

Chapter 6

Comparing Mean Vectors from Two Populations

Analysts may have data on samples from different populations and may be interested in answering questions as to whether the populations under consideration have different means (or mean vectors in the multivariate case)

Two independent mean differences

① When the two populations have equal variances

Ⓐ Univariate case: (small n_1, n_2)

$X_{11}, X_{12}, X_{13}, \dots, X_{1n}$ is a random sample of size n_1 from $N(\mu_1, \sigma_1^2)$
 $X_{21}, X_{22}, X_{23}, \dots, X_{2n}$ is a random sample of size n_2 from $N(\mu_2, \sigma_2^2)$

If $\sigma_1^2 = \sigma_2^2 = \sigma^2$ (Perform Levene's ^{F test} to check this)

Test statistic for the test $H_0: \mu_1 - \mu_2 = \delta_0$ vs $H_a: \mu_1 - \mu_2 \neq \delta_0$ is given as

$$t^* = \frac{\bar{X}_1 - \bar{X}_2 - \delta_0}{\sqrt{s_{\text{pooled}}^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim t(n_1 + n_2 - 2)$$

$$\text{where } S_{\text{pooled}}^2 = \frac{S_1^2(n_1-1) + S_2^2(n_2-1)}{n_1+n_2-2}$$

with

$$\bar{x}_i = \frac{\sum_{j=1}^{n_i} x_{ij}}{n_i} \quad \text{and} \quad S_i^2 = \frac{1}{n_i-1} \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2$$

The $100(1-\alpha)\%$ confidence interval for $\mu_1 - \mu_2$ is given by

$$\bar{x}_1 - \bar{x}_2 \pm t(\alpha/2, n_1+n_2-2) \sqrt{S_{\text{pooled}}^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

(B) Multivariate case: Assume $\underline{x}_{11}, \underline{x}_{12}, \dots, \underline{x}_{1n} \sim N_p(\mu_1, \Sigma_1)$
(small n_1, n_2) $\underline{x}_{21}, \underline{x}_{22}, \dots, \underline{x}_{2n} \sim N_p(\mu_2, \Sigma_2)$

To test equality of variance covariance matrices use the Box M test which test $H_0: \Sigma_1 = \Sigma_2$

Will be demonstrated in R.

If the test concludes that $\Sigma_1 = \Sigma_2 = \Sigma$, we have

$\sum_{j=1}^{n_1} (\underline{x}_{1j} - \bar{\underline{x}}_1)(\underline{x}_{1j} - \bar{\underline{x}}_1)^T$ as an estimate of $(n_1-1)\Sigma$ and

$\sum_{j=1}^{n_2} (\underline{x}_{2j} - \bar{\underline{x}}_2)(\underline{x}_{2j} - \bar{\underline{x}}_2)^T$ as an estimate of $(n_2-1)\Sigma$

so that the pooled covariance is ^{matrix} calculated as

$$S_p = S_{\text{pooled}} = \frac{\sum_{j=1}^{n_1} (\tilde{x}_{1j} - \bar{\tilde{x}}_1)(\tilde{x}_{1j} - \bar{\tilde{x}}_1)^T + \sum_{j=1}^{n_2} (\tilde{x}_{2j} - \bar{\tilde{x}}_2)(\tilde{x}_{2j} - \bar{\tilde{x}}_2)^T}{n_1 + n_2 - 2}$$

$$= \frac{n_1 - 1}{n_1 + n_2 - 2} S_1 + \frac{n_2 - 1}{n_1 + n_2 - 2} S_2$$

where $S_1 = \frac{1}{n_1 - 1} \sum_{j=1}^{n_1} (\tilde{x}_{1j} - \bar{x}_1)(\tilde{x}_{1j} - \bar{x}_1)^T$ and
 $S_2 = \frac{1}{n_2 - 1} \sum_{j=1}^{n_2} (\tilde{x}_{2j} - \bar{x}_2)(\tilde{x}_{2j} - \bar{x}_2)^T$

To test

$$H_0: \tilde{\mu}_1 - \tilde{\mu}_2 = \delta_0 \quad \text{vs} \quad H_a: \tilde{\mu}_1 - \tilde{\mu}_2 \neq \delta_0$$

Test statistic :

$$\bar{T}^2 = (\bar{\tilde{x}}_1 - \bar{\tilde{x}}_2 - (\tilde{\mu}_1 - \tilde{\mu}_2)) \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_{\text{pooled}} \right]^{-1} (\bar{\tilde{x}}_1 - \bar{\tilde{x}}_2 - (\tilde{\mu}_1 - \tilde{\mu}_2))$$

$$(\bar{T}^*)^2 = \frac{(n_1 + n_2 - p - 1)}{(n_1 + n_2 - 2) p} \bar{T}^2 \sim F(\alpha, p, n_1 + n_2 - p - 1)$$

Reject H_0 if $(\bar{T}^*)^2 > F(\alpha, p, n_1 + n_2 - p - 1)$

The $100(1-\alpha)\%$ confidence region for $\tilde{\mu}_1 - \tilde{\mu}_2$ is given by all $\tilde{\mu}_1 - \tilde{\mu}_2$ satisfying:

$$(\bar{\tilde{x}}_1 - \bar{\tilde{x}}_2 - (\tilde{\mu}_1 - \tilde{\mu}_2))^T \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_{\text{pooled}} \right]^{-1} ((\bar{\tilde{x}}_1 - \bar{\tilde{x}}_2) - (\tilde{\mu}_1 - \tilde{\mu}_2)) \leq \frac{(n_1 + n_2 - 2) p}{(n_1 + n_2 - p - 1)} F(\alpha, p, n_1 + n_2 - p - 1)$$

Note that this region is an ellipsoid centered at the observed $\bar{\bar{x}}_1 - \bar{\bar{x}}_2$ and whose axes are determined by the eigen values and eigen vectors of S_{pooled} .

Refer to example 6.3 (Problem 4 in example document)

③ Simultaneous confidence intervals for $\bar{\mu}_1 - \bar{\mu}_2$

We construct these confidence intervals by considering all possible linear combinations of the differences in mean vectors.

The $100(1-\alpha)\%$ simultaneous confidence interval for $\bar{a}^T(\bar{\mu}_1 - \bar{\mu}_2)$ for all a is given as :

$$\bar{a}^T(\bar{\bar{x}}_1 - \bar{\bar{x}}_2) \pm c \sqrt{\bar{a}^T \left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_{\text{pooled}} a}$$

where $c^2 = \frac{(n_1+n_2-2)p}{n_1+n_2-p-1} F(\alpha, p, n_1+n_2-p-1)$

$\overline{T^2}$ simultaneous

In particular, we have that the $100(1-\alpha)\%$ confidence interval for $\mu_{1,i} - \mu_{2,i}$ is given by :

$$\bar{x}_{1,i} - \bar{x}_{2,i} \pm c \sqrt{S_{ii, \text{pooled}} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} \quad \text{for } i=1, 2, \dots, p$$

The Bonferroni $100(1-\alpha)\%$ simultaneous confidence intervals for the p population mean differences for $\mu_{1,i} - \mu_{2,i}$ is given as :

$$(\bar{x}_{1,i} - \bar{x}_{2,i}) \pm t\left(\frac{\alpha}{2p}, n_1+n_2-2\right) \sqrt{S_{ii, \text{pooled}} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

② Independent samples with $\Sigma_1 \neq \Sigma_2$

A) Univariate case

$$X_{1,i} \quad i=1, 2, \dots, n \sim N(\mu_1, \sigma_{11})$$

$$X_{2,i} \quad i=1, 2, \dots, n \sim N(\mu_2, \sigma_{22})$$

$\sigma_{11} \neq \sigma_{22}$ (from Levene's test)

The test of $H_0: \mu_1 - \mu_2 = \delta_0$ vs $H_a: \mu_1 - \mu_2 \neq \delta_0$
has test statistic

$$\frac{\bar{X}_1 - \bar{X}_2 - \delta_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \sim t(\alpha/2, v)$$

$$v = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2} \right)^2}{\frac{\left(\frac{s_1^2}{n_1} \right)^2}{n_1-1} + \frac{\left(\frac{s_2^2}{n_2} \right)^2}{n_2-1}} \quad (\text{Round this down ALWAYS})$$

(B) Multivariate case: Independent samples, $\Sigma_1 \neq \Sigma_2$

Result: (When n_1 and n_2 are large)

Let the sample sizes be such that n_1-p and n_2-p are large. Then an approximate $100(1-\alpha)\%$ confidence ellipsoid for $\mu - \mu_2$ is given by all $\mu - \mu_2$ satisfying

$$[\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)]^T \left[\frac{1}{n_1} S_1 + \frac{1}{n_2} S_2 \right]^{-1} [\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)] \leq \chi^2(\alpha, p)$$

Test statistic. Reject H_0 if T.S > $\chi^2(\alpha, p)$

- $100(1-\alpha)\%$ simultaneous confidence intervals for all linear combinations $\underline{a}^T (\mu_1 - \mu_2)$ is given by

$$\underline{a}^T (\bar{x}_1 - \bar{x}_2) \pm \sqrt{\chi^2(\alpha, p)} \sqrt{\underline{a}^T \left(\frac{1}{n_1} S_1 + \frac{1}{n_2} S_2 \right) \underline{a}}$$

Book example 6.5 (Problem 5 in example document)

When Sample sizes are not large

2 populations need to be MVN
 n_1 and n_2 each have to be greater than p

$$H_0: \mu_1 - \mu_2 = \underline{s}_0 \quad H_a: \mu_1 - \mu_2 \neq \underline{s}_0$$

Test statistic

$$T^2 = (\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2))^T \left[\frac{1}{n_1} S_1 + \frac{1}{n_2} S_2 \right]^{-1} (\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2))$$

$$(T^*)^2 = \frac{v-p+1}{vp} T^2$$

Reject H_0 if $(T^*)^2 > F(\alpha, p, v-p+1)$

v can be calculated as

$$\frac{p+p^2}{\sum_{i=1}^2 \frac{1}{n_i} \left\{ \text{tr} \left[\left(\frac{1}{n_i} S_i \left(\frac{1}{n_1} S_1 + \frac{1}{n_2} S_2 \right)^{-1} \right)^2 \right] + \left(\text{tr} \left[\frac{1}{n_i} S_i \left(\frac{1}{n_1} S_1 + \frac{1}{n_2} S_2 \right)^{-1} \right] \right)^2 \right\}}$$

where $\min(n_1, n_2) \leq v \leq n_1 + n_2$

The approximate $100(1-\alpha)\%$ confidence region is given by all $\underline{\mu}_1 - \underline{\mu}_2$ such that

$$(\bar{x}_1 - \bar{x}_2 - (\underline{\mu}_1 - \underline{\mu}_2))^T \left[\frac{1}{n_1} S_1 + \frac{1}{n_2} S_2 \right]^{-1} (\bar{x}_1 - \bar{x}_2 - (\underline{\mu}_1 - \underline{\mu}_2)) \leq \frac{\sqrt{p}}{v-p+1} F(\alpha, p, v-p+1)$$

where v is as given before.

Book example 6.6 (Problem 6 in class example)

Dependent Populations

A) Two dependent populations:

- Previously, we conducted hypothesis testing about the difference between two population means when the two samples are drawn independently from two different populations
- Here we deal with inferences for the mean of dependent populations which are also called paired populations.
- Use this test when observations from the two populations of interest are collected in pairs.

Check examples document for an example on when you have paired data.

Let X_{j1} denote the response to treatment 1 (or the response before treatment), and let X_{j2} denote the response to treatment 2 (or response after treatment) for the j th trial. Note: (X_{j1}, X_{j2}) are measurements recorded on the j th unit or j th pair of like units.

We analyze the n differences:

$$D_j = X_{j1} - X_{j2}, \quad j=1, 2, \dots, n$$

The differences D_j represent independent observations from $N(\delta, \sigma_\delta^2)$ where δ is the population mean difference between the 2 treatments, ie $E(X_{j1} - X_{j2})$

To test:

$H_0: \delta = 0$ (zero mean difference for treatments)

$H_1: \delta \neq 0$

Test Statistic: $t = \frac{\bar{D} - \delta}{S_\delta / \sqrt{n}} \sim t(\alpha/2, n-1)$

where $\bar{D} = \frac{1}{n} \sum_{j=1}^n D_j$ and $S_\delta^2 = \frac{1}{n-1} \sum_{j=1}^n (D_j - \bar{D})^2$

Reject H_0 if $|t| > t(\alpha/2, n-1)$

A $100(1-\alpha)\%$ confidence interval for the mean differences $\delta = E(X_{j1} - X_{j2})$ is given by:

$$\bar{D} \pm t(\alpha/2, n-1) \frac{S_\delta}{\sqrt{n}}$$

Check example (Problem 3) in example document

b) Multivariate Case

- p variables, n experimental units.
- The j th experimental unit has p observations given as

$X_{1,j1}$ = variable 1 under treatment 1 } Population 1
 $X_{1,j2}$ = " 2 " " "
 \vdots

$X_{1,jp}$ = " p under treatment 1

$$\begin{aligned}
 X_{2j1} &= \text{Variable 1 under treatment 2} \\
 X_{2j2} &= \quad " \quad 2 \quad " \quad 2 \\
 &\vdots \\
 X_{2jp} &= \quad " \quad p \quad \text{under treatment 2}
 \end{aligned}
 \quad \left. \right\} \text{Population 2}$$

The paired difference random variables becomes:

$$D_{j1} = X_{1j1} - X_{2j1}$$

$$D_{j2} = X_{1j2} - X_{2j2}$$

 \vdots

$$D_{jp} = X_{1jp} - X_{2jp}$$

So for the j th experimental unit, we have

$$\underline{D_j}^T = (D_{j1} \ D_{j2} \ \dots \ D_{jp}),$$

$$E(\underline{D_j}) = \underline{\delta} = \begin{pmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_p \end{pmatrix} \quad \text{and} \quad \text{Cov}(\underline{D_j}) = \sum_d \quad j \quad j=1, 2, \dots, n$$

$\underline{D_1}, \underline{D_2}, \dots, \underline{D_n}$ are independent $N_p(\underline{\delta}, \Sigma_d)$ random vectors.
Given the observed differences $\underline{d_j}^T = (d_{j1} \ d_{j2} \ \dots \ d_{jp}), j=1, 2, \dots, n$
from $\underline{D_1}, \dots, \underline{D_n}$, to test:

$$H_0: \underline{\delta} = \underline{0} \quad \text{vs} \quad H_a: \underline{\delta} \neq \underline{0}$$

$$\text{Test statistic: } \bar{T}^2 = n (\bar{\underline{d}} - \bar{\underline{\delta}})^T S_d^{-1} (\bar{\underline{d}} - \bar{\underline{\delta}})$$

where

$$\text{If } \bar{T}^2 > \frac{(n-1)p}{n-p} f(\alpha, p, n-p), \text{ reject } H_0.$$

A $100(1-\alpha)\%$ confidence region for $\bar{\delta}$ consists of all $\bar{\delta}$ such that

$$(\bar{D} - \bar{\delta})^\top S_d^{-1} (\bar{D} - \bar{\delta}) \leq \frac{(n-1)p}{n(n-p)} F(\alpha, p, n-p)$$

Also, $100(1-\alpha)\%$ simultaneous confidence intervals for the individual mean differences δ_i are given by

$$\bar{d}_i \pm \sqrt{\frac{(n-1)p}{(n-p)} F(\alpha, p, n-p)} \sqrt{\frac{s_{d,i}^2}{n}} ; i=1, 2, \dots, p$$

where \bar{d}_i is the i th element of \bar{D} and $s_{d,i}^2$ is the i th diagonal element of S_d

- The Bonferroni $100(1-\alpha)\%$ simultaneous confidence intervals for individual mean differences are

$$\bar{d}_i \pm t(\alpha/2q, n-1) \sqrt{\frac{s_{d,i}^2}{n}}$$

Note: ① If n and $n-p$ are both large, T^2 is approximately $\chi^2(\alpha, p)$ regardless of the form of the underlying population of the differences.

- ② For $n-p$ large, $\frac{(n-1)p}{n-p} F(\alpha, p, n-p) \stackrel{d}{=} \chi^2(\alpha, p)$ and normality need not be assumed.

Refer to example document (Problem 7) - Book example

6.1

Question: What if we have more than 2 populations?

Comparing several population means

Univariate ANOVA

Suppose we are interested in comparing 3 or more population means. As usual, we can draw independent samples from each population, calculate their sample means and use them in the analyses.

Eg. ① We want to compare average weight of chickens from 4 different poultry farms [One variable of interest: weight, 4 populations: Location 1-4]

② Checking differences in average maths scores for high schools in the Madison county.

Now, suppose we wish to compare means of g different populations, we can test

$$H_0: \mu_1 = \mu_2 = \dots = \mu_g \text{ against}$$

H_a : At least one group mean is different from the others.

Assumption: $x_{11}, x_{12}, \dots, x_{1n_1}$ is a random sample from $N(\mu_1, \sigma^2)$ and the random samples from each population are independent.

- Each population mean, μ_l , can be decomposed into the overall mean (grand mean) and a component which is due to the specific population, that is:

$$\mu_l = \mu + \tau_l$$

(lth population mean) (Grand mean) (lth population treatment effect)

H_0 and H_a can be written in terms of the individual treatment effects as

Constant
Variation

$H_0: \tau_1 = \tau_2 = \dots = \tau_g = 0$ (There is no treatment effect)

$H_a: \text{At least one of the } \tau_i \neq 0$ (There is a possible treatment effect)

- Each individual observation X_{ij} can be decomposed into $(\mu + \tau_i + e_{ij})$

$$X_{ij} = \mu + \tau_i + e_{ij} = \underset{\substack{\text{(overall)} \\ \text{mean}}}{\mu} + \underset{\substack{\text{(treatment)} \\ \text{effect}}}{\tau_i} + \underset{\substack{\text{(random)} \\ \text{error}}}{e_{ij}}$$

Subject to $\sum_{i=1}^g n_i \tau_i = 0$. $e_{ij} \sim N(0, \sigma^2)$
independent

- We estimate the model from the sample as follows:

$$X_{ij} = \bar{x} + (\bar{x}_i - \bar{x}) + (x_{ij} - \bar{x})$$

(Observation) $\begin{pmatrix} \text{Overall} \\ \text{Sample mean} \\ - \text{estimate of } \mu \end{pmatrix}$ $\begin{pmatrix} \text{estimated treatment} \\ \text{effect} - \\ \text{estimate of } \tau_i \end{pmatrix}$ $\begin{pmatrix} \text{residual} \\ \text{estimate of } e_{ij} \end{pmatrix}$

- Question: How much of the total variation in the data comes from variation within the group or variation between the group?

Example (Check example in examples document)

Population 1: observations
 $9 \ 6 \ 9$

$n_1 = 3$

$\bar{x}_1 = 8$

Population 2: $0 \ 2$

$n_2 = 2$

$\bar{x}_2 = 1$

Population 3: $3 \ 1 \ 2$

$n_3 = 3$

$\bar{x}_3 = 2$

Overall mean: $\frac{9+6+9+0+2+3+1+2}{8} = 4 = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + n_3 \bar{x}_3}{n_1 + n_2 + n_3}$

We wish to find the 3 sum of squares

a) $SST = \text{Sum of squares total}$

= Sum of squared deviation of each observation from
the overall mean

$$= (9-4)^2 + (6-4)^2 + (9-4)^2 + (0-4)^2 + (2-4)^2 + (3-4)^2 + (1-4)^2 + (2-4)^2$$

$$= \boxed{88}$$

$$= \sum_{l=1}^g \sum_{j=1}^{n_l} (x_{lj} - \bar{x})^2$$

$$= SS_{\text{cor}} \quad (\text{Total corrected } SS)$$

$$df(SS_{\text{cor}}) = \left(\sum_{l=1}^g n_l \right) - 1 \quad - \text{Why?}$$

b) $SSW = \text{"Sum of squares within the individual groups"}$

= Sum of squared deviation of each observation
from its corresponding population mean.

$$= (9-8)^2 + (6-8)^2 + (9-8)^2 + (0-2)^2 + (2-2)^2 + (3-2)^2 + (1-2)^2 + (2-2)^2$$

$$= 10$$

$$= \sum_{l=1}^g \sum_{j=1}^{n_l} (x_{lj} - \bar{x}_l)^2$$

$$= SS_{\text{res}}$$

$$df(SS_{\text{res}}) = \left(\sum_{l=1}^g n_l \right) - g \quad - \text{Why?}$$

$$\begin{aligned}
 \text{c) } SSB &= \text{"Sum of squares between the groups"} \\
 &= (8-4)^2 + (8-4)^2 + (8-4)^2 + (2-4)^2 + (2-4)^2 + (3-4)^2 + (1-4)^2 + (2-4)^2 \\
 &= 78 \\
 &= \sum_{l=1}^g n_l (\bar{x}_l - \bar{x})^2 \\
 &= SS_{tr}
 \end{aligned}$$

$$df(SS_{tr}) = g-1$$

$$\begin{aligned}
 \sum_{l=1}^g \sum_{j=1}^{n_l} x_{lj}^2 &= \sum_{l=1}^g n_l \bar{x}^2 &+ \sum_{l=1}^g n_l (\bar{x}_l - \bar{x})^2 &+ \sum_{l=1}^g \sum_{j=1}^{n_l} (x_{lj} - \bar{x}_l)^2 \\
 (SS_{obs}) &\quad (SS_{mean}) &\quad (SS_{tr}) &\quad (SS_{res})
 \end{aligned}$$

ANOVA Table

Source of Variation	Sum of squares (SS)	Degrees of freedom (df)
Treatments	$SS_{tr} = \sum_{l=1}^g n_l (\bar{x}_l - \bar{x})^2$	$g-1$
Residual (error)	$SS_{res} = \sum_{l=1}^g \sum_{j=1}^{n_l} (x_{lj} - \bar{x}_l)^2$	$\sum_{l=1}^g n_l - g$
Total	$SS_{tot} = \sum_{l=1}^g \sum_{j=1}^{n_l} (x_{lj} - \bar{x})^2$	$\sum_{l=1}^g n_l - 1$

$$F^* = \frac{SS_{tr}/(g-1)}{SS_{res}/(\sum_{l=1}^g n_l - g)} \sim F(\alpha, g-1, \sum_{l=1}^g n_l - g)$$

Reject H_0 if $F^* > F(\alpha, g-1, \sum_{l=1}^g n_l - g)$

Back to our example (Testing at $\alpha=0.01$)

$$F^* = \frac{78/2}{10/5} = 19.5$$

$F^* = 19.5 > F(0.01, 2, 5) = 13.27$ so reject H_0 .

Multivariate ANOVA

- Used when comparing more than 2 populations.
- Random samples from each of g populations are arranged as

Population 1 : $\tilde{X}_{11}, \tilde{X}_{12}, \dots, \tilde{X}_{1m}$,

Population 2 : $\tilde{X}_{21}, \tilde{X}_{22}, \dots, \tilde{X}_{2n_2}$

⋮ ⋮ ⋮

Population g : $\tilde{X}_{g1}, \tilde{X}_{g2}, \dots, \tilde{X}_{gn_g}$

- We use MANOVA to investigate whether the population mean vectors are the same. If not, there are methods to investigate further which mean components differ significantly.

Assumptions about the structure of the data for one-way MANOVA

1. $\tilde{x}_{11}, \tilde{x}_{12}, \dots, \tilde{x}_{1n_1}$ is a random sample of size n_1 from a population with mean μ_l , $l=1, 2, \dots, g$. The random samples from the different populations are independent.
 2. All populations have a common covariance matrix, Σ .
 3. Each population is multivariate normal. (Can be relaxed when sample sizes n_l are large.)
- Refer to problem 9 (Book example 6.9)

Model

$$\tilde{x}_{lj} = \underline{\mu} + \tilde{\tau}_l + \tilde{e}_{lj}, \quad j=1, 2, \dots, n_l \quad \text{and} \quad l=1, 2, \dots, g$$

subject to $\sum_{l=1}^g n_l \tilde{\tau}_l = 0$

- $\tilde{e}_{lj} \sim N_p(0, \Sigma)$, $\underline{\mu}$ is the overall mean,
- $\tilde{\tau}_l$ is the l th treatment effect

The model is estimated from the sample as:

$$\tilde{x}_{lj} = \bar{\tilde{x}} \quad \left(\begin{array}{c} \text{observation} \\ \text{overall mean} \end{array} \right) + (\bar{\tilde{x}}_l - \bar{\tilde{x}}) \quad \left(\begin{array}{c} \text{estimated} \\ \text{treatment effect} \\ \hat{\tau}_l \end{array} \right) + (\tilde{x}_{lj} - \bar{\tilde{x}}_l) \quad \left(\begin{array}{c} \text{residual} \\ \hat{e}_{lj} \end{array} \right)$$

As in the univariate case, the total sum of squares is decomposed into the sum of squares between groups and the sum of squares within groups.

MANOVA Table for Comparing Population Mean Vectors

Source of variation	Matrix of sum of squares and cross products	Degrees of freedom
Treatment (Between)	$B = \sum_{l=1}^g n_l (\bar{X}_{l\cdot} - \bar{X})(\bar{X}_{l\cdot} - \bar{X})^T$ $= (n_1-1)S_1 + (n_2-1)S_2 + \dots + (n_g-1)S_g$	$g-1$
Residual (Error)	$W = \sum_{l=1}^g \sum_{j=1}^{n_l} (\bar{X}_{lj} - \bar{X}_{l\cdot})(\bar{X}_{lj} - \bar{X}_{l\cdot})^T$	$\sum_{l=1}^g n_l - g$
Total	$B + W = \sum_{l=1}^g \sum_{j=1}^{n_l} (\bar{X}_{lj} - \bar{X})(\bar{X}_{lj} - \bar{X})^T$	$\sum_{l=1}^g n_l - 1$

Testing

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

$$H_a: \mu_i \neq \mu_j \text{ for at least one } i \neq j$$

We calculate

$$\Lambda^* = \frac{|W|}{|B+W|} = \frac{\left| \sum_{l=1}^g \sum_{j=1}^{n_l} (\tilde{x}_{lj} - \bar{\tilde{x}}_l)(\tilde{x}_{lj} - \bar{\tilde{x}}_l)^T \right|}{\left| \sum_{l=1}^g \sum_{j=1}^{n_l} (x_{lj} - \bar{x}) (x_{lj} - \bar{x})^T \right|}$$

Λ^* is known as Wilks' lambda

- The exact distribution of Λ^* can be derived for the special cases as listed in the table below:

i.e. ($\sum_{i=1}^g n_i = n$)

For any other case and large sample sizes, a modification of λ^* due to Bartlett can be used in the hypothesis test.

$$-\left(n-1 - \frac{p+g}{2}\right) \ln \lambda^* = -\left(n-1 - \frac{p+g}{2}\right) \ln \left(\frac{|W|}{|B+W|} \right)$$

$$\sim \chi^2(\alpha, p(g-1))$$

Reject H_0 if

$$-\left(n-1 - \frac{p+g}{2}\right) \ln \left(\frac{|W|}{|B+W|} \right) > \chi^2(\alpha, p(g-1))$$