

## Unit II (Multivariate Analysis)

### The generalised $T^2$ Statistic

In the univariate case, a random sample of size  $n$  observations  $x_1, x_2, \dots, x_n$  are drawn from a normal population with mean  $\mu$  and unknown variance  $\sigma^2$ , the test statistic under the hypothesis

$$H_0: \mu = \mu_0 \quad \text{is} \quad t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t(\alpha) \text{ with } (n-1) \text{ d.f.}$$

$$\text{where } \bar{x} = \frac{1}{n} \sum x_i \quad \text{and} \quad s^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2 \quad \rightarrow \textcircled{1}$$

The decision rule is

accept  $H_0$  if  $|t| \leq t(\alpha)$  with  $(n-1)$  d.f.

reject  $H_0$  if  $|t| > t(\alpha)$  with  $(n-1)$  d.f.

The multivariate analog of the square of 't' given  $\textcircled{1}$  is

$$T^2 = N (\bar{x} - \mu_0)' S^{-1} (\bar{x} - \mu_0)$$

where,  $\bar{x}$  is a mean vector of the sample of size  $N$  and  $S$  is the Sample Co-variance matrix,

Hotelling proposed the  $T^2$  statistic for two samples and derived the distribution when  $\mu$  is the population mean.

## Derivation of the generalised Hotelling $T^2$ Statistic

The likelihood ratio test of the hypothesis

$H_0: \mu = \mu_0$  on the basis of the sample from  $N(\mu, \Sigma)$  is based on the  $T^2$  statistic,

$$T^2 = N(\bar{x} - \mu)' S^{-1}(\bar{x} - \mu)$$

Suppose we have  $N$  observations  $x_1, x_2, \dots, x_N$  ( $N > p$ ),

the likelihood function is

$$L(\mu, \Sigma) = (2\pi)^{-\frac{NP}{2}} |\Sigma|^{-N/2} \exp \left\{ -\frac{1}{2} \sum_{\alpha=1}^N (x_\alpha - \mu)' \Sigma^{-1} (x_\alpha - \mu) \right\}$$

The likelihood ratio criterion for testing the hypothesis

$H_0: \mu = \mu_0$  is

$$\lambda = \frac{\max_{\Sigma} L(\mu_0, \Sigma)}{\max_{\mu, \Sigma} L(\mu, \Sigma)} \rightarrow \textcircled{2}$$

when the parameters are unrestricted, the maximum occurs when  $\mu, \Sigma$  are defined by the maximum likelihood estimators of  $\mu$  and  $\Sigma$ .

$$\text{ie } \hat{\mu}_{ML} = \bar{x} \text{ and } \hat{\Sigma}_{ML} = \frac{1}{N} \sum_{\alpha=1}^N (x_\alpha - \bar{x})' (x_\alpha - \bar{x})'$$

when  $\mu = \mu_0$  the likelihood function is maximized at

$$\hat{\Sigma}_{\omega} = \frac{1}{N} \sum_{\alpha=1}^N (x_\alpha - \mu_0)' (x_\alpha - \mu_0)' \rightarrow \textcircled{3}$$

under  $\omega$ , the likelihood function reduced to

$$\begin{aligned} \max_{\Sigma} L(\mu_0, \Sigma) &= (2\pi)^{-\frac{NP}{2}} |\Sigma_{\omega}|^{-N/2} \exp \left\{ -\frac{1}{2} \sum_{\alpha=1}^N (x_\alpha - \mu_0)' \Sigma_{\omega}^{-1} (x_\alpha - \mu_0) \right\} \\ &= (2\pi)^{-\frac{NP}{2}} |\Sigma_{\omega}|^{-N/2} \exp \left\{ -\frac{1}{2} \text{tr} \left[ \sum_{\alpha=1}^N (x_\alpha - \mu_0)' \Sigma_{\omega}^{-1} (x_\alpha - \mu_0) \right] \right\} \\ &= \text{exp} \left\{ -\frac{1}{2} \sum_{\alpha=1}^N \text{tr} \left[ (x_\alpha - \mu_0)' \Sigma_{\omega}^{-1} (x_\alpha - \mu_0) \right] \right\} \\ &= \text{exp} \left\{ -\frac{1}{2} \sum_{\alpha=1}^N \text{tr} \left[ \Sigma_{\omega}^{-1} (x_\alpha - \mu_0)' (x_\alpha - \mu_0) \right] \right\} \end{aligned}$$

$$= \exp \left\{ -\frac{1}{2} \ln \left| \sum_{\alpha=1}^N \left[ (x_{\alpha} - \mu_0)' (x_{\alpha} - \mu_0) \right] \right| \right\}$$

$$= \exp \left\{ -\frac{1}{2} \ln \left| \sum_{\alpha=1}^N \left[ N \hat{\Sigma}_{\omega} \right] \right| \right\}$$

$$= \exp \left\{ -\frac{1}{2} \ln |N I| \right\}$$

$$= \exp \left\{ -\frac{1}{2} \ln |N P| \right\}$$

$$\therefore \text{Max}_{\Sigma} L(\mu_0, \Sigma) = (2\pi)^{-N/2} \left| \hat{\Sigma}_{\omega} \right|^{-N/2} \exp \left( -\frac{NP}{2} \right) \rightarrow (4)$$

Similarly,

$$\text{Max}_{\mu, \Sigma} L(\mu, \Sigma) = (2\pi)^{-N/2} \left| \hat{\Sigma}_{\Omega} \right|^{-N/2} \exp \left( -\frac{NP}{2} \right) \rightarrow (5)$$

Thus the likelihood ratio criterion is

$$\lambda = \frac{\text{Max}_{\Sigma} L(\mu_0, \Sigma)}{\text{Max}_{\mu, \Sigma} L(\mu, \Sigma)}$$

$$= \frac{\left| \hat{\Sigma}_{\Omega} \right|^{N/2}}{\left| \hat{\Sigma}_{\omega} \right|^{N/2}}$$

$$= \frac{\frac{1}{N} \left| \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})' \right|^{N/2}}{\frac{1}{N} \left| \sum_{\alpha=1}^N (x_{\alpha} - \mu_0)(x_{\alpha} - \mu_0)' \right|^{N/2}}$$

$$\therefore \lambda = \frac{|A|^{N/2}}{|A + N(\bar{x} - \mu_0)(\bar{x} - \mu_0)'|^{N/2}}$$

$$\lambda^{2/N} = \frac{|A|}{|A + N(\bar{x} - \mu_0)(\bar{x} - \mu_0)'|}$$

$$= \frac{|(N-1)S|}{|(N-1)S + N(\bar{x} - \mu_0)(\bar{x} - \mu_0)'|} \quad \text{where } S = \frac{1}{N-1}(A)$$

$$= \frac{|(N-1)S|}{|(N-1)S| \left| 1 + \frac{N(\bar{x} - \mu_0)(\bar{x} - \mu_0)'}{(N-1)S} \right|}$$

$$= \frac{|(N-1)S|}{|(N-1)S| \left| 1 + \frac{N(\bar{x} - \mu_0) \bar{S}'(\bar{x} - \mu_0)'}{N-1} \right|}$$

$$\lambda^{\frac{2}{N}} = \frac{1}{1 + \frac{T^2}{N-1}} \rightarrow \textcircled{6}$$

$$\text{where } T^2 = N(\bar{x} - \mu_0) \bar{S}'(\bar{x} - \mu_0)'$$

The likelihood ratio test is defined by the critical region  $\lambda \leq \lambda_0 \rightarrow \textcircled{7}$

where  $\lambda_0$  is chosen so that the Prob. of  $\textcircled{7}$  when the null hypothesis  $H_0$  is true.

$$\therefore \textcircled{5} \Rightarrow \lambda^{\frac{2}{N}} \leq \lambda_0^{\frac{2}{N}}$$

$$\Rightarrow \lambda^{-\frac{2}{N}} \geq \lambda_0^{-\frac{2}{N}}$$

$$\therefore \cancel{(N-1)} \lambda^{-\frac{2}{N}} - 1 \geq \lambda_0^{-\frac{2}{N}} - 1$$

$$(N-1)(\lambda^{\frac{2}{N}} - 1) \geq (N-1)(\lambda_0^{-\frac{2}{N}} - 1)$$

$$T^2 \geq T_0^2$$

$$\therefore T_0^2 = (N-1)(\lambda_0^{-\frac{2}{N}} - 1)$$

# The distribution of Hotelling's $T^2$

W.K.T  $T^2 = N(\bar{x} - \mu_0)' S^{-1}(\bar{x} - \mu_0) \rightarrow \textcircled{1}$

let  $T^2 = Y' S^{-1} Y \rightarrow \textcircled{2}$

where  $Y \sim N(\mu, \Sigma)$  and  $nS$  is distributed independently

as  $\sum_{\alpha=1}^n z_{\alpha} z_{\alpha}'$  with  $z_1, z_2, z_3 \dots z_n$  are independent, each

with distribution  $N(0, \Sigma)$ . The  $T^2$  defined with

$Y = \sqrt{N}(\bar{x} - \mu_0)$  and  $\mu = \sqrt{N}(\mu - \mu_0)$  and  $n = N-1$ .

$T^2 = N(\bar{x} - \mu_0)' S^{-1}(\bar{x} - \mu_0)$

$T^2 = \sqrt{N}(\bar{x} - \mu_0)' S^{-1} \sqrt{N}(\bar{x} - \mu_0) \rightarrow \textcircled{3}$

From  $\textcircled{2}$  and  $\textcircled{3}$

$Y = \sqrt{N}(\bar{x} - \mu_0)$

$E(Y) = E\{\sqrt{N}(\bar{x} - \mu_0)\}$

$E(Y) = \sqrt{N}\{E(\bar{x} - \mu_0)\}$

$\mu = \sqrt{N}(\mu - \mu_0)$  and  $n = N-1$

Let  $D$  be a non singular matrix such that

$D' \Sigma D = I$  and define

$\left. \begin{aligned} Y^* &= DY \\ S^* &= DS D' \\ Y^* &= DY \end{aligned} \right\} \rightarrow \textcircled{4}$

Now  $Y^* = DY$

$\Rightarrow DY = Y^*$

$Y = D^{-1} Y^*$

$Y' = (D^{-1} Y^*)'$

$Y' = Y^{*'} (D^{-1})'$

$S^* = DS D'$

$DS D' = S^*$

$S = D^{-1} S^* (D^{-1})'$

$= (D^{-1}) [S^*]^{-1} (D^{-1})'$

$S = (D^{-1} S^* (D^{-1})')^{-1}$

$S^{-1} = D' S^* D$



Then from (2)

$$T^2 = Y' S^{-1} Y$$

$$= (Y^*)' (D^{-1})' (D^{-1} S^* D^{-1}) (D^{-1} Y^*)$$

$$= (Y^*)' (D^{-1})' (D^*)^{-1} Y^*$$

$$T^2 = (Y^*)' (S^*)^{-1} Y^*$$

where  $Y^* \sim N(Y^*, I)$  since  $Y \sim N(Y, \Sigma)$   
 $Y^* = DY \sim N(DY, DSD')$   
 $\therefore Y^* \sim N(Y^*, I)$

$$nS = \sum_{\alpha=1}^n z_{\alpha} z_{\alpha}' \quad \text{where } z_{\alpha} \sim N(0, \Sigma)$$

$$nS^* = n(DSD')$$

{∴ from (4)}

$$= D nS D'$$

$$= D \left( \sum_{\alpha=1}^n z_{\alpha} z_{\alpha}' \right) D'$$

$$= \sum_{\alpha=1}^n (Dz_{\alpha}) (z_{\alpha}' D')$$

$$= \sum_{\alpha=1}^n (Dz_{\alpha}) (Dz_{\alpha})'$$

$$\text{①} \leftarrow \sum_{\alpha=1}^n z_{\alpha}^* (z_{\alpha}^*)'$$

where  $Dz_{\alpha} = z_{\alpha}^*$

where  $z_{\alpha}^*$  are independent each with distribution  $N(0, DSD')$

$$\therefore z_{\alpha} \sim N(0, \Sigma)$$

$$z_{\alpha}^* = Dz_{\alpha} \sim N(0, DSD')$$

$$\sim N(0, I)$$

Let the first row of a  $p \times p$  orthogonal matrix

$Q$  be defined by

$$q_{1i} = \frac{Y_i^*}{\sqrt{Y^{*1} Y^*}}; \quad i=1, 2, \dots, p \quad \rightarrow (5)$$

This is permissible since  $\sum_i q_{1i}^2 = 1$

now let  $U = Q Y^*$  and  $B = Q^* S^* Q^1$

$$Q \text{ was defined by } U_1 = \sum q_{1i} Y_i^* = \frac{Y_i^{*1} Y_i^*}{\sqrt{Y^{*1} Y^*}} \\ = \sqrt{Y^{*1} Y^*}$$

$$\text{and } U_j = \sum q_{ji} Y_i^*$$

$$= \sum q_{ji} q_{1i} \sqrt{Y^{*1} Y^*} \quad \text{using (5)}$$

$$= \sqrt{Y^{*1} Y^*} \sum q_{ji} q_{1i} \quad \text{⊗}$$

$$U_j = 0 \quad \text{for } j \neq 1$$

$$\text{since } q_{ji} q_{1i} = 0$$

$$\text{Then } \frac{T^2}{n} = \frac{Y^{*1} S^{*-1} Y^*}{n}$$

$$\text{Substitute } U = Q Y^*$$

$$\Rightarrow Y^* = Q^{-1} U$$

$$\therefore \frac{T^2}{n} = \frac{(Q^{-1} U)^1 (S^*)^{-1} (Q^{-1} U)}{n}$$

$$= \frac{U' (Q^{-1})' (S^*)^{-1} (Q^{-1} U)}{n}$$

$$= \frac{U' (Q')' (S^*)^{-1} Q^{-1} U}{n} \quad (\because Q^{-1} = Q' \text{ if } Q \text{ is orthogonal})$$

$$\frac{\Gamma^2}{n} = \frac{U' Q (S^*)^{-1} Q' U}{n} \rightarrow \textcircled{b}$$

and  $B = Q_n S^* Q'$

$$B^{-1} = (Q_n S^* Q')^{-1}$$

$$= (Q')^{-1} (S^*)^{-1} Q^{-1}$$

$$B^{-1} = \frac{(Q')^{-1} (S^*)^{-1} Q^{-1}}{n}$$

$$\because Q' = Q^{-1}$$

$$\boxed{B^{-1} = \frac{Q (S^*)^{-1} Q'}{n}}$$

$$\therefore \textcircled{b} \Rightarrow \frac{\Gamma^2}{n} = U' B^{-1} U$$

Since  $U_i = \sqrt{Y_i^{*1} Y_i^{*2}}$  and  $U_j = 0$  for  $j \neq i$

The vector  $U = \begin{bmatrix} U_i \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$



$$\frac{1}{2} T^2 = U B U$$

$$= (U, 0, 0, \dots, 0)$$

$$\begin{bmatrix} b^{11} & b^{12} & \dots & b^{1p} \\ b^{21} & b^{22} & \dots & b^{2p} \\ \vdots & \vdots & \ddots & \vdots \\ b^{p1} & b^{p2} & \dots & b^{pp} \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\frac{1}{2} T^2 = U_1^2 b^{11}$$

$$b^{ij} = \frac{-1}{B}$$

$$T \approx \frac{1}{T} \approx \frac{1}{T}$$

## Applications of $T^2$ in tests on mean vector

(i) Tests the hypothesis that the mean vector is a given vector

$T^2$  test statistic is used to testing the hypothesis that the mean vector  $\mu = \mu_0$  when the dispersion matrix  $\Sigma$  is known.

The likelihood ratio test of the hypothesis

$H_0: \mu = \mu_0$  on the basis of the sample  $N$  from  $N(\mu, \Sigma)$  is equivalent to  $T^2 \geq T_0^2$  with the level of significance  $\alpha$ .  
the hypothesis is rejected if

$$T^2 \geq \frac{(N-1)P}{N-P} F_{(P, N-P)} \text{ d.f. } (\alpha)$$

(ii) Finding the confidence region for the mean vector

Let  $X_\alpha = (X_{\alpha 1}, X_{\alpha 2}, \dots, X_{\alpha p})'$   $\alpha = 1, 2, \dots, N$   
be a sample of size  $N$  from  $p$ -variate normal distn. with unknown mean  $\mu$  and unknown positive definite Covariance matrix  $\Sigma$ .

$$\text{WKT } \bar{X} = \frac{1}{N} \sum_{\alpha=1}^N X_\alpha$$

$$\text{and } S = \frac{1}{N-1} \sum_{\alpha=1}^N (X_\alpha - \bar{X})(X_\alpha - \bar{X})'$$

$N(\bar{X} - \mu)' S^{-1} (\bar{X} - \mu)$  is distributed as Hotelling  $T^2$  distribution with  $(N-1)$  d.f.

let  $T_0^2(\alpha)$ ;  $0 < \alpha < 1$  such that  $\Pr(T^2 \geq T_0^2) = \alpha$ ,  
then the probabilities of drawing a sample of  $X_\alpha$  with mean vector  $\bar{X}$  and sample covariance matrix  $S$  such that

$$N(\bar{X} - \mu_0)' S^{-1} (\bar{X} - \mu_0) \leq T_0^2(\alpha) \rightarrow \textcircled{1}$$

Thus if we compute ① for a particular sample, we have confidence  $(1-\alpha)$  the equation ① is true Statement concerning  $\mu$ .

### (ii) Test for the equality of two mean vectors

The  $T^2$  statistic is used to test the null hypothesis that the mean vector of one normal population is equal to the mean vector of the other normal population, where the Covariance matrices are assumed to be equal but unknown.

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

Suppose  $Y_1^{(1)}, Y_2^{(1)}, \dots, Y_{N_1}^{(1)}$  be a sample from  $N_p(\mu^{(1)}, \Sigma)$

$$\text{The sample mean } \bar{Y}^{(1)} = \frac{1}{N_1} \sum_{\alpha=1}^{N_1} Y_{\alpha}^{(1)}$$

$$\text{then } \bar{Y}^{(1)} \sim N\left(\mu^{(1)}, \frac{\Sigma}{N_1}\right)$$

$$\text{Similarly } \bar{Y}^{(2)} \sim N\left(\mu^{(2)}, \frac{\Sigma}{N_2}\right)$$

$$\therefore (\bar{Y}^{(1)} - \bar{Y}^{(2)}) \sim N\left(\mu^{(1)} - \mu^{(2)}, \frac{\Sigma}{N_1} + \frac{\Sigma}{N_2}\right)$$

$$\sim N\left(0, \left(\frac{N_1 + N_2}{N_1 N_2}\right) \Sigma\right) \quad \left\{ \because \text{by } H_0 \right.$$

$$\sqrt{\frac{N_1 N_2}{N_1 + N_2}} (\bar{Y}^{(1)} - \bar{Y}^{(2)}) \sim N(0, \Sigma)$$

If we let

$$S = \frac{1}{N_1 + N_2 - 2} \left\{ \sum_{\alpha=1}^{N_1} (Y_{\alpha}^{(1)} - \bar{Y}^{(1)}) (Y_{\alpha}^{(1)} - \bar{Y}^{(1)})' + \sum_{\alpha=1}^{N_2} (Y_{\alpha}^{(2)} - \bar{Y}^{(2)}) (Y_{\alpha}^{(2)} - \bar{Y}^{(2)})' \right\}$$

Then  $(N_1 + N_2) S$  is distributed as

$$\sum_{\alpha=1}^{N_1 + N_2 - 2} Z_{\alpha} Z_{\alpha}' \quad \text{where } Z_{\alpha} \sim N(0, \Sigma)$$

Thus  $T^2 = \frac{N_1 N_2}{N_1 + N_2} (\bar{y}^{(1)} - \bar{y}^{(2)})' S^{-1} (\bar{y}^{(1)} - \bar{y}^{(2)})$  is

distributed as  $T^2$  with  $(N_1 + N_2 - 2)$  d.f.

The critical region is

$$T^2 > \frac{(N_1 + N_2 - 2) P}{N_1 + N_2 - P - 1} F_{(P, N_1 + N_2 - P - 1)}(\alpha)$$

(iv) Two Sample Problem with unequal Covariance matrices

Let  $\{X_\alpha^{(i)}\}$   $\alpha=1, 2, \dots, N_i$  be samples from  $N(\mu^{(i)}, \Sigma_i)$

we wish to test the hypothesis that  $i=1, 2$

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

The mean of the first sample is normally distributed with expected value  $E(\bar{x}^{(1)}) = \mu^{(1)}$  and Covariance matrix

$$\Sigma_1 = \frac{1}{N_1} \sum_{\alpha=1}^{N_1} (X_\alpha^{(1)} - \bar{x}^{(1)}) (X_\alpha^{(1)} - \bar{x}^{(1)})'$$

Similarly  $E(\bar{x}^{(2)}) = \mu^{(2)}$

$$\text{and } \Sigma_2 = \frac{1}{N_2} \sum_{\alpha=1}^{N_2} (X_\alpha^{(2)} - \bar{x}^{(2)}) (X_\alpha^{(2)} - \bar{x}^{(2)})'$$

Thus  $(\bar{x}^{(1)} - \bar{x}^{(2)})$  has mean  $(\mu^{(1)} - \mu^{(2)})$  and

$$\text{Covariance matrix } \frac{\Sigma_1}{N_1} + \frac{\Sigma_2}{N_2}$$

If  $N_1 = N_2 = N$  (say)

$$\text{and let } Y_\alpha = X_\alpha^{(1)} - X_\alpha^{(2)}$$

$$\begin{aligned} \bar{Y} &= \frac{1}{N} \sum Y_\alpha \\ &= \bar{x}^{(1)} - \bar{x}^{(2)} \end{aligned}$$

$$\therefore \bar{Y} \sim N((\mu^{(1)} - \mu^{(2)}), \frac{1}{N} (\Sigma_1 + \Sigma_2))$$

$$S = \frac{1}{N-1} \sum_{\alpha=1}^N (Y_{\alpha} - \bar{Y})(Y_{\alpha} - \bar{Y})'$$

$\therefore$  our test statistic  $T^2 = N\bar{Y}' S^{-1} \bar{Y}$  is distributed  $T^2$  distribution with  $(N-1)$  d.f.