

Believe it or not, the semester begins.

1.22

o Tests = Applied Multivariate Statistical Analysis R Johnson

Week 1-2

o Test = 2.126 Test 1

Hum 6122

6.1 Test 2

Multivariate

o Final project = 2-3, suitable dataset for analysis

A Dataset

Chapter 1 Introduction to Multivariate Statistics

Intro:

measures

Def: Variables
Cases

1
2
3
4
5
6
7

Types: Analysis of interdependence \rightarrow no X or Y

Principal component analysis

Factor analysis

Multidimensional Scaling

Cluster analysis

- dependence

Multivariate

Structural equation models

Canonical Correlation

Multivariate analysis of variance MANOVA

Covariance MANOVA

Discriminatory / Classification analysis

Example:

Data layout

- n experimental, p variables / characteristics / observed / measured
- $X \rightarrow n$ rows, p columns
- X_{jk} k^{th} variable in the j^{th} unit

HM.

mean vector in $R \rightarrow$ Descriptive Statistics

- Variances: $S_k^2 = S_{kk} = \frac{1}{n} \sum_{j=1}^n (x_{jk} - \bar{x}_k)^2 = \frac{W_{kk}}{n}$, $k=1, \dots, p$

n , not $n-1$

\downarrow bias estimator \downarrow not bias

- Covariance: $S_{jk} = \frac{1}{n} \sum_{j=1}^n (x_{ji} - \bar{x}_i)(x_{jk} - \bar{x}_k) = \frac{W_{jk}}{n}$, $i=1, \dots, p$, $k=1, \dots, p$

n , not $n-1$

~~Means~~

~~Sum of Squares~~

Homework 6/22

1.29

Week 2. -)

Ex 1

$$b'd = [2, -1, 4, 0] \begin{bmatrix} -1 \\ 3 \\ -2 \\ 1 \end{bmatrix}$$

$$= 2(-1) + (-1) \cdot 3 + 4 \cdot (-2) + 0 = -13$$

$$= -13$$

$$b'b = 2^2 + (-1)^2 + 4^2 + 0^2 = 21$$

$$d'd = (-1)^2 + 3^2 + (-2)^2 + 1^2 = 15$$

$$(-13)^2 \leq 21 \times 15$$

$$169 \leq 315$$

Ex 2

a) Graph: $n=3$ $p=2$ 3×2 data matrix

1.

$$x_1 = [9, 1]$$

$$y_1 = \begin{bmatrix} 9 \\ 5 \\ 1 \end{bmatrix}$$

$$y_2 = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix}$$

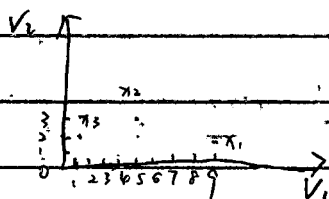
Step 1: n, p

$$x_2 = [5, 2]$$

Step 2: n rows (people)

$$x_3 = [1, 2]$$

Scatter plot.



2.

① Vector

$$\bar{x}_1 = \frac{9+5+1}{3} = 5$$

② 2 element

$$\bar{x}_2 = \frac{1+3+2}{3} = 2$$

$$\bar{x} = [5, 2]$$

↓ somewhere in the middle. 在图中.

b). deviation vector

subtract average ^{by element} \bar{v}_j from each column.

$$d_1 = \begin{bmatrix} 9 \\ 5 \\ 1 \end{bmatrix} - \begin{bmatrix} 5 \\ 5 \\ 5 \end{bmatrix} = \begin{bmatrix} 4 \\ 0 \\ -4 \end{bmatrix}$$

How far is each observation from its mean.

Vector in 2 dimensions

$$d_2 = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix} - \begin{bmatrix} 2 \\ 2 \\ 2 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix}$$

always add up to 0

Correlation, variance-covariance

length

$$Ld_1 = \sqrt{d_1 \cdot d_1} = \sqrt{4^2 + 0^2 + (-4)^2} = \sqrt{32}$$

$$Ld_2 = \sqrt{d_2 \cdot d_2} = \sqrt{(-1)^2 + 1^2 + 0^2} = \sqrt{2}$$

$$d_1 \cdot d_2 = \begin{bmatrix} 4 \\ 0 \\ -4 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix} = -4 + 0 + 0 = -4$$

inner product

4	-1
0	1
-4	0

$$\cos \theta = \frac{d_1 \cdot d_2}{\sqrt{d_1 \cdot d_1} \cdot \sqrt{d_2 \cdot d_2}} = \frac{d_1 \cdot d_2}{Ld_1 \cdot Ld_2} = \frac{-4}{\sqrt{32} \sqrt{2}} = -\frac{1}{2} \quad | 120^\circ$$

$\theta = 120^\circ$ (θ between vector)

or between column $\begin{pmatrix} 1 \\ 3 \\ 2 \end{pmatrix}$ and $\begin{pmatrix} 9 \\ 5 \\ 1 \end{pmatrix}$

(对 @)

Variance $\rightarrow S$

$$S_{11} = \frac{Ld_1}{n} = \frac{32}{3}$$

Standard deviation $= \sqrt{S} = \sqrt{Var}$

$$S_{22} = \frac{2}{3}$$

length
 n (or $n-1$)

$$= \text{variance} = S_{ii}$$

$$r_{12} = -\frac{1}{2}$$

as of angle = correlation

H7-1

GM { Uni
Multi

Estimation and Hypothesis Testing

Hw 6122

2.5

Week 3-2

Univariate case

Suppose $X_1, \dots, X_n \sim N(\mu, \sigma^2)$

Then $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \hat{\mu}$ is MLE of μ

$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$ is MLE of σ^2

Then $\bar{X} \sim N(\mu, \frac{\sigma^2}{n})$

$$(n-1)S^2 \sim \chi_{n-1}^2, \text{ where } S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Multivariate:

$X_1, \dots, X_n \sim N_p(\mu, \Sigma)$

Then $\bar{X} = \hat{\mu}$ is p.d.f. vector

$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$ is p.d.f. matrix

$\hat{\mu}$ is MLE of μ

$\hat{\Sigma}$ is MLE of Σ

Σ covariance matrix

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$$

$\hat{\mu} \sim N_p(\mu, \frac{\Sigma}{n})$

no chi χ^2 in multi for Σ instead Wp

$(n-1)S^2 \sim W_p(n-1, \Sigma)$

is Wishart distribution with $n-1$ df and dimension $p \times p$ and scale Σ

Recall cheat sheet

$\bar{X} \rightarrow$

Hypotheses - when Σ is known $\rightarrow Z$ test

$S \rightarrow$ matrix

Recall the univariate case (Imperial)

$X_1, \dots, X_n \sim N(\mu, \sigma^2)$

$H_0: \mu = \mu_0$ vs $H_a: \mu \neq \mu_0$

Test statistic:

$$Z_{\text{test}} = \frac{\bar{X} - \mu_0}{\sqrt{\frac{\hat{\sigma}^2}{n}}} = \frac{\bar{X} - \mu_0}{\frac{\hat{\sigma}}{\sqrt{n}}} \sim N(0,1)$$

Z-test around 0. RR: we reject H_0 if $|Z_{test}| > Z_{\alpha/2}$ (2-side)

no reject

when α is the given significant level (usually 0.05)

Z-test

Slightly change the rule for multi

reject

RR: rejection rule

Note: RR can be written as $Z_{test}^2 > Z_{\alpha/2}^2 = (Z_{\alpha/2}^*)^2$

↳ Multivariate $Z'Z$

(inner product of it self)

P-value would be same

(Multivariate)

Generic case

Multivariate case: Assume $X_1, \dots, X_k \sim N_p(\mu, \Sigma)$

Σ not given

P-variate: if $\hat{\Sigma} = \Sigma$ then

$$\text{This means } \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix}$$

$$\text{Multivariate } Z^2 \left(\hat{\Sigma} \text{ is known} \right)$$

Function of Vectors

$$\text{Then } H_0: \mu = \mu_0 \Rightarrow M = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_p \end{pmatrix} = \begin{pmatrix} \mu_{10} \\ \vdots \\ \mu_{p0} \end{pmatrix} = \mu_0$$

ATI: ① Hypo

② T-test, Z^2 test

③ Ordered Value q

$$H_1: \mu \neq \mu_0$$

④ p-values

⑤ RR

为啥???

Why not (分开算)

Type I error at α for each of p tests

multiple comparison

Note: Why don't we do p univariate tests

$$H_0: \mu_i = \mu_{i0}, i = 1, \dots, p$$

A: The type I error will become inflated

Ex: test $p=3$

$$\frac{\mu_1}{\sigma} = \frac{\mu_2}{\sigma} = \frac{\mu_3}{\sigma} = 0$$

Each test, $P(\text{Type I error}) = 0.05$

$P(\text{At least one Type I error}) = 1 - P(\text{no error}) = 1 - 0.95^2 = 0.14$

Benjamini

Correction ???

independent

Σ is level

Null hypothesis:

$H_0: \mu = \mu_0$ vs $H_a: \mu \neq \mu_0$ (μ is $p \times 1$)

T.S.

$$U_{n1} = Z^2 = \frac{(\bar{X} - \mu_0)^2}{\frac{S^2}{n}}$$

1/R 2

$$Z^2 = n(\bar{X} - \mu_0)' \Sigma^{-1} (\bar{X} - \mu_0)$$

Scalar $1 \times p$ $p \times p$ $p \times 1$

X' transpose

= scalar (number)

$X' = t(X)$

follow χ^2 with p df.

$\Rightarrow Z^2_{test} \sim \chi^2_p$

$$= n(\bar{X} - \mu_0)' (S^2)^{-1} (\bar{X} - \mu_0)$$

Hypotheses

distance

$\Sigma \rightarrow$ direction

distance χ^2

direction χ^2

RR: Reject H_0 if $Z^2_{test} > \chi^2_{p(\alpha)}$

$pchisq(\chi^2)$

2.2.2

Ex: $p=2$, $n=20$

$x_1 = \text{Height}$

$x_2 = \text{Weight}$

Variance is

in 2 units

Assume $\Sigma = \begin{pmatrix} 20 & 100 \\ 100 & 1000 \end{pmatrix}$

(Assume 6)

Height

variance of

covariance

G_{12}

variance of weight

Test: $H_0: \mu = \begin{pmatrix} 70 \\ 170 \end{pmatrix}$ vs $H_1: \mu \neq \begin{pmatrix} 70 \\ 170 \end{pmatrix}$

$R \rightarrow$ correlation matrix

Given: $\bar{x} = \begin{pmatrix} 71.45 \\ 164.7 \end{pmatrix}$

(invert Σ) $\begin{pmatrix} 10.2 \\ 10.2 \end{pmatrix}$

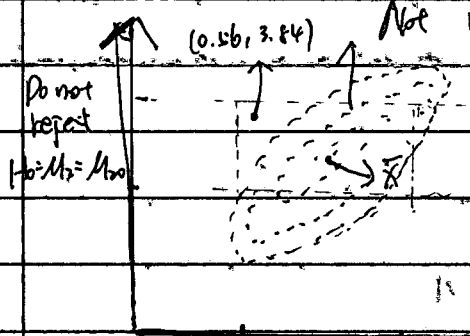
T.S. = $Z_{stat} = 5.4$

$X_{crit} = 1.99$

Conclusion: $5.4 > 1.99$

\Rightarrow reject H_0

P-value = 0.015



$(0.56, 3.84)$

Not rejecting

rectangular \rightarrow do not reject both univariate

rectangular: multi reject, uni reject

But multi reject, uni reject

Do not reject univariate

$H_0: \mu_1 = \mu_2$

$Ax = c$ solve \Rightarrow solve (A, c) given in R

solve $(A) \rightarrow$ inverse of A

What if G^2 is unknown

$A =$ the test

Univariate are:

$x_1, \dots, x_n \sim N(\mu, G^2)$

$H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$ (default in R)

$$t_{\text{test}} = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t_{n-1}$$

Sample Variance

- Covariance

RR: Reject H_0 if $|t_{\text{test}}| > t_{n-1}(\frac{\alpha}{2})$

Variance -

$H_0: \mu \leq \mu_0$

Covariance -

• Multivariate \rightarrow no one tail test.

H7-2

Σ known χ^2
 Σ unknown $\left\{ \begin{array}{l} \text{small sample } T^2 \\ \text{large sample } F \end{array} \right.$ $\hat{\Sigma} = \frac{n-1}{n} S$

HWIDM 6122

2.10

Week 4-1

Test for μ when Σ is unknown

- Hotelling's T^2 test

Mult: T^2

Let $x_1, \dots, x_n \sim N_p(\mu, \Sigma)$

$H_0: \mu = \mu_0$ vs $H_1: \mu \neq \mu_0$

where μ is $p \times 1$

Σ is $p \times p$ (unknown)

T.S.

$$T^2_{\text{test}} = n(\bar{x} - \mu_0)' S^{-1} (\bar{x} - \mu_0)$$

Recall the univariate case:

$$t^2_{\text{test}} = \frac{(\bar{x} - \mu_0)^2}{s^2/n} = n(\bar{x} - \mu_0)^2 / s^2$$

properties

near 0 \Rightarrow OK

big \Rightarrow reject null

under H_0 , $T^2_{\text{test}} \sim T^2_{p, n-1}$

similar F, recode χ^2 . Hotelling

RR: reject $T^2_{\text{test}} > T^2_{p, n-1, \alpha}$

where $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ is $p \times 1$

$$S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$$

is $p \times p$

tables in optional book

Properties:

1) We must have $n-1 > p$, otherwise S is singular and S^{-1} doesn't exist.

more features than observations

2) In multivariate case, we don't consider one-tail alternatives like $H_1: \mu > \mu_0$

3) T^2 can be converted to F distribution

$$T^2_{p, n-1} = \frac{(n-1)p}{n-p} \cdot F_{p, n-p}$$

4) When $p=1$, T^2 is the same as t^2_{n-1} distribution

5) When $n \rightarrow \infty$, $T^2_{p, n-1} \rightarrow \chi^2_p$ (univariate)

Ex: Table 5.1 on p215

$p=3$, $n=20$, $H_0: \mu = \begin{pmatrix} 4 \\ 10 \\ 10 \end{pmatrix}$

x_1 = sweat rate

$H_1: \mu \neq \begin{pmatrix} 4 \\ 10 \\ 10 \end{pmatrix}$

x_2 = sodium

x_3 = potassium

$\alpha=0.1$

$$(\text{crit. value}) = \frac{(n-1)p}{n-p} F_{p, n-p, \alpha}$$

AAA? why
 no $\begin{pmatrix} 1 \\ 0 \end{pmatrix} < \dots$
 matrix $\neq \text{I}$

*** F table, 1 page 1x

Ques: * Critical

(F-table)

draw

$$F_{0.05} = 9.74$$

From the table

$$\text{Crit. value} = \frac{(n-1)p}{n-p} F_{3, 13} (0.1) = 2.17$$

$$p\text{-value} = 2.44$$

* F related to t, cross-multiply \Rightarrow p-value

$$p\text{-value} = 0.065 < 0.1$$

reject H_0

*** ~~reject~~ test *** 不若 ???

Idea:

$$T.S. \lambda = \frac{\max_{H_0} L(\mu, \Sigma)}{\max L(\mu, \Sigma)}$$

algebra + Calculus

R.R: If $\lambda < \text{crit. value}$, then reject H_0

MANOVA

[R] \Rightarrow Hotelling

$$\text{It can be shown that } \lambda = \left(\frac{\det(\hat{\Sigma})}{\det(\hat{\Sigma}_0)} \right)^{1/2}$$

$$\text{where } \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$$

$$\hat{\Sigma}_0 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_0)(x_i - \mu_0)'$$

$$\text{Fact: } \lambda_{\frac{1}{n}}^2 = \left(1 + \frac{T^2}{n-1} \right)^{-1}$$

paired

① 2 sample

② Which variable?

MANOVA

+ Comparing two means from independent sample

Male vs Female: Same score or not?

Is the difference one?

Company 2 means

① H_0

② T.S. - t-test - pooled

③ D.F. value $df = n_1 + n_2 - 2$

④ R.R. - normal

Normal?

- The two samples are independent

Normal?

- Equal variances: $G_1^2 = G_2^2 = 6$

- Normal Populations

Assume:

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

Standard deviation pooled

Test: $H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$

$$t\text{-test: } \frac{\bar{x}_1 - \bar{x}_2}{S_{\text{pooled}} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{n_1 + n_2 - 2}$$

RR: If $|t_{test}| > t_{n_1+n_2-2}(\frac{\alpha}{2})$, then reject H_0



In multivariate case

$X_{11}, \dots, X_{1p} \sim N_p(\mu_1, \Sigma_1)$ & $\mu_1 = \mu_2$ same dimensions

$X_{21}, \dots, X_{2p} \sim N_p(\mu_2, \Sigma_2)$

Assumptions:

- The two samples are independent
- Equal cov matrices: $\Sigma_1 = \Sigma_2 = \Sigma$
- Normal populations

$H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$

S_{pooled} - 1 sample

μ_1 & μ_2 are $p \times 1$ vectors.

~~Def~~ $S_{pooled} = \frac{1}{n_1+n_2-2} ((n_1-1)S_1 + (n_2-1)S_2)$

I.S. $T_{test} = \frac{n_1 n_2}{n_1+n_2} (\bar{X}_1 - \bar{X}_2)' S_{pooled}^{-1} (\bar{X}_1 - \bar{X}_2)$ ★★

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

Under H_0 , $T_{test}^2 \sim T_{p, n_1+n_2-2}^2$

It can be shown that

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

RR: Reject H_0 when

$T_{test}^2 > \text{crit value}$

Hw:

$$(A'A)' = A'(A')' = A'A$$

$M' = M \Rightarrow$ symmetric

Rule: $(AB)' = B'A'$

$$(A')' = A$$

$n \times p$

2 group's $\left\{ \begin{array}{l} \text{Cov.} \\ \text{Diff.} \end{array} \right.$ $\left\{ \begin{array}{l} \text{Cov.} \\ \text{Diff.} \end{array} \right.$

Assumptions

H7-3

μ_0

Hand M 6122

$$W = \frac{x_1 + x_2 + x_3}{3}$$

0.12

Week 4-2

$$= \frac{1}{3}x_1 + \frac{1}{3}x_2 + \frac{1}{3}x_3$$

$$= a'x$$

$$a = \begin{bmatrix} \frac{1}{3} \\ \frac{1}{3} \\ \frac{1}{3} \end{bmatrix}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

OL

for 2 groups' multi difference

paired 2 groups' multi

$$E(a'x) =$$

$$Var(a'x) =$$

Group = μ_0

Hypotheses - Part 3

- Simulation CI's

Confidence Interval

Permutation Adjustment

Given: $x_{11}, \dots, x_{1n_1} \sim N_p(\mu_1, \Sigma)$

$x_{21}, \dots, x_{2n_2} \sim N_p(\mu_2, \Sigma)$

Multidimensional Confidence Region

Goal: CI for $\mu_1 - \mu_2$, which is $p \times 1$

Recall: Univariate case

Variance of Covariance

$H_0: \mu_1 - \mu_2 = \mu_0$ vs $H_1: \mu_1 - \mu_2 \neq \mu_0$

$$t_{test} = \frac{\bar{x}_1 - \bar{x}_2 - \mu_0}{\text{SE}}$$

$$SE = \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Table value

\Rightarrow CI for $\mu_1 - \mu_2$ is

$\alpha \rightarrow t \rightarrow$

CI

$$\bar{x}_1 - \bar{x}_2 \pm t_{\alpha/2} \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

Why $\mu_1 \neq \mu_2$?

Table

Univariate case

$$T_{test} = \frac{n_1 n_2}{n_1 + n_2} (\bar{x}_1 - \bar{x}_2 - \mu_0)' S_p^{-1} (\bar{x}_1 - \bar{x}_2 - \mu_0)$$

$$T_{test} \sim \frac{(n_1 + n_2 - 2)P}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}$$

$$\text{Let } Q = \frac{(n_1 + n_2 - 2)P}{n_1 + n_2 - p - 1} \sim F_{p, n_1 + n_2 - p - 1}$$

Linear Transformation

Scale: F

where α is given (0.05, 0.01, ...)

For any $p \times 1$ vector a

① valid for any a

100(1- α)% CI for $a'(\mu_1 - \mu_2)$ to Pa diff test - plug in a

$$\text{is } a'(\bar{x}_1 - \bar{x}_2) \pm C \sqrt{a' S_p a (\frac{1}{n_1} + \frac{1}{n_2})}$$

② pooled variance

In particular we can choose

$$A = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow i^{\text{th}} \text{ row} \rightarrow p \times 1$$

whether variable has a significant difference

x_{ji}
group variables

$$\text{then } a'(u_i - u_j) = \mu_{ii} - \mu_{jj}$$

\Rightarrow linear transformation

and the CI for $\mu_{ii} - \mu_{jj}$ is

$$\frac{\mu_{ii} + \mu_{jj}}{2} - \frac{\mu_{ii} + \mu_{jj}}{2}$$

$$\left[\bar{x}_{ii} - \bar{x}_{jj} \pm \sqrt{s_{ii} \left(\frac{1}{n_i} + \frac{1}{n_j} \right)} \right]$$

test of variable

$$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \text{not variable} \rightarrow \bar{\mu}_{ii} - \bar{\mu}_{jj}$$

$$\begin{bmatrix} 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \text{1st variable} \rightarrow \mu_{12} - \mu_{21}$$

Remark:

B: When testing

$$H_0: \mu_1 - \mu_2 = 0 \text{ vs } H_a: \mu_1 - \mu_2 \neq 0$$

and we reject H_0 , how do we know which variable contributed the most to rejecting H_0 .

A: Calculate

$$\hat{a} = S_p^{-1} (\bar{x}_1 - \bar{x}_2)$$

Then the largest in abs value element of \hat{a} is corresponding to the var. contributing the most to rejecting the H_0 .

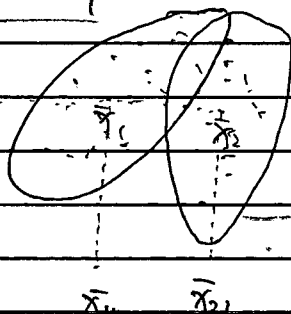
Covariance Matrix

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

$$\frac{1}{n_1 + n_2 - 2} \left[\sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2 \right]$$

(Which variable contribute most to reject H_0 .)

$$\bar{x}_{12} = \bar{x}_{21}$$



*** $\mu_1 - \mu_2 \neq 0$ reject Difference

$\left. \begin{array}{l} \text{Indep test } \mu \text{ vs } \mu - \mu \\ \text{dep test } x_1 - x_2 \text{ vs } \bar{d} \end{array} \right\}$

Independent -

H₀ to
 reject Null
 follow up

paired obs

Assume one -

Suppose we have two sample:

Sample size the same!

$$x_{11} \quad x_{21} \rightarrow d_1 = x_{11} - x_{21}$$

\vdots

$$x_{1n} \quad x_{2n} \rightarrow d_n = x_{1n} - x_{2n}$$

Assume: there is a natural pairing between x_{1i} & x_{2i} $i=1, \dots, n$

Before & After pre-post test.

Dependence

that is the two samples are dependent

One Sample

How to link them

Lab1 Lab2

Still 1 - n \rightarrow independence

but 1-2 group \rightarrow dependent

$$x_{1i} \quad x_{2i}$$

We reduce the two samples to one sample by

calculating the difference: $d_i = x_{1i} - x_{2i} \quad i=1, \dots, n$

Assume:

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix} \right)$$

T.s. $t_{test} = \frac{\bar{d}}{s.d. / \sqrt{n}}$, where \bar{d} is average difference one sample t-test
 s.d. is st. dev of difference

Under H_0 , $t_{test} \sim t_{n-1}$

RR: Reject H_0 if $|t_{test}| > t_{\alpha/2}$

$H_0: \mu_1 - \mu_2 = 0$ vs $H_1: \mu_1 - \mu_2 \neq 0$

Check the normality distribution \Rightarrow post hoc? $\star \star \star$!

Alternative case:

$$\begin{array}{lcl} x_{11} - x_{21} = d_1 & \left. \begin{array}{l} \vdots \\ \vdots \end{array} \right\} & \text{vector} \\ \vdots & \rightarrow & \vdots \\ x_{1n} - x_{2n} = d_n \end{array}$$

Assumption $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim N_{2p} \left(\begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \begin{pmatrix} \Sigma_1 & \Sigma_{12} \\ \Sigma_{21} & \Sigma_2 \end{pmatrix} \right)$

Hypotheses:

$H_0 = \mu_1 = \mu_2$ vs $\mu_1 \neq \mu_2$ (p x 1 vector)

Cheat Sheet

T.S.

$T_{\text{test}}^2 = n \bar{d}' S_d^{-1} \bar{d}$

where \bar{d} is mean diff (p x 1)

S_d is cov. matrix of diff (p x p)

Too large \rightarrow reject

\hookrightarrow How large \rightarrow crit val \rightarrow f-stats

$T_{\text{test}}^2 > \frac{(n-1)p}{n-p} F_{p, n-p}(\alpha)$

$n \rightarrow$ sample size always

*** 2 groups?

separate?

\downarrow
both

d is normally distributed

HWDA 6122

~~HWDA 6122~~ - Part 2

2.19

$$\text{Model: } x_{ij} = \mu + \tau_i + \epsilon_{ij}$$

Week 5-2

$$i = 1, \dots, g$$

$$j = 1, \dots, n_i \quad n = n_1 + \dots + n_g$$

For X'

$x_{ij}, \mu, \tau_i, \epsilon_{ij}$ are all $p \times 1$ vectors

$\epsilon_{ij} \stackrel{\text{ind.}}{\sim} N_p(0, \Sigma)$

with $\sum_{i=1}^g n_i \tau_i = 0$

$$H_0: \tau_1 = \tau_2 = \dots = \tau_g = 0$$

Idea:

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

$$\Rightarrow \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x})(x_{ij} - \bar{x})' \quad (p \times p)$$

everything is constant (why not inner product $A'A$)
~~inner product $C'A'A'$~~

$$B = \sum_{i=1}^g n_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})' \quad \text{Between}$$

$$+ \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)' \quad \text{Within}$$

Define:

$$W = \sum_{i=1}^g \sum_{j=1}^{n_i} (\bar{x}_{ij} - \bar{x}_i)(\bar{x}_{ij} - \bar{x}_i)'$$

$$= (n-1)S + \dots + (n_g-1)S_g$$

Generalized

$$B = \sum_{i=1}^g n_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})'$$

***?

$$\text{Total} = B + W$$

Generalized

$$\text{T.S. } \lambda_{\text{test}} = \frac{\det(W)}{\det(B+W)}$$

Convenient?

Re. Reject H_0 if λ_{test} is "small"

df ~~***~~

n-g

B.t.

Determinate \uparrow

W.d.

W \Rightarrow residual

error

within

g-1

n-1

with

If not in 6.3 $\rightarrow \chi^2$.

Otherwise, use Bartlett approximation

$$-(n-1 - \frac{(p+g)}{2}) \cdot \ln(\lambda_{test}) \approx \chi^2_{p(g-1)}$$

Likelihood Ratio Test

$\lambda_{test} = \prod_{i=1}^g (1 + \frac{1}{\lambda_i})$, where $\lambda_1, \dots, \lambda_s$ eigenvalues of $W^{-1}B$

not full rank $\Rightarrow S$

don't know eigenvalue

\Rightarrow the number of the non zero eigenvalue

Ex (6.9 on p.304): $p=2, g=3, n_1=3, n_2=2, n_3=3, n=8, \alpha=0.01$

$X_i \rightarrow$ Measurement from Group i

π_1

π_2

9 3

0 4 3 8

6 2

2 0 1 9

Whenever is called

9 7

2 7

column

①

$$\bar{X}_1 = \begin{bmatrix} 9 \\ 6 \\ 9 \\ 4 \end{bmatrix}$$

$$\bar{X}_2 = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\bar{X}_3 = \begin{bmatrix} 2 \\ 8 \end{bmatrix}$$

Covariance of

Midterm (bivariate)

$$S = \begin{bmatrix} 78 & -12 \\ -12 & 48 \end{bmatrix}$$

Sum of square treatment

Sum of square residual

Always symmetric

Covariance

between groups

$$B+W = \begin{bmatrix} 88 & -11 \\ -11 & 32 \end{bmatrix}$$

$$W = \begin{bmatrix} 10 & 1 \\ 1 & 24 \end{bmatrix}$$

Covariance within group (all three)

②

$$\lambda_{test} = \frac{|W|}{|B+W|} = \frac{\begin{vmatrix} 10 & 1 \\ 1 & 24 \end{vmatrix}}{\begin{vmatrix} 88 & -11 \\ -11 & 32 \end{vmatrix}} = \frac{10(24) - 1(1)}{88(-11) - (-11)(-11)} = \frac{239}{6215} = 0.0385$$

Midterm.

2×2 determinant.

2×2 matrix

③

$p=2, g=3 \Rightarrow$ table $(p=2, g=3), n=8, \alpha=$

Table

Strive to reject

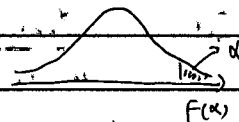
the null \rightarrow variance

$$\Rightarrow \frac{p-g-1}{g-1} \cdot \frac{1-\sqrt{0.0385}}{\sqrt{0.0385}} = 8.19$$

Alt. value:

$$F_{4,5}(0.01) = 7.01$$

$$q^2_{\alpha}(0.99, 4, 8)$$



Converted Statistic

④

Conclusion: Since $8.19 > 7.01$, we reject H_0 @ $\alpha=0.01$ level. reject H_0 the three group doesn't have the same average of 3 parameter

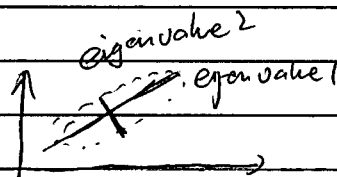
That is, the three groups do not have the same mean vectors

n \rightarrow don't worry about it.

Third version Use eigenvalue for A-test.

det is product of eigenvalues

eigen = direction not change.



~~With statistics~~

$q \rightarrow$ cr. values

$$\lg = \log_{10} x \quad \text{to} \quad \ln = \log_e x$$

Common logarithm

Natural logarithm

All tests: ~~***~~

$\lambda_1, \dots, \lambda_r$ eigenvalues of $W^{-1}B$

Wilks: $\lambda = \prod_{i=1}^r \left(\frac{1}{1 + \lambda_i} \right) \rightarrow \text{Small}$

Ray: $\frac{\lambda_i}{1 + \lambda_i}$, where λ_i is the largest

Pillai: $\sum_{i=1}^r \frac{\lambda_i}{1 + \lambda_i} = \text{tr}[(B+W)^{-1}B]$

Hottelling: $\sum_{i=1}^r \lambda_i = \text{tr}(W^{-1}B)$

① \bar{x}_i, \bar{x}

② $B, W, B+W$

③ $A_{\text{test}} = \frac{|W|}{|B+W|} \rightarrow \text{Wilks' } \lambda$

④ Find:

P.g. $n \Rightarrow T_{0.3}$

Wilks' $\lambda \rightarrow A_{\text{test}}$
Critic. Value

$A_{\text{test}} > \text{Critic. Value}$

Wilks' (lambda)

⑤ Cochran

F \rightarrow Table 6.3

$\chi^2 \rightarrow \text{pg-1}$

χ^2
F

Looking
 χ^2 , F table

Mon
Note
Matrix

Kor

Analysis of Variance (ANOVA)

t-test for more than 2 groups

- ~~variance~~ are [ANOVA]

Suppose we have "g" groups (populations)

From each we have a sample

Group 1 ... Group g

 x_{11} x_{g1} n_1, \dots, n_g x_{12} x_{g2}

01

 \vdots \vdots x_{1n_1} x_{gn_g} Assume: For each $i=1, \dots, g$ \bar{x} $x_{i1}, \dots, x_{in_i} \sim N(\mu_i, \sigma^2)$

Normal

and the groups are independent groups independent.

This set-up is known as one-way ANOVA 1 factor separate groups

It can be written as a probability model

$$x_{ij} = \underbrace{\mu}_{\text{True}} + \underbrace{\tau_i}_{\text{Specific effect of } i\text{th group}} + \underbrace{\epsilon_{ij}}_{\text{Group effect}} = \mu_i + \epsilon_{ij} \quad \begin{matrix} i=1, \dots, g \\ j=1, \dots, n_j \end{matrix} \Rightarrow \text{unobservable.}$$

 μ_i is the population mean of group i ϵ_{ij} are random errors, normally distributed

We want to test:

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

$$H_1: \text{At least two } \mu_i \text{ are different and } \mu \text{ is related to } \tau_i$$

Alternatively, we can test

$$H_0: \tau_1 = \tau_2 = \dots = \tau_g$$

$$\text{where } \sum_{i=1}^g n_i \tau_i = 0$$

Idea

estimate of τ_i

Between & Within Group

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

group

obs

 \Rightarrow observable

observation

overall average

i-th group average

residual/error

no group

Treatment effect

$$\Leftrightarrow (x_{ij} - \bar{x})^2 = (\bar{x}_i - \bar{x})^2 + (x_{ij} - \bar{x}_i)^2 + 2(\bar{x}_i - \bar{x})(x_{ij} - \bar{x}_i)$$

$$\Leftrightarrow \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 = n_i (\bar{x}_i - \bar{x})^2 + \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 + \underbrace{2(\bar{x}_i - \bar{x}) \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)}_{=0}$$

$$\Leftrightarrow \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 = \sum_{i=1}^g n_i (\bar{x}_i - \bar{x})^2 + \underbrace{\sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2}_{\text{adding every single obs}}$$

SS_{Total}
(Corrected)

$SS_{\text{Treatment}}$

(between groups sum of squares)

R output

SS_{Resid} or SSE

(within groups sum of squares)

Group Different \uparrow B W

Scale different. \rightarrow Mean.

ANOVA table:

Source	SS	df	MS	F
Treatment	SS_{Treat}	$g-1$	$SS_{\text{Treat}}/(g-1)$	MS_{Treat}
Residual	SS_{Resid}	$n-g$	$SS_{\text{Resid}}/(n-g)$	MS_{Resid}
Total	SS_{Total}	$n-1$		

rej!!
bigger F ratio stats, null hyp

where $n = n_1 + n_2 + \dots + n_g$

(not a matrix, different length)

$$T.S. F_{\text{test}} = \frac{MS_{\text{Treat}}}{MS_{\text{Resid}}} \sim F_{g-1, n-g}$$

2.R. Reject H_0 if $F_{\text{test}} > F_{g-1, n-g}(\alpha)$

(larger value for F statistics, reject)

$\bar{x}_1 \quad \bar{x}_2$
1 1
1 1
 $x_{11}, x_{12}, x_{13} \quad x_{21}, x_{22}, x_{23}$
?? Noise
 \bar{x}_1, \bar{x}_2
 $x_{12} \dots$

Note

We reject H_0 if SS_{Treat} is "larger" than SS_{Resid}

$$\Leftrightarrow \frac{SS_{\text{Treat}}}{SS_{\text{Resid}}} \uparrow \text{ is "large"}$$

SS Total is

$\Rightarrow 1 + \frac{SS_{Treat}}{SS_{Resid}}$ is "large"

① $p=1, q=n$

② \bar{x}_i, \bar{x}

$\Rightarrow \frac{SS_{Total}}{SS_{Resid}}$ is "large"

③ $SS_T, SS_{Treat}, SS_{Resid}$

④ Table

$\Rightarrow \frac{SS_{Resid}}{SS_{Total}}$ is "small" \rightarrow Analogy to Multivariate

⑤ $F_{q-1, n-q}$

⑥ Δ_{un}

Ex 1) 6.7 on p. 298

Group 1	Group 2	Group 3	$g=3$	(In R: length must be diff aa 1)
9	0	3		
6	2	1		
9		2		
$n_1=3$	$n_2=2$	$n_3=3$		
In R:				

x	Group	$n = 3+2+3=8$	Sum of Squares
9	g_1	$\bar{x}_1 = \frac{9+6+9}{3} = 8$	
6	g_1	$\bar{x}_1 = 8$	
9	g_1	$\bar{x}_1 = 8$	
0	g_2	$\bar{x}_2 = \frac{0+2}{2} = 1$	
2	g_2	$\bar{x}_2 = 1$	
3	g_3	$\bar{x}_3 = \frac{3+1+2}{3} = 2$	
1	g_3	$\bar{x}_3 = 2$	
2	g_3	$\bar{x}_3 = 2$	

Corrected?

• $SS_{Total} = SS - SS_{Mean}$ (Incorrect) $\sum (x_{ij} - \bar{x})^2$
 \downarrow \downarrow
Sum of Squares $n \cdot \bar{x}^2$ $= \sum (x_{ij}^2 - 2\bar{x}x_{ij} + \bar{x}^2)$
 $= 9^2 + 6^2 + \dots + 2^2 - n \cdot \bar{x}^2 = 216 - 128 = 88$ $= \sum x_{ij}^2 - 2\bar{x} \sum x_{ij} + n\bar{x}^2$
 \downarrow \downarrow
 $SS_{Treat} = 3(8-4)^2 + 2(1-4)^2 + 3(2-4)^2$ $n\bar{x}$
 $= 78$

• $SS_{Resid} = 88 - 78 = 10$

$H_0: \mu_1 = \mu_2 = \mu_3$ $\alpha = 0.01$

Source	SS	df	MS	F	P-value
Treatment	78	2	39	$\frac{39}{2} = 19.5$	\geq in R
Resid	10	5	2		
Total	88	7			

$F_{test} = 19.5$

gf \rightarrow critical value -

Crit Value: $F_{2,5}(0.01) = 19.2$

page

Conclusion: Since $9.5 > 13.2$

we reject H_0

That is, at least two groups means are different

In the multivariate case, the "between" and "within" SS are matrices!

$$SS_{\text{Total}} B = \sum_{i=1}^g n_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})'$$

$$W = \sum_{i=1}^g \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)'$$

$\bar{x}_i, \bar{x} \rightarrow p \times 1$ vectors

SS_{Total}

Main idea: W should be "smaller" compared to Total = $B + W$

$$T.S. \frac{\det(W)}{\det(B+W)} \rightarrow \text{"small"} \rightarrow \text{reject } H_0$$

*** Statistic Test \rightarrow only F table

no Hotelling's table

one M
one H_0

genetics: Multivariate + Normal Distribution

MANOVA - Part 3 Equal Covariance Assumption Testing

H4DM b122

2.24

Week 6-1

- Testing for equal covariance matrices

Recall: MANOVA assumptions

$$X_{ij} = \mu + \gamma_i + \epsilon_{ij}, \quad i=1, \dots, g \\ j=1, \dots, n_i$$

① g, n

Assume: $\epsilon_{ij} \sim N_p(0, \Sigma)$

② within each

some covariance for all groups

g, N

Groups are independent

$$H_0: \gamma_1 = \dots = \gamma_g = 0$$

③ Same

$$H_0: \Sigma_1 = \Sigma_2 = \dots = \Sigma_g$$

where Σ_i is the covariance of group i

Box's M test:

$$M = \left[\sum_{i=1}^g (n_i - 1) \right] \ln(\det(S_{pool})) - \sum_{i=1}^g (n_i - 1) \ln(\det(S_i))$$

Define:

$$C = \frac{1}{M} \left[\sum_{i=1}^g \frac{1}{n_i - 1} \right] \frac{2p^2 + 3p - 1}{6(p+1)(g-1)}$$

p, g, n

Under H_0 ,

$$C = (1 - U) M \approx \chi^2_U$$

$$\text{where df. } U = \frac{p(p+1)(g-1)}{2}$$

Newton
Iteration

From

Likelihood Ratio test.

Under H_0 ,

RR:

Distribution

Reject H_0 if
 $C > \chi^2_U(\alpha)$

Remark: The test works well if $n_i > 20, i=1, \dots, g$

and $p, g \leq 5$

Remark: Very sensitive to violations of normality and as a result produces very small p-values

High sensitivity
Reject Null low power