects will be in a single cluster at termination of the algorithm,.) Record the identity of cluster that are merged and the levels(distances or similarities) at which the mergers take place.

### **Average Linkage Clustering method:**

In this method two objects with lowest distance are merged into a cluster. This cluster is the first one. At the second step, either a third object is merged to the first formed cluster, or the two closest unclustered objects are joined to form a second cluster. The decision rests on whether the distance from one of the unclustered objects to the first formed cluster is shorter than the distances between the two closest unclustered objects. The process continues until all objects belong to a single cluster. At any step, if a new cluster is formed and if the distance between the new cluster and the old one is shortest, then both the clusters are combined. **The distance between two clusters is defined as the average of the distances between all pairs of objects, one from the first cluster and another from the second cluster.** After clustering, the objects in the clusters can be represented by a diagram. This diagram is known as Dendrogram.

The following are the steps in the complete linkage clustering algorithm for grouping  $N$  objects (items or variables) .

1. Starts with  $N$  clusters, each containing a single entity and an  $N \times N$  symmetric matrix of distance (or similarities)  $D = \{d_{ik}\}$ .
2. Search the distance matrix for the smallest distance (other than zero) which gives the distance between nearest (most similar) pair of clusters. Let the distance between “most similar” clusters  $U$  and  $V$  be  $d_{UV}$ .
3. Merge clusters  $U$  and  $V$ . Label the newly formed cluster ( $UV$ ). Update the entries in the distance matrix by
  - (a) deleting the rows and columns corresponding to cluster  $U$  and  $V$
  - and (b) adding a row and column giving the distances between the cluster

( $UV$ ) and the remaining clusters. Here, the distance between the cluster ( $UV$ ) and any other cluster  $W$  is computed by

$$d_{(UV)W} = (d_{UW} + d_{VW}) / 2$$

In general, if we take two clusters, one is with  $m$  objects and another is with  $n$  objects, then we will get ' $mn$ ' pairs of objects, one is from the first cluster and another is from the second cluster. Thus, we get totally ' $mn$ ' distances. **Now the distance between these two clusters is defined as the average (mean) of these ' $mn$ ' distances.**

4. Repeat steps 2 and 3 a total of  $N-1$  times.(All objects will be in a single cluster at termination of the algorithm,.) Record the identity of cluster that are merged and the levels(distances or similarities) at which the mergers take place.

## **Nonhierarchical clustering methods**

Nonhierarchical clustering techniques are designed to groups items, rather than variables, into a collection of  $K$  clusters. The number of clusters,  $K$ , may either be specified in advance or determined as part of the clustering procedure. Because a matrix of distances (similarities) does not have to be determined and the basic data do not have to be stored during the computer run, nonhierarchical methods can be applied to much larger data sets than hierarchical techniques.

Nonhierarchical methods start from either (1) an initial partition of items into groups or (2) an initial set of seed points, which will form the nuclei of clusters. Good choice for starting configurations should be free of overt biases. One way to start is to randomly select seed points from among the items or to randomly partition the items into initial groups.

In this section we discuss one of the more popular nonhierarchical procedure, the  $K$ -means method.

### **K-means method:**

MacQueen suggests the term  $K$ -means for describing his algorithm that assigns each item to the cluster having the nearest centroid (mean). In its simplest version, the process is composed of these three steps.

1. Partition the items into  $K$  initial clusters.
2. Proceed through the list of items, assigning an item to the cluster whose centroid (mean) is nearest. (Distance is usually computed using Euclidean distance with

either standardized or unstandardized observations.) Recalculate the centroid for the cluster receiving the new item and for the cluster losing the item.

3. Repeat Step 2 until no more reassignments take place.

Rather than starting with a partition of all items into  $K$  preliminary groups in Step 1, we could specify  $K$  initial centroid (seed points) and then proceed to Step 2.

The final assignment of items to clusters will be, to some extent, dependent upon the initial partition or the initial selection of seed points. Experience suggests that most major changes in assignment occur with the first reallocation step.

## UNIT-1

1. (2001)

- (a) Define the generalized inverse of an  $m \times n$  matrix and derive four important properties.
- (b) How do you define the joint, conditional and marginal densities of vectors of random variables.

(OR)

- (c) Explain and distinguish uncorrelated, orthogonal and independent random vectors.
- (d) Define eigen values and eigen vectors. Establish their properties.

(20002)

- (a) Define the generalized inverse of an  $m \times n$  matrix and indicate a method of obtaining the generalized inverse of a square matrix whose determinant is zero.
- (b) For a given random vectors X and Y, prove that

$$(1) \text{cov}(AX, BY) = A[\text{cov}(X, Y)] B' \text{ and}$$

$$(2) \text{cov}(X, Y) = E[XY'] - E(X)[E(Y)]' .$$

(OR)

- (c) Show that every matrix satisfies its own characteristic function.
- (d) If X is a  $p \times 1$  vector of random variables with mean vector,  $\theta$  and Dispersion matrix,  $\Sigma$  show that  $E[X'AX] = \text{tr}(A\Sigma) + \theta' A\theta$  where A is a symmetric matrix of order  $p \times p$ .

(2003)

- (a) Define non-singular multivariate normal (MVN) distribution and obtain its characteristic function.



(b) Let  $X \sim N_p(\mu, \Sigma)$ . If  $X' = (X^{(1)'}, X^{(2)'})$ , where  $X^{(1)'}$  has  $k$ -components and  $X^{(2)'}$  has  $(p-k)$  components and further  $\mu$  and  $\Sigma$  are partitioned appropriately. Then find the conditional distribution of  $X^{(1)} / X^{(2)} = x^{(2)}$ .

(OR)

(c) Let  $X \sim N_p(\mu, \Sigma)$  and suppose that  $X$  is partitioned in the form

$X = (X^{(1)}, X^{(2)})'$ , where  $X^{(1)} : (q \times 1)$  and  $X^{(2)} : ((p-q) \times 1)$ . Then

$X^{(1)}$  and  $X^{(2)}$  are statistically independent if and only if  $\text{cov}(X^{(1)}, X^{(2)}) = 0$ .

(d) Obtain the maximum likelihood estimators of the parameters of a MVN distribution.

(2004)

(a) Define the multivariate normal distribution and obtain its characteristic function.

(b) Define :

(1) Marginal and

(2) Conditional distributions

Prove that the conditional distribution obtained from the multivariate normal distribution is normal.

(OR)

(c) Show that if  $X$  is  $N_p(\mu, \xi)$ , then the marginal distribution of any subset of  $X$  is multivariate normal with mean. Variances and covariances obtained by taking proper components of  $\mu$  and  $\xi$ .

(d) Let  $\mu$  and  $\xi$  be the parameters of a multivariate normal distribution.

Assuming a random sample from this distribution, obtain the maximum Likelihood. Estimate for  $\mu$  and  $\xi$ .

(2006)

- (a) Define the characteristic function of a p-dimensional random variables.  
Establish its properties.
- (b) Find the covariance matrix fore the two random variables  $X_1$  and  $X_2$   
when their joint probability function,  $P_{12}(x_1, x_2)$  is represented by the  
entries in the body of the following table:

$x_1 \backslash x_2$	0	1	$p_i(x_i)$
-1	0.24	0.06	0.3
0	0.16	0.14	0.3
1	0.40	0.00	0.4
$p_j(x_j)$	0.8	0.2	1

(OR)

- (c) Obtain the maximum likelihood estimators of the parameters of a  
multivariate normal distribution.
- (d) Let  $X$  be  $N_3(\mu, \xi)$  with

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

Are  $X_1$  and  $X_2$  independent? What about  $(X_1, X_2)$  and  $X_3$ ?

(Model paper)

- (a) Define the multivariate normal distribution and obtain its characteristic  
function.
- (b) Prove that the conditional distribution obtained from the multivariate  
normal distribution is normal.

(OR)

- (c) Show that if  $X$  is  $N_p(\mu, \Sigma)$  then the marginal distribution of any subset of  $X$  is multivariate normal with mean, variances and co-variances obtained by taking proper components of  $\mu$  and  $\Sigma$ .
- (d) Let  $\mu$  and  $\xi$  be the parameters of a multivariate normal distribution. Assuming a random sample from this distribution, obtain the maximum likelihood estimates for  $\mu$  and  $\Sigma$ .

(Model paper)

- (a) Define multivariate normal distribution and obtain its characteristic function.
- (b) Define marginal and conditional distributions. Prove that the conditional distribution obtained from the multivariate normal distribution is also multivariate normal distribution.

(2007)

- (a) Derive the density function of a p-variate normal distribution.
- (b) Let  $X$  has a trivariate normal distribution with  $E(X) = \tilde{0}$  and variance-

covariance matrix  $\Sigma = \begin{bmatrix} 1/2 & -1/2 & 1/2 \\ -1/2 & 1 & -1/2 \\ 1/2 & -1/2 & 1 \end{bmatrix}$ . Find the conditional

distribution of  $X_1$  given  $X_2 = x_2$  and  $X_3 = x_3$ .

(OR)

- (c) Obtain the characteristic function of multivariate normal distribution.
- (d) Derive the ML estimators of the mean vector and covariance matrix in the case of p-variate normal distribution.

2. (2001)

- (a)  $X$  is a p-variate random vector having the distribution  $N_p(\mu, \Sigma)$

where  $\mu$  is  $(p \times 1)$  and  $\Sigma$  is  $(p \times p)$ .  $X$  is partitioned as

$$X = \begin{pmatrix} X_{(1)} \\ X_{(2)} \end{pmatrix} \quad \text{where } X_{(1)} \text{ is of order } (q \times 1). \mu \text{ and } \Sigma \text{ are also}$$

conformably partitioned. Derive (1) the marginal distribution of  $X_{(1)}$

and (2) the conditional distribution of  $X_{(1)}$  given  $X_{(2)} = x_{(2)}$ .

- (b) Define the  $p$ -variate normal distribution with mean vector  $\mu$  and dispersion matrix  $\Sigma$ . Derive three important properties of the multivariate normal distribution.

(OR)

- (c) In the  $p$ -variate normal case, show that the sample mean vector and the sample covariance matrix are independently distributed.
- (d) Derive the M.L. estimators of the mean vector and the covariance matrix in the multivariate normal distribution  $N_p(\mu, \Sigma)$ .

(2002)

- (a) Define multivariate normal distribution and obtain its characteristic function.
- (b) If  $\mathbf{X}$  is  $N_p(\mu, \Sigma)$ , prove that  $X'\Sigma^{-1}X$  has a chi-square distribution with  $P$  degrees of freedom.

(OR)

- (c) State and prove a necessary and sufficient condition for  $X'AX$  and  $X'BX$  to be independently distributed, given that  $X$  is  $N_p(\mu, \Sigma)$ .
- (d) Derive the sampling distribution of the sample correlation coefficient ( $\gamma$ ) when the population correlation coefficient ( $\rho$ ) is zero.

(2003)

- (a) Define Hotelling's  $T^2$  and derive its probability density function. Explain the uses of this distribution.
- (b) Discuss a test procedure for testing the equality of mean vectors of two MVN populations, assuming equal dispersion matrices.

(OR)

- (c) stating the assumptions clearly, discuss the problem of comparing several multivariate normal population means.
- (e) Explain the following:
- (1) Wilks lambda criterion.
  - (2) Need for simultaneous confidence intervals.

(2004)

- (a) Derive the null distribution of Hotelling's  $T^2$  statistic.
- (b) Show that Hotelling's  $T^2$  statistic can be used to test equality of means of corresponding variables in two multivariate normal populations having the same variance-covariance matrix.

(OR)

- (c) Explain in detail the procedure of carrying out MANOVA of one way classification.
- (d) Define simultaneous confidence regions and illustrate broad steps to obtain the same.

(2006)

- (a) Derive the null distribution of Hotelling's  $T^2$  statistic.
- (b) Explain in detail the likelihood ratio principle.

(OR)

- (c) Explain in detail the one way classification of multivariate data. Write the MANOVA table.
- (d) Discuss the Behern-Fisher problem.

(Model paper)

- (a) Derive the null distribution of Hotelling's  $T^2$  statistic.
- (b) Show that Hotelling's  $T^2$  statistic can be used to test the equality of means of corresponding variables in two multivariate normal populations having the same variances-covariance matrix.

(OR)

- (c) Explain in detail the procedure of carrying out MANOVA of one way classification.
- (d) Define simultaneous confidence regions and illustrate broad steps to obtain the same.

(Model paper)

- (a) Derive the MLEs of the mean vector,  $\mu$  and the variance-covariance matrix,  $\Sigma$  based on a random sample of size  $n$  drawn from the normal population  $N_p(\mu, \Sigma)$ .
- (b) Show that the sample mean vector and sample variance-covariance matrix obtained based on a random sample of size  $n$  from a normal population  $N_p(\mu, \Sigma)$  are independently distributed and also mention their distributions.

(2007)

- (a) Define Hotelling's  $T^2$  statistic. What is its distribution? What are the applications of  $T^2$ ? Explain the relationship between Hotelling's  $T^2$  and Mahalanobis  $D^2$ .
- (b) Discuss a test procedure for testing the equality of mean vector of two multivariate normal populations having equal dispersion matrix.

(OR)

- (c) Explain MANOVA one way classification with an example.
- (d) Explain in detail the likelihood ratio principle.

3. (a) How do you test for the goodness of fit of a linear model?

(b) What is path analysis? What purposes are served by it?

(OR)

(c) What is multicollinearity? Discuss the methods of overcoming multicollinearity .

(d) Define autocorrelation . What are the consequences of the presence of autocorrelation?

(2002)

(a) State and prove Gauss-Markoff theorem for a standard general linear

model ,  $n \times 1 = n \times k \quad K \times 1 \quad n \times 1$ .

(b) Explain the concept of hetroskedasticity and describe a test procedure for testing the same.

(OR)

(c) Explain the concept of multicollinearity and indicate the consequences of OLS estimation of the parameters in the presences of collinearity.

(d) Explain the concept of auto-correlation and describe Durbin-Watson d-statistic for detecting the same.

(2003)

(a) What do you understand by “ dimension reduction”? describe a technique that is used for this purpose.

(b) Define:

(1) Canonical variables and

(2) Canonical correlation

Explain how do you estimate canonical correlation under the normality assumption.

(OR)

- (c) What are various properties of principal components? Show that principal Components are all uncorrelated.
- (d) Prove that canonical correlations are invariant under non-singular transformations and any function of the variance-covariance matrix that is invariant is a function of the canonical correlations.

(2004)

- (a) Describe the principal components analysis.
- (b) What are canonical variates and canonical correlations? How do you compute them?

(OR)

- (c) Define principal components and discuss their use in statistical analysis. If  $N_p(\mu, \xi)$  then explain how you would compute the various principals components.
- (d) Show that the canonical correlations are invariant under non-singular linear transformations of  $X^{(1)}, X^{(2)}$  variables of the form  $\underset{(p \times p)}{C} \underset{(p \times 1)}{X^{(1)}}$  and

$$\underset{(q \times q)}{D} \underset{(p \times 1)}{X^{(2)}}.$$

(2006)

- (a) What is principal component analysis? How are they useful?
- (b) Suppose the random variables  $X_1, X_2$  and  $X_3$  have the covariance matrix

$$\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

The corresponding eigen values are 5.8285, 2.00, 0.1716. obtain the three principal components and their variances.

(OR)

- (c) What are canonical variables? Explain how these are useful in the analysis



of multivariate data.

(d) For the following covariance matrix

$$\begin{array}{cc|cc} 100 & 0 & 0 & 0 \\ 0 & 1 & 0.95 & 0 \\ \hline 0 & 0.95 & 1 & 0 \\ 0 & 0 & 0 & 100 \end{array}$$

Obtain the first pair of canonical variables and the canonical correlation between them.

(Model paper)

- Describe the principle component analysis.
- What are canonical variates and canonical correlations? How do you compute them.

(OR)

- Define principle components and discuss their uses in statistical analysis .  
If  $N_p(\mu, \Sigma)$  then explain how you would compute the various principle components.
- Show that canonical correlations are invariant under non-singular transformations of  $X^{(1)}, X^{(2)}$  variables of the form  $\begin{matrix} C & X^{(1)} & D & X^{(2)} \\ (p \times p) & (p \times 1) & (q \times q) & (p \times 1) \end{matrix}$ .

(Model paper)

- Derive the null distribution of Hotelling's  $T^2$  statistic.
- If  $X$  is  $N_p(\mu, \Sigma)$ , prove that  $X'\Sigma^{-1}X$  has a chi-square distribution with  $p$  degrees of freedom .

(2007)

- What do you mean by dimension reduction? Describe a technique that is used for this purpose.
- Let the random variables  $X_1, X_2, X_3$  have the covariance matrix

$$\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

Obtain the three principal components.

(OR)

- (c) Define canonical variables and canonical correlations. Explain how do you estimate canonical correlation under the normality assumption.
- (d) Show that the canonical correlations are invariant under non-singular linear transformations of  $X^{(1)}, X^{(2)}$  variables of the form  $C_{p \times p}, X_{p \times 1}^{(1)}$  and  $D_{q \times q}, X_{q \times 1}^{(2)}$ .

4. (2001)

- (a) Explain principal component analysis. Give the uses of principal components in the factor analysis.
- (b) Describe cluster analysis and explain its uses.

(OR)

- (c) Explain canonical correlations. Show that the multiple correlation coefficient is a special case of canonical correlation coefficient. How will you obtain the latter?
- (d) What are canonical variables? Explain how these are useful in the analysis of multivariate data.

(2002)

- (a) Define principal components and explain the procedure of obtaining them.
- (b) What are canonical variates and canonical correlations? How do you compute them?

(OR)

- (c) distinguish between factor analysis and principal component analysis.

- (d) Given two sets of variates  $X^{(1)}$  &  $X^{(2)}$ , show that the canonical correlation are invariant under non-singular linear transformations of the form

$$Y^{(1)} = AX^{(1)}, Y^{(2)} = BX^{(2)} \text{ where A and B are non-singular matrices.}$$

(2003)

- (a) Let  $f_i(x) \sim N_p(\mu_i, \Sigma), i = 1, 2, \dots$ . Derive the optimal classification rule.

Also obtain the probabilities of misclassification.

- (b) Describe the following hierarchical clustering methods:

(1) single linkage method.(SLINK)

(2) complete linkage method (CLINK)

(OR)

- (c) Discuss Fisher's method for discrimination among several populations.

- (d) Explain the concept of clustering. What is a Dendogram? The following of sample correlations for five stocks:

$$D = \begin{bmatrix} 1 & & & & \\ .58 & 1 & & & \\ .51 & .60 & 1 & & \\ .39 & .39 & .44 & 1 & \\ .46 & .32 & .43 & .52 & 1 \end{bmatrix}$$

Treating the sample correlation coefficients as similarity measures, cluster analyze the stocks using the nearest neighbour method and also draw the dendogram.

(2004)

- (a) Explain the problem of classification and derive fishers linear discriminant function.

- (b) Explain the importance of cluster analysis and describe of hierarchial clustering.

(OR)

- (c) Explain the problem of classification into one of the two known multivariate normal populations.
- (d) Describe the method of classification of an individual into one of several p-variate normal populations having a common dispersion matrix.  $\xi$  where all the parameters are known.

(2006)

- (a) What is meant by discriminant analysis? Establish the relationship between Fishers discriminant function and Mahalanobi's  $D^2$  statistic.
- (b) Explain the importance of cluster analysis and describe the Hierarchical method of clustering.

(OR)

- (c) Explain Fisher's method of discriminating between two multivariate populations.
- (d) A researcher has enough data available to estimate the density functions  $f_1(\mathbf{x})$  and  $f_2(\mathbf{x})$  associated with populations  $\pi_1$  and  $\pi_2$  respectively. Suppose  $C(2/1)=5$  units,  $C(1/2)=10$  units and it is known that about 20% of all objects (for which the measurement  $\mathbf{x}$  can be recorded) belonging  $\pi_2$ .

Suppose the density function evaluated at a new observation  $\mathbf{x}_0$  give  $f_1(\mathbf{x}_0)=0.3$  and  $f_2(\mathbf{x}_0)=0.4$ . Do you classify the new observation as  $\pi_1$  and  $\pi_2$ .

(Model paper)

- (a) Explain the problem of classification and derive Fishers linear discriminant function.
- (b) Let  $f_i(X) \sim N_p(\mu_i, \Sigma)$ ,  $i=1,2,\dots$ . Derive the optimal classification rule also obtain the probabilities of misclassification.
- (c) Explain the problem of classification into one of the known multivariate normal populations.

(d) Describe the method of classification of an individual in to one of several p-variate populations having a common dispersion matrix  $\Sigma$  where all the parameters are known .

(Model paper)

(a) Stating the assumptions clearly , discuss the problem of comparing several multivariate normal population means.

(b) Explain the following:

(1) Likelihood ratio test

(2) Need for simultaneous confidence intervals.

(2007)

(a) Explain the problem of classification and derive Fisher's linear discriminant function.

(b) Explain the problem of classification into one of the two known multivariate normal populations.

(OR)

(c) What are the standards of good classification? Describe a linear discriminant function as a tool for classification.

(d) Describe the method of classification of an individual into one of several p-variate normal populations having a common dispersion matrix when the parameters are known.

5. (2001)

(a) What is the problem of classification and what are the standards of good classification?

(b) Discuss the problem of classification into one of two known multivariate normal populations.

(OR)

(c) Explain the concept of discriminant function and discuss,with an example , its use in discriminatory analysis.

(d) When do you apply cluster analysis? Discuss single linkage method with an example.

(2002)

- (a) Explain the problem of classification into one of the two known multivariate normal populations.
- (b) Distinguish between cluster analysis and discriminant analysis.

(OR)

- (c) Describe the method of classification of an individual into one of several  $p$ -variate normal populations having a common dispersion matrix,  $\Sigma$  where all the parameters are known.
- (d) Explain the importance of cluster analysis and describe the method of Hierarchical clustering.