

Reading: 2.5 & 2.6

A random variable is a variable that assigns a numerical value to each outcome in the sample space.

The values of a r.v are determined by the outcomes of a random experiment.

Ex

Toss a die once

r.v: $X = \# \text{ of dots on the die}$
possible values: $S = \{1, 2, 3, 4, 5, 6\}$.

* A function that maps each outcome in a sample space to a real number line.

Random Vector

- A vector whose elements are random variables X_1, X_2, X_3
- Eq Biology: (weight, height, blood pressure) of an individual of the same day for different locations at the same time
- ② Spatial data: Record the temperature at different locations at the same time for different days $I = (t_1, t_2, \dots, t_5)$
- ③ Survival data: Failure time ($\text{time}_1, \text{time}_2, \text{time}_3$) of different objects.

A random vector $\underline{X} = (X_1, X_2, \dots, X_p)_{px1}$ and each element of X_i is a r.v on its own with its own marginal distribution.

- Define the following marginal parameters:
 - ① $E(X_i) = \begin{cases} \int_{-\infty}^{\infty} x_i f_i(x_i) dx_i & \text{if } X_i \text{ is a continuous r.v with pdf } f_i(x_i) \\ \sum_{\text{all } x} x_i p_i(x_i) & \text{if } X_i \text{ is discrete pmf } p_i(x_i) \end{cases}$
 - ② $V(X_i) = E[(X_i - \mu_i)^2] = \sigma_{ii} = \sigma_i^2$
 - ③ For 2 r.variables X_i, X_j , $\text{Cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)] = \sigma_{ij}$; if $i \neq j$
- A measure of linear association between the variables X_i, X_j .
- If X_i and X_j are independent, $\text{Cov}(X_i, X_j) = 0$. The converse is generally not true except in the case of Multivariate normal distribution
- $\text{Cov}(X_i, X_j) = \text{Cov}(X_j, X_i)$
- $\sigma_{ij} = \sigma_{ji}$
- The behaviour of any pair of random variables, such as X_i, X_j is described by their joint pdf.

Revise joint distributions.

* For a random vector $\underline{X} = (X_1, X_2, \dots, X_p)_{px1}$

$$\textcircled{1} \quad \underline{\mu} = E(\underline{X}) = \begin{bmatrix} E(X_1) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_p \end{bmatrix} = \text{population mean vector} \quad \underline{\mu} = \sum \underline{x}$$

\textcircled{2} Population Variance Covariance matrix

$$\sum = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \dots & \sigma_{pp} & \end{pmatrix}_{p \times p} = E((\underline{x} - \underline{\mu})(\underline{x} - \underline{\mu})^T) = E \begin{bmatrix} (x_1 - \mu_1)^2 & (x_1 - \mu_1)(x_2 - \mu_2) & \dots & (x_1 - \mu_1)(x_p - \mu_p) \\ (x_2 - \mu_2)(x_1 - \mu_1) & (x_2 - \mu_2)^2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ (x_p - \mu_p)(x_1 - \mu_1) & \dots & \dots & (x_p - \mu_p)^2 \end{bmatrix}$$

$$= \left[E(X_1 - \mu_1)^2 \right]$$

- In a dataset of p variables, there will be p variances and $\frac{p(p-1)}{2}$ distinct covariances.

③ From the Σ , we can obtain the correlation matrix as

$$\rho = \begin{bmatrix} \frac{\sigma_{11}}{\sqrt{\sigma_{11}\sigma_{11}}} & \frac{\sigma_{12}}{\sqrt{\sigma_{11}\sigma_{22}}} & \dots & \frac{\sigma_{1p}}{\sqrt{\sigma_{11}\sigma_{pp}}} \\ \vdots & \ddots & & \vdots \\ \frac{\sigma_{p1}}{\sqrt{\sigma_{pp}\sigma_{11}}} & \dots & \frac{\sigma_{pp}}{\sqrt{\sigma_{pp}\sigma_{pp}}} & \end{bmatrix} = \begin{bmatrix} 1 & \rho_{12} & \dots & \rho_{1p} \\ \vdots & 1 & \dots & \vdots \\ \rho_{p1} & \dots & \dots & 1 \end{bmatrix}$$

Note that if $V^{1/2} = \begin{pmatrix} \sqrt{\sigma_{11}} & 0 & \dots & 0 & 0 \\ 0 & \sqrt{\sigma_{22}} & & & \\ \vdots & & \ddots & & \\ 0 & \dots & 0 & \sqrt{\sigma_{pp}} & \end{pmatrix}_{pxp}$ then

↑ Standard deviation matrix

$$\rho = V^{-1/2} \sum V^{1/2} = (V^{1/2})^{-1} \sum (V^{1/2})^{-1}$$

and

$$\Sigma = V^{1/2} \rho V^{1/2}$$

Example 1

Find the covariance matrix for the random vector (X_1, X_2) , when their pmf $P_a(x_1, x_2)$ is represented by

				$P_a(x_1)$
$P_a(x_2)$	-1	0.24	0.06	0.3
	0	0.16	0.14	0.3
	1	0.4	0.00	0.4
		0.8	0.2	1

Solution.

$$E(X_1) = -1(0.3) + 0(0.3) + 1(0.4) = 0.1$$

$$E(X_2) = 0(0.8) + 1(0.2) = 0.2$$

$$\sigma_{11} = E(X_1 - \mu_1)^2 = \sum_{x_1 \in X_1} (x_1 - \mu_1)^2 P_a(x_1) = (-1-0.1)^2(0.3) + (0-0.1)^2(0.3) + (1-0.1)^2(0.4) = 0.69$$

$$\sigma_{22} = E(X_2 - \mu_2)^2 = \sum_{x_2 \in X_2} (x_2 - \mu_2)^2 P_a(x_2) = (0-0.2)^2(0.8) + (1-0.2)^2(0.2) = 0.16$$

$$\begin{aligned} \sigma_{12} = \sigma_{21} &= E[(X_1 - \mu_1)(X_2 - \mu_2)] = \sum_{\substack{\text{all pairs} \\ (x_1, x_2)}} (x_1 - \mu_1)(x_2 - \mu_2) P_a(x_1, x_2) \\ &= (-1-0.1)(0-0.2)(0.24) + (-1-0.1)(1-0.2)(0.06) + \dots + (1-0.1)(1-0.2)(0.00) \\ &= -0.08 \end{aligned}$$

$$\text{So } \underline{\mu} = E(\underline{x}) = \begin{pmatrix} E(X_1) \\ E(X_2) \end{pmatrix} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \begin{pmatrix} 0.1 \\ 0.2 \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix} = \begin{pmatrix} 0.69 & -0.08 \\ -0.08 & 0.16 \end{pmatrix}$$

Example 2

Given the covariance matrix below, find $\sqrt{V^{1/2}}$ and ρ .

$$\Sigma = \begin{bmatrix} 4 & 1 & 2 \\ 1 & 9 & -3 \\ 2 & -3 & 25 \end{bmatrix}$$

Solution

$$V^{1/2} = \begin{bmatrix} \sqrt{\sigma_{11}} & 0 & 0 \\ 0 & \sqrt{\sigma_{22}} & 0 \\ 0 & 0 & \sqrt{\sigma_{33}} \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 5 \end{bmatrix}$$

$$(V^{1/2})^{-1} = \begin{bmatrix} 1/2 & 0 & 0 \\ 0 & 1/3 & 0 \\ 0 & 0 & 1/5 \end{bmatrix}$$

$$\begin{aligned} \rho &= (V^{1/2})^{-1} \sum (V^{1/2})^{-1} \\ &= \begin{bmatrix} 1/2 & 0 & 0 \\ 0 & 1/3 & 0 \\ 0 & 0 & 1/5 \end{bmatrix} \begin{bmatrix} 4 & 1 & 2 \\ 1 & 9 & -3 \\ 2 & -3 & 25 \end{bmatrix} \begin{bmatrix} 1/2 & 0 & 0 \\ 0 & 1/3 & 0 \\ 0 & 0 & 1/5 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 1/6 & 1/5 \\ 1/6 & 1 & -1/5 \\ 1/5 & -1/5 & 1 \end{bmatrix} \end{aligned}$$

Random Matrix

A random matrix is a matrix whose elements are random variables.

Define

$$\bar{X}_{n \times p} = \{X_{ij}\} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1p} \\ \vdots & & & \\ X_{n1} & \dots & & X_{np} \end{pmatrix}_{n \times p}$$

a) $E(\bar{X}) = \{E(X_{ij})\}$ if each expectation exists.

b) Let \bar{X} and \bar{Y} be random matrices of the same dimension and let \bar{A} and \bar{B} be be $n \times m$ or $m \times n$ matrices of constants

$E(\bar{X}^+) = E(\bar{X}) + E(\bar{Y})$ - Linear property of expectation

$$E(\bar{A} \bar{X} \bar{B}) = \underbrace{\bar{A} E(\bar{X}) \bar{B}}_{K \times L} \quad (\text{See example 2.39, 2.40})$$

Partitioning random vectors

Given a large dataset, characteristics measured on individual trials could naturally fall into two or more groups.

Example 1: Consumption and income

Imagine a survey conducted on a group of individuals, where you collect data on both consumption habits - spending on groceries, entertainment and utilities, while income would encompass variables like salary, investments and other sources of earnings.

Categorize into 2 groups: consumption habits \times income.

Example 2 : Personality traits and physical characteristics.

Variables measured : Extraversion, agreeableness, conscientiousness, height, weight, body mass index. \sum^2 groups
 a) personality traits
 b) physical characteristics

To handle situations like these, where distinct groups of characteristics arise naturally, a useful approach is to consider these groups as subsets of a larger collection of characteristics.

- If the total collection is represented by a $p \times 1$ dimension random vector \tilde{X}_{px} . We can regard the subsets as components of \tilde{X} and can be dealt with by partitioning

$$\tilde{X} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_q \\ - \\ X_{q+1} \\ \vdots \\ X_p \end{bmatrix} \begin{array}{l} \text{q}x1 \\ \text{(p-q)x1} \end{array} = \begin{bmatrix} \tilde{X}^{(1)} \\ \tilde{X}^{(2)} \end{bmatrix}$$

② Then $\tilde{\mu} = E(\tilde{X}) = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_q \\ \dots \\ \mu_{q+1} \\ \vdots \\ \mu_p \end{bmatrix} = \begin{bmatrix} \tilde{\mu}^{(1)} \\ \tilde{\mu}^{(2)} \end{bmatrix}$

$$\begin{aligned} b) \quad \Sigma_{pxp} &= E[(\tilde{X} - \tilde{\mu})(\tilde{X} - \tilde{\mu})^T] = E \left[\begin{array}{ccc} (\tilde{X}^{(1)} - \tilde{\mu}^{(1)}) & (\tilde{X}^{(1)} - \tilde{\mu}^{(1)})^T & \tilde{X}^{(2)} - \tilde{\mu}^{(2)} \\ (\tilde{X}^{(2)} - \tilde{\mu}^{(2)}) & & \end{array} \right] \\ &= E \left[\begin{array}{cccc} (\tilde{X}^{(1)} - \tilde{\mu}^{(1)}) & (\tilde{X}^{(1)} - \tilde{\mu}^{(1)})^T & & \\ & \ddots & \ddots & \\ & & \ddots & \\ & & & (\tilde{X}^{(2)} - \tilde{\mu}^{(2)})^T \end{array} \right] \\ &= \begin{array}{c} q \\ \hline \sum_{11} \\ p-q \end{array} \left[\begin{array}{cc} \sum_{11}^{qq} & \sum_{12}^{qq} \\ \sum_{21}^{qq} & \sum_{22}^{(p-q)(p-q)} \end{array} \right] \end{aligned}$$

Where \sum_{11}^{qq} = Variance covariance matrix for $\tilde{X}^{(1)}$
 $\sum_{22}^{(p-q)(p-q)}$ = " $\tilde{X}^{(2)}$

$$\sum_{12}^{qq} = E[(\tilde{X}^{(1)} - \tilde{\mu}^{(1)})(\tilde{X}^{(2)} - \tilde{\mu}^{(2)})^T] = \sum_{12}^{T}$$

Linear Combinations of random vectors

- Let X_1 be a single random variable and c be a constant with mean and variance
 - a) $E(cX_1) = c\mu_1$
 - b) $V(cX_1) = c^2\sigma_{11} = c^2V(X_1)$

with $E(X_1) = \mu_1, E(X_2) = \mu_2, V(X_1) = \sigma_{11}, V(X_2) = \sigma_{22}, \text{Cov}(X_1, X_2) = \sigma_{12}$

- For 2 random variables X_1 and X_2 , and for constants a and b ,

(a) $E(aX_1 + bX_2) = a\mu_1 + b\mu_2$ Note $aX_1 + bX_2$ is a linear combination of X_1 and X_2

(b) $\text{Cov}(aX_1, bX_2) = ab \text{Cov}(X_1, X_2) = ab\sigma_{12}$

Proof:

$$\textcircled{c} \quad V(aX_1 + bX_2) = a^2V(X_1) + b^2V(X_2) + 2ab\text{Cov}(X_1, X_2)$$

$$= a^2\sigma_{11} + b^2\sigma_{22} + 2ab\sigma_{12}$$

Prove in HW 2

- $aX_1 + bX_2 = (a \ b) \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} = \underline{c}^T \underline{X}$

$$\rightarrow E(aX_1 + bX_2) = E\left[\left(a \ b\right) \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}\right] = (a \ b) \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \underline{c}^T \underline{\mu}$$

$$\rightarrow \text{If we let } \Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}$$

then

$$\begin{aligned} V(aX_1 + bX_2) &= V(\underline{c}^T \underline{X}) \\ &= V\left[(ab) \begin{pmatrix} X_1 \\ X_2 \end{pmatrix}\right] \\ &= \underline{c}^T \Sigma \underline{c} \\ &= (ab) \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} \\ &= a^2\sigma_{11} + 2ab\sigma_{12} + b^2\sigma_{22} \end{aligned}$$

- We can extend this result for q linear combinations of p variables

$$\underline{c}^T \underline{X} = c_1X_1 + c_2X_2 + \dots + c_pX_p$$

$$\text{Then } E(\underline{c}^T \underline{X}) = \underline{c}^T \underline{\mu} ; \underline{\mu} = E(\underline{X}) \text{ and } \Sigma = \text{Cov}(\underline{X})$$

$$V(\underline{c}^T \underline{X}) = \underline{c}^T \Sigma \underline{c}$$

- Even more generally, we can consider q linear combinations of p random variables

$$X_1, \dots, X_p$$

$$Z_1 = c_{11}X_1 + c_{12}X_2 + \dots + c_{1p}X_p$$

$$Z_2 =$$

⋮

$$Z_q = c_{q1}X_1 + c_{q2}X_2 + \dots + c_{qp}X_p$$

So that

$$\underline{Z} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & & & \\ c_{q1} & c_{q2} & \dots & c_{qp} \end{bmatrix} \begin{pmatrix} X_1 \\ \vdots \\ X_p \end{pmatrix} = \underline{C} \underline{X}$$

$$\text{Mean of } \tilde{z} = E(\tilde{z}) = E(\bar{c}x) = \bar{c}\mu$$

$$\sum_z = \text{Cov}(z) = \text{Cov}(\bar{c}x) = \bar{c} \sum_x \bar{c}^T$$

$$(AB)^T = B^T A^T$$

Proof

$$\begin{aligned} \text{Cov}(\bar{c}x) &= E[(\bar{c}x - \bar{c}\mu)(\bar{c}x - \bar{c}\mu)^T] \\ &= E[\bar{c}(x - \mu)(x - \mu)^T c^T] \\ &= \bar{c}E[(x - \mu)(x - \mu)^T] c^T \\ &= \bar{c} \sum_x c^T \end{aligned}$$

* We will rely on these results when discussing principal component analysis and factor analysis.

Example

Find μ_z and \sum_z for $\tilde{z} = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix}$ if

$$\begin{aligned} z_1 &= x_1 - x_2 & \text{if } \tilde{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \mu_x = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \\ z_2 &= x_1 + x_2 \end{aligned}$$

$$\sum_x = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix}$$

Solution

$$\tilde{z} = \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$E(\tilde{z}) = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix} = \begin{pmatrix} \mu_1 - \mu_2 \\ \mu_1 + \mu_2 \end{pmatrix}$$

$$\sqrt{(\tilde{z})} = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{pmatrix} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$

$$= \begin{pmatrix} \sigma_{11} - \sigma_{12} & \sigma_{12} - \sigma_{22} \\ \sigma_{11} + \sigma_{12} & \sigma_{12} + \sigma_{22} \end{pmatrix} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$

$$= \begin{pmatrix} \sigma_{11} - \sigma_{12} - \sigma_{12} + \sigma_{22} & \sigma_{11} - \sigma_{12} + \sigma_{12} - \sigma_{22} \\ \sigma_{11} + \sigma_{12} - \sigma_{12} - \sigma_{22} & \sigma_{11} + \sigma_{12} + \sigma_{12} + \sigma_{22} \end{pmatrix}$$

$$\begin{pmatrix} \sigma_{11} - 2\sigma_{12} + \sigma_{22} & \sigma_{11} - \sigma_{22} \\ \sigma_{11} - \sigma_{22} & \sigma_{11} + 2\sigma_{12} + \sigma_{22} \end{pmatrix}$$

Note

Just as we partitioned the population variance covariance matrix and population mean vector, we can also partition the sample mean vector and sample covariance matrix the same way to distinguish quantities corresponding to groups of variables. Thus

Let

$$\bar{\underline{x}}_{px1} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_q \\ \bar{x}_{q+1} \\ \vdots \\ \bar{x}_p \end{pmatrix}_{p \times 1}^q = \begin{bmatrix} \bar{\underline{x}}^{(1)} \\ \bar{\underline{x}}^{(2)} \\ \vdots \\ \bar{\underline{x}}^{(p)} \end{bmatrix}_p^q$$

be a vector of sample averages constructed from n observations on p variables x_1, \dots, x_p .

and let

$$\bar{S}_n = \begin{bmatrix} S_{11} & S_{12} & S_{13} & \dots & S_{1p} \\ S_{21} & S_{22} & S_{23} & \dots & S_{2p} \\ S_{31} & S_{32} & S_{33} & \dots & S_{3p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ S_{p1} & S_{p2} & S_{p3} & \dots & S_{pp} \end{bmatrix}$$

be the corresponding sample variance covariance matrix.

matrix:

Section 3.3 Sample Statistics.

$$\text{Define } \bar{\underline{x}}_{n \times p} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} =$$

- If the row vectors $\bar{\underline{x}}_1^T, \bar{\underline{x}}_2^T, \dots, \bar{\underline{x}}_n^T$ represent independent observations from a common joint distribution with density $f(\underline{x}) = f(x_1, \dots, x_p)$, then x_1, x_2, \dots, x_n are said to form a random sample from $f(\underline{x})$.
 $f(\underline{x}) = f(x_1) f(x_2) \dots f(x_n)$ where $f(x_j) = f(x_{j1}, x_{j2}, \dots, x_{jp})$ is the density function for the j th row.
- Let x_1, \dots, x_n be a random sample from a joint distribution that has mean μ and variance cov matrix Σ .

① Define $\bar{\underline{x}}_{px1} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_p \end{pmatrix} = \frac{1}{n} \sum_{i=1}^n \bar{\underline{x}}_i$ = Sample mean vector

ⓐ Then $\bar{\underline{x}}_{px1}$ is an unbiased estimator of μ that is,

$$E(\bar{\underline{x}}_{px1}) = \underline{\mu}_{px1}$$

ⓑ $Cov(\bar{\underline{x}}_{px1}) = \frac{1}{n} \sum_x$ (ie the population variance covariance matrix divided by sample size)

② Define the sample covariance matrix as S_n

$$S_n = \begin{pmatrix} \frac{1}{n} \sum_{j=1}^n (x_{j1} - \bar{x}_1)^2 & \dots & \frac{1}{n} \sum_{j=1}^n (x_{j1} - \bar{x}_1)(x_{jp} - \bar{x}_p) \\ \vdots & \ddots & \vdots \\ \frac{1}{n} \sum_{j=1}^n (x_{j1} - \bar{x}_p)(x_{jp} - \bar{x}_p) & \dots & \frac{1}{n} \sum_{j=1}^n (x_{ji} - \bar{x}_p)^2 \end{pmatrix}$$

ⓐ $E(S_n) = \frac{n-1}{n} \sum = \Sigma - \frac{1}{n} \sum$ so that

$E\left(\frac{n}{n-1} S_n\right) = \sum$ so $\frac{n}{n-1} S_n$ is an unbiased estimator of Σ

$\frac{1}{n-1} \sum_{j=1}^n (x_{ji} - \bar{x}_p)(x_{jp} - \bar{x}_p)$ for the n th row.

Sample mean vector, covariance and correlation as matrix operations

Define $\bar{X}_{n \times p} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & & & \\ \vdots & & & \\ x_{n1} & & & x_{np} \end{pmatrix}$

$$\bar{X}^T = \begin{pmatrix} x_{11} & x_{21} & \dots & x_{n1} \\ x_{12} & & & \\ \vdots & & & \\ x_{1p} & \dots & \dots & x_{np} \end{pmatrix}$$

$$① \bar{x}_{px1} = \frac{1}{n} \bar{X}^T \mathbf{1}$$

$$② \text{Recall } S = \frac{n-1}{n} S_n \quad \xrightarrow{\bar{X}^T A \bar{X} \text{ quadratic forms}}$$

$$= \frac{1}{n-1} \bar{X}^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \bar{X}$$

where S_n is the biased sample covariance matrix.

$$\frac{1}{n} \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & & & \\ \vdots & & & \\ x_{n1} & & & x_{np} \end{pmatrix} \left(I - \frac{1}{n} \begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix}^T \right) \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & & & \\ \vdots & & & \\ x_{n1} & & & x_{np} \end{pmatrix} = \begin{pmatrix} \frac{1}{n-1} \sum (x_{ii} - \bar{x}_i)^2 \\ \vdots \\ \frac{1}{n-1} \sum (x_{ij} - \bar{x}_i)(x_{ik} - \bar{x}_k) \\ \vdots \\ \frac{1}{n-1} \sum (x_{ii} - \bar{x}_i)^2 \end{pmatrix}$$

$$③ R = D^{-1/2} S_n D^{-1/2}$$

$$D^{1/2} = \begin{pmatrix} \sqrt{s_{11}} & 0 & \dots & 0 \\ 0 & \sqrt{s_{22}} & & \\ \vdots & & \ddots & \\ 0 & 0 & \dots & \sqrt{s_{nn}} \end{pmatrix}$$

$$\bar{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & & & \\ \vdots & & & \\ x_{n1} & & & x_{np} \end{pmatrix} - \frac{1}{n} \bar{X}^T \mathbf{1} \mathbf{1}^T \begin{pmatrix} \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \\ \bar{x}_1 & \bar{x}_2 & & \bar{x}_p \\ \bar{x}_1 & & \ddots & \\ \vdots & & & \bar{x}_p \\ \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_p \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_1 & \dots & x_{1p} - \bar{x}_1 \\ x_{21} - \bar{x}_1 & & & \\ \vdots & & & \\ x_{n1} - \bar{x}_1 & & & x_{np} - \bar{x}_p \end{pmatrix}$$

$$\begin{pmatrix} \bar{X} - \frac{1}{n} \bar{X}^T \mathbf{1} \mathbf{1}^T \end{pmatrix}^T \begin{pmatrix} \bar{X} - \frac{1}{n} \bar{X}^T \mathbf{1} \mathbf{1}^T \end{pmatrix}$$

$$= \bar{X}^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \bar{X}$$

$$(n-1) S = \bar{X}^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \bar{X}$$

If $A^T A = A$, A is an idempotent matrix

Sample values of linear combinations of variables

Consider a linear combination of vectors:

$C^T \bar{X} = c_1 \bar{X}_1 + c_2 \bar{X}_2 + \dots + c_p \bar{X}_p$ where $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p$ are independent realizations from the random vector.

$$\text{Recall } E(\tilde{\zeta}^T \tilde{x}) = c^T \mu \quad \text{Var}(\tilde{\zeta}^T \tilde{x}) = \tilde{\zeta}^T \Sigma \tilde{\zeta}$$

Suppose we want to estimate these,

$$\text{Let } Y_i = \tilde{\zeta}^T \tilde{x}_i$$

① An estimate of the mean $\tilde{\zeta}^T \mu$ is given by

$$\frac{1}{n} \sum_{i=1}^n Y_i = \frac{1}{n} \sum_{i=1}^n \tilde{\zeta}^T \tilde{x}_i = \tilde{\zeta}^T \left(\frac{1}{n} \sum_{i=1}^n \tilde{x}_i \right) = \tilde{\zeta}^T \bar{\tilde{x}}$$

② An estimate of the variance $\tilde{\zeta}^T \Sigma \tilde{\zeta}$ is given by

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2 &= \frac{1}{n} \sum_{i=1}^n (\tilde{\zeta}^T \tilde{x}_i - \tilde{\zeta}^T \bar{\tilde{x}})^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\tilde{\zeta}^T (\tilde{x}_i - \bar{\tilde{x}}))^2 \\ &= \frac{1}{n} \sum_{i=1}^n (\tilde{\zeta}^T (\tilde{x}_i - \bar{\tilde{x}})(\tilde{x}_i - \bar{\tilde{x}})^T \tilde{\zeta}) \\ &= \tilde{\zeta}^T \frac{1}{n} \sum_{i=1}^n (\tilde{x}_i - \bar{\tilde{x}})(\tilde{x}_i - \bar{\tilde{x}})^T \tilde{\zeta} \\ &= \tilde{\zeta}^T S_n \tilde{\zeta} \end{aligned}$$

③ An estimate for the population covariance matrix of $b^T \tilde{x}$, $\tilde{\zeta}^T \tilde{x}$ which is $\text{Cov}(b^T \tilde{x}, \tilde{\zeta}^T \tilde{x}) = b^T \Sigma \tilde{\zeta}$ is given by $\tilde{\zeta}^T S_n b$.

Example Given $\bar{\tilde{x}} = \begin{pmatrix} x_1 & x_2 & x_3 \\ 1 & 2 & 5 \\ 4 & 1 & 6 \\ 4 & 0 & 4 \end{pmatrix}_{3 \times 3}$. Find the following:

a) $\bar{x}_{3 \times 1}$

b) S_n

c) Consider the linear combinations $2x_1 + 2x_2 - x_3$ and $x_1 - x_2 + 3x_3$

i) Find their means and variances respectively.

ii) Determine $\text{Cov}(2x_1 + 2x_2 - x_3, x_1 - x_2 + 3x_3)$

Solution

a) $\bar{x}_{3 \times 1} = \frac{1}{n} \bar{\tilde{x}}^T \mathbf{1} = \frac{1}{3} \begin{pmatrix} 1 & 4 & 4 \\ 2 & 1 & 0 \\ 5 & 6 & 4 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \frac{1}{3} \begin{pmatrix} 9 \\ 3 \\ 15 \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \\ 5 \end{pmatrix}$

b) $S_n = \frac{n-1}{n} S = \frac{1}{n} \bar{\tilde{x}}^T (I - \frac{1}{n} \mathbf{1} \mathbf{1}^T) \bar{\tilde{x}}$

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

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

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

$$\begin{pmatrix} 2/3 & -1/3 & -1/3 \\ -1/3 & 2/3 & -1/3 \\ -1/3 & -1/3 & 2/3 \end{pmatrix}$$

$$c) \text{ Let } \underline{b}^T \underline{x} = 2x_1 + 2x_2 - x_3 = (2 \ 2 \ -1) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \text{ and } \underline{c}^T \underline{x} = x_1 - x_2 + 3x_3 = (1 \ -1 \ 3) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

$$\begin{aligned} \text{Sample mean of } \underline{b}^T \underline{x} &= \underline{b}^T \bar{\underline{x}} \\ &= (2 \ 2 \ -1) \begin{pmatrix} 3 \\ 5 \end{pmatrix} = 3 \end{aligned}$$

$$\begin{aligned} \text{II } \underline{c}^T \underline{x} &= \underline{c}^T \bar{\underline{x}} \\ &= (1 \ -1 \ 3) \begin{pmatrix} 3 \\ 5 \end{pmatrix} = 17 \end{aligned}$$

$$\begin{aligned} \text{Sample variance of } \underline{b}^T \underline{x} &= \underline{b}^T S_n \underline{b} \\ &= (2 \ 2 \ -1) \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2/3 & 1/3 \\ 0 & 1/3 & 2/3 \end{pmatrix} \begin{pmatrix} 2 \\ 2 \\ -1 \end{pmatrix} \\ &= (2 \ 2 \ -1) \begin{pmatrix} 2 \\ -1 \\ 0 \end{pmatrix} = 2 \end{aligned}$$

$$\begin{aligned} \text{Sample variance of } \underline{c}^T \underline{x} &= \underline{c}^T S_n \underline{c} \\ &= (1 \ -1 \ 3) \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2/3 & 1/3 \\ 0 & 1/3 & 2/3 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \\ 3 \end{pmatrix} = (1 \ -1 \ 3) \begin{pmatrix} 3 \\ -2/3 \\ 5/3 \end{pmatrix} \\ &= 26/3 \end{aligned}$$

Sample Cov of $\underline{b}^T \underline{x}$ and $\underline{c}^T \underline{x}$

$$\begin{aligned} \text{Cov}(\underline{b}^T \underline{x}, \underline{c}^T \underline{x}) &= \underline{b}^T S_n \underline{c} \\ &= (2 \ 2 \ -1) \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2/3 & 1/3 \\ 0 & 1/3 & 2/3 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \\ 3 \end{pmatrix} \\ &= (2 \ 2 \ -1) \begin{pmatrix} 3 \\ -2/3 \\ 5/3 \end{pmatrix} \\ &= 3 \end{aligned}$$

- The results discussed for a single linear combination can be extended to any number of linear combinations.

- Consider q linear combinations which is expressed in matrix notation as

$$\begin{pmatrix} c_{11} \underline{x}_1 + c_{12} \underline{x}_2 + \dots + c_{1p} \underline{x}_p \\ c_{21} \underline{x}_1 + \dots \\ \vdots \\ c_{q1} \dots + c_{qp} \underline{x}_p \end{pmatrix} = \begin{pmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & & & \\ c_{q1} & \dots & & c_{qp} \end{pmatrix} \begin{pmatrix} \underline{x}_1 \\ \underline{x}_2 \\ \vdots \\ \underline{x}_p \end{pmatrix} = \underline{C} \underline{x}$$

Let the i th row of \underline{C} be \underline{c}_i^T and the j th row be \underline{c}_j^T . From previous results, the i th row has sample mean $\underline{c}_i \bar{\underline{x}}$ and the i,j th element of the sample variance covariance, $\underline{c}_i^T S_n \underline{c}_j$.

- $\underline{C} \underline{x}$ will have sample mean $\underline{C} \bar{\underline{x}}$ and sample covariance $\underline{C} S_n \underline{C}^T$

Ch 4

Univariate normal distribution

Density:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad -\infty < x < \infty$$

CDF

$$P(X \leq x) = \Phi(x) = \int_{-\infty}^x \phi(t) dt$$

Properties

It is symmetric so $\Phi(-x) = 1 - \Phi(x)$ for every value x .

Visualize normality:

histogram, qqplot.

Tests for normality

Shapiro wilk, Kolmogorov smirnov

H₀: Data is normally distributed

H_a: Data is not normally distributed.

Multivariate normal

- Many multivariate techniques to be discussed in this class depend on the MVN assumption
- $p \geq 2$ dimensions
- Density function for a p -dimensional MVN random vector $X = (X_1, \dots, X_p)$, has the form

$$f(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}; \quad -\infty < x_i < \infty \text{ for } i=1, \dots, p$$

We say that $X \sim N_p(\mu, \Sigma)$

- When $p=2$, we have a bivariate normal distribution

Given a random vector $X = (X_1, X_2)$, let $E(X_1) = \mu_1, E(X_2) = \mu_2, \Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \sigma_{22} \end{pmatrix}$. Suppose Σ is invertible ie,

$$\Sigma^{-1} = \frac{1}{\sigma_{11}\sigma_{22} - \rho_{12}^2} \begin{pmatrix} \sigma_{22} & -\rho_{12} \\ -\rho_{12} & \sigma_{11} \end{pmatrix}. \text{ We know } \rho_{12} = \frac{\sigma_{12}}{\sqrt{\sigma_{11}} \sqrt{\sigma_{22}}} \text{ so } \sigma_{12} = \rho_{12} \sqrt{\sigma_{11}} \sqrt{\sigma_{22}}$$

$$= \frac{1}{\sigma_{11}\sigma_{22} - \rho_{12}^2 \sigma_{11}\sigma_{22}} \begin{pmatrix} \sigma_{22} - \rho_{12} & 0 \\ 0 & \sigma_{11} \end{pmatrix}$$

We can show that the density function for the bivariate normal is given by

$$f(x_1, x_2) = \frac{1}{2\pi\sqrt{\sigma_{11}\sigma_{22}(1-\rho_{12}^2)}} \exp \left\{ -\frac{1}{2(1+\rho_{12})^2} \left[\left(\frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} \right)^2 + \left(\frac{x_2 - \mu_2}{\sqrt{\sigma_{22}}} \right)^2 - 2\rho_{12} \left(\frac{x_1 - \mu_1}{\sqrt{\sigma_{11}}} \right) \left(\frac{x_2 - \mu_2}{\sqrt{\sigma_{22}}} \right) \right] \right\}$$

Homework:

Show that when X_1 and X_2 are uncorrelated, this joint density can be written as the product of 2 univariate normal densities.

Properties of MVN distribution

Suppose $\underline{X} \sim N_p(\underline{\mu}_{px}, \Sigma)$, then

- ① Any linear combination of variables $\underline{a}^T \underline{X} = a_1 X_1 + a_2 X_2 + \dots + a_p X_p$ is said to follow $N(a^T \underline{\mu}, a^T \Sigma a)$
- ② If $a^T \underline{X} \sim N(a^T \underline{\mu}, a^T \Sigma a)$ for every a , then \underline{X} must be $N_p(\underline{\mu}, \Sigma)$
- ③ For q linear combinations $\underline{\underline{A}}^T \underline{X}_{px}$, if $\underline{X} \sim N_p(\underline{\mu}, \Sigma)$,

$$\underline{\underline{A}}^T \underline{X} \sim N_p(\underline{\underline{A}}^T \underline{\mu}, \underline{\underline{A}}^T \Sigma \underline{\underline{A}})$$

- ④ $\underline{X}_{px} + \underline{d}_{px} \sim N_p(\underline{\mu} + \underline{d}, \Sigma)$ where \underline{d} is a vector of constants

- ⑤ All subsets of \underline{X} are multivariate normally distributed
Eg if we partition \underline{X} , we know that its mean vector and covariance matrix will be

$$\underline{X}_{px} = \begin{bmatrix} \underline{X}^{(1)} \\ \vdots \\ \underline{X}^{(q)} \end{bmatrix}_{p \times q}, \quad \underline{\mu} = \begin{bmatrix} \underline{\mu}^{(1)} \\ \vdots \\ \underline{\mu}^{(q)} \end{bmatrix}_{p \times q}, \quad \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}_{(p-q) \times q \times (p-q)}$$

then $\underline{X}^{(1)} \sim N_q(\underline{\mu}^{(1)}, \Sigma_{11})$

Independence:

⑥ If \tilde{X}_1 and \tilde{X}_2 are independent then $\text{Cov}(\tilde{X}_1, \tilde{X}_2) = \mathbf{0}_{q_1 \times q_2}$

⑦ If $\tilde{X} = \begin{bmatrix} \tilde{X}_1 \\ \tilde{X}_2 \end{bmatrix}_{q_1+q_2} \sim N_{q_1+q_2}\left(\begin{bmatrix} \tilde{\mu}_1 \\ \tilde{\mu}_2 \end{bmatrix}, \begin{bmatrix} \tilde{\Sigma}_{11} & \tilde{\Sigma}_{12} \\ \tilde{\Sigma}_{21} & \tilde{\Sigma}_{22} \end{bmatrix}\right)$ then \tilde{X}_1 and \tilde{X}_2 are independent iff $\tilde{\Sigma}_{12} = \tilde{\Sigma}_{21} = \mathbf{0}_{q_1 \times q_2}$

⑧ If \tilde{X}_1 and \tilde{X}_2 are independent and distributed as $N_{q_1}(\tilde{\mu}_1, \tilde{\Sigma}_{11})$ and $N_{q_2}(\tilde{\mu}_2, \tilde{\Sigma}_{22})$ respectively, then

$$\begin{pmatrix} \tilde{X}_1 \\ \tilde{X}_2 \end{pmatrix} \sim N_{q_1+q_2}\left(\begin{pmatrix} \tilde{\mu}_1 \\ \tilde{\mu}_2 \end{pmatrix}, \begin{pmatrix} \tilde{\Sigma}_{11} & \mathbf{0} \\ \mathbf{0} & \tilde{\Sigma}_{22} \end{pmatrix}\right)$$

⑨ If $\tilde{X} \sim N_p(\tilde{\mu}, \tilde{\Sigma})$ with $|\tilde{\Sigma}| > 0$ and if $\tilde{\Sigma}^{-1}$ exists then

$$(\tilde{X} - \tilde{\mu})^T \tilde{\Sigma}^{-1} (\tilde{X} - \tilde{\mu}) \sim \chi^2_{(p)}$$

Proof

Recall: If Z_1, Z_2, \dots, Z_p are independent $\mathcal{N}(0, 1)$, then $\sum_{i=1}^p Z_i^2 \sim \chi^2_{(p)}$
So in general, if $\tilde{Z} \sim N_p(0, I)$ then $\tilde{Z}^T \tilde{Z} \sim \chi^2_{(p)}$

$\tilde{X} - \tilde{\mu} \sim N_p(0, \tilde{\Sigma})$ and $\tilde{Z} = \tilde{\Sigma}^{-1/2}(\tilde{X} - \tilde{\mu}) \sim N(0, I)$

$$\begin{aligned} \text{Since } \tilde{Z}^T \tilde{Z} &= (\tilde{\Sigma}^{-1/2}(\tilde{X} - \tilde{\mu}))^T (\tilde{\Sigma}^{-1/2}(\tilde{X} - \tilde{\mu})) \\ &\equiv (\tilde{X} - \tilde{\mu})^T \tilde{\Sigma}^{-1/2} \tilde{\Sigma}^{-1/2} (\tilde{X} - \tilde{\mu}) \\ &\equiv (\tilde{X} - \tilde{\mu})^T \tilde{\Sigma}^{-1} (\tilde{X} - \tilde{\mu}) \sim \chi^2_{(p)} \end{aligned}$$

Sampling Distribution of \bar{X} and S

① $X \sim N(\mu, \sigma^2)$, then $\bar{X} \sim N(\mu, \sigma^2/n)$

② $(n-1) \frac{s^2}{\sigma^2} \sim \chi^2_{(n-1)} \Rightarrow (n-1)s^2 \sim \sigma^2 \chi^2_{(n-1)}$ $s^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$

③ If $\tilde{X} \sim N_p(\tilde{\mu}, \tilde{\Sigma})$, then a) $\bar{X} \sim N_p(\tilde{\mu}, \frac{1}{n} \tilde{\Sigma})$

b) $(n-1) \underline{S} \sim \mathcal{W}_{(n-1)} \rightarrow$ Wishart distribution

c) \bar{X} and S are independent [HW: Prove for the univariate case]

2 important theorems:

① Law of large numbers:

Y_1, \dots, Y_n are independent observation from a population with $E(Y_i) = \mu$. Then $\bar{Y} \xrightarrow{P} \mu$, that is $P(|\bar{Y} - \mu| < \varepsilon) = 0 \forall \varepsilon > 0$.

As a consequence of LLN we can also say

$$S^2 \xrightarrow{P} \sigma^2 \quad \text{or} \quad \frac{1}{n} \sum (Y_i - \bar{Y})^2 \xrightarrow{P} \sigma^2$$

for the multivariate case

a) Each $\bar{X}_i \xrightarrow{P} \mu_i ; i=1, \dots, p$ so that

$$\bar{\underline{X}} \xrightarrow{P} \underline{\mu}$$

b) Each sample covariance $s_{ik} \xrightarrow{P} \sigma_{ik} ; k=1, 2, \dots, p$ so that

$$\underline{S} (\text{or } S_n) \xrightarrow{P} \Sigma$$

② Central Limit Theorem

Let X_1, \dots, X_n be independent observation from a population with mean μ and finite (non singular) covariance Σ , then $\sqrt{n}(\bar{X} - \mu) \sim N_p(0, \Sigma)$ when n is large.

It follows from here that

$$n(\bar{X} - \bar{Y})^T \Sigma^{-1} (\bar{X} - \bar{Y}) \sim \chi_{(p)}^2$$

Properties of Wishart Distribution

① If $\bar{Z}_1, \bar{Z}_2, \dots, \bar{Z}_m$ are iid $N_p(\bar{\Omega}, \Sigma)$

then $\sum_{j=1}^m \bar{Z}_j \bar{Z}_j^T \sim W_m(\Sigma)$

② If $A_1 \sim W_{m_1}(\Sigma)$ independent of $A_2 \sim W_{m_2}(\Sigma)$

then

$$A_1 + A_2 \sim W_{m_1 + m_2}(\Sigma)$$

③ If $A \sim W_p(\Sigma)$, then $CAC^T \sim W_p(C\Sigma C^T)$

$\underset{K \times p}{K \times p}$ $\underset{p \times p}{p \times p}$ $\underset{p \times K}{p \times K}$

The Wishart distribution is a family of distributions for symmetric positive definite matrices.

Let $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n$ be independent $N_p(\bar{\Omega}, \Sigma)$, we can form a $n \times p$ data matrix $\bar{X} = (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_n)$. Then the distribution of the $p \times p$ $M = \bar{X}^T \bar{X} = \sum_{i=1}^n \bar{X}_{ii} \bar{X}_{ii}^T$ is said to follow the Wishart distribution.

Density function

For a random matrix M which follows a Wishart distribution with n degrees of freedom and covariance matrix Σ , $M \underset{p \times p}{\sim} W_n(\Sigma)$, if $n \geq p$

$$W_p(n, \Sigma) \leftarrow$$

$$f(\mathbf{M}) = \frac{1}{2^{np/2} \Gamma_p(\gamma_2) |\Sigma|^{n/2}} |\mathbf{M}|^{(n-p+\gamma_2)} \exp\left[-\frac{1}{2} \text{trace}(\Sigma^{-1} \mathbf{M})\right]$$

↑
multivariate gamma

The Wishart distribution is a multivariate extension of the χ^2 distribution so that

$$\mathbf{M} \sim W_p(n, \sigma^2) \underset{\text{def}}{\approx} \sigma^2 \mathcal{X}_p(n)$$

Assessing the assumption of normality.

- Several of the multivariate statistical techniques to be discussed assume that each vector of observations $\mathbf{x}_j \sim N_p(\mu, \Sigma)$
- This assumption is not so important in cases of large n when we are making inferences related to \bar{x} .
- We would like to check the normality assumption for p -dimensions, but in practice, it suffices to investigate univariate and bivariate normality.
- Given a particular data set, we may be interested in addressing the following questions:

- ① Do the marginal distributions of the elements of $\bar{\Sigma}$ appear to be normal? What about a few linear combinations of the components x_i ?
- ② Do the scatter plots of pairs of observations on different characteristics give elliptical appearances as expected from normal populations?

③ Are there any extreme observations that you should check for their accuracy?

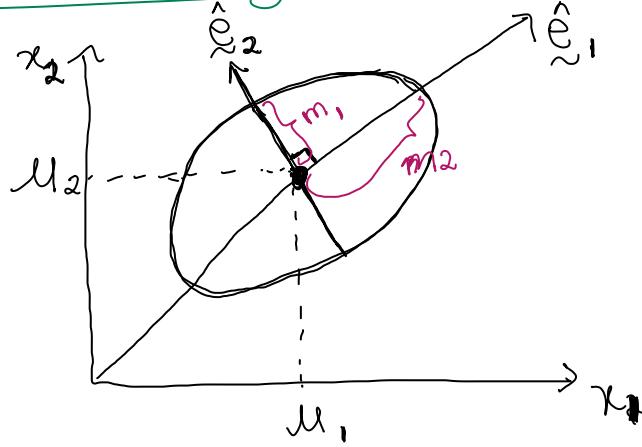
Evaluating normality in the univariate marginal distribution

* Check R file on Multivariate Normal.

Evaluating Bivariate Normality

Recall: If the observations were generated from a MVN distribution, each bivariate distribution would be normal and the contours of constant density would be ellipses.

> Understanding ellipses)



- \hat{e}_1, \hat{e}_2 ; \hat{e}_1 and \hat{e}_2 are the normalized eigen vectors

$$m_1 = \pm c \sqrt{\lambda_1} \hat{e}_1$$

$$m_2 = \pm c \sqrt{\lambda_2} \hat{e}_2$$

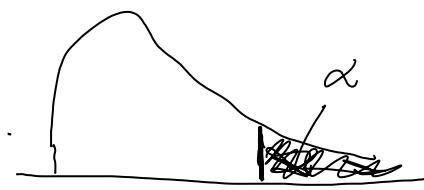
λ_1, λ_2 are eigen values

We can show that

$$c^2 = (\underline{x} - \underline{\mu})^\top \Sigma^{-1} (\underline{x} - \underline{\mu})$$

$$c^2 = (\underline{x} - \bar{\underline{x}})^\top S^{-1} (\underline{x} - \bar{\underline{x}})$$

Ellipsoid centered at $\underline{\mu}$ $= \chi^2(p, \alpha)$



a) Center = mean vector. Change in mean will shift the contours without altering the shape.

b) The orientation of the ellipses are determined by Σ . If the ellipse is vertically or horizontally oriented, it suggest

little to no covariance (or correlation) between X_1, X_2 .

If inclined at an angle, there is covariance / correlation between X_1, X_2

c) The size of the ellipse in terms of its major and minor axes provide information about the variance of the individual variables. A larger ellipse \Rightarrow higher variance in both variables
Smaller ellipse \Rightarrow lower variance - - -

d) Aspect ratio (ratio of major axis to the minor axis length) gives an idea of how the variances compare between X_1, X_2 .

$$\text{Aspect ratio} = \begin{cases} \text{Close to 1} & V(X_1) \text{ almost equal to } V(X_2) \\ \text{Large dev. from 1} & V(X_1) > V(X_2) \end{cases}$$

e) Change in Σ changes the shape.

f) if $\text{Cov}(X_1, X_2) = 0$, ellipse becomes circles indicates X_1 and X_2 are independent

Mahalanobis

① Given bivariate outcomes \underline{x} , such that we expect at least 0.5 or 50% of the observations to lie within the ellipse given by

$$\{\underline{x} : (\underline{x} - \bar{\underline{x}})^T S^{-1} (\underline{x} - \bar{\underline{x}}) \leq \chi^2_{\alpha}(0.5)\}$$

See example document and R code for illustration.

② Chi square plot (Bivariate & Multivariate)

a) Order squared distances $d_j^2 = (\underline{x} - \bar{\underline{x}})^T S^{-1} (\underline{x} - \bar{\underline{x}})$; $j=1, \dots, n$ from the smallest ; that $d_{(1)}^2 \leq d_{(2)}^2 \leq \dots \leq d_{(n)}^2$

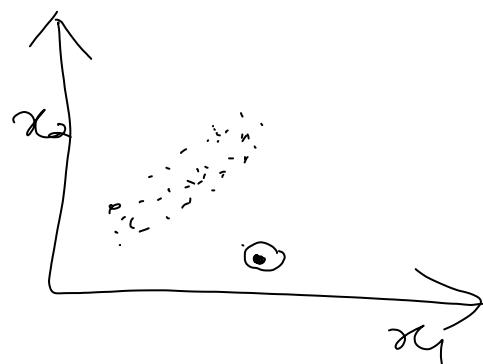
b) Graph the pairs $(g_{c,p}((j-1/2)/h), d_{(j)}^2)$

We expect the plot to resemble a straight line through the origin
 Any systematic curve suggest a lack normality

Detecting outliers in your data..

- Outliers are best detected visually when possible.
- Univariate : dot plot, boxplot , $(Q_1 - 1.5 \text{ IQR}, Q_3 + 1.5 \text{ IQR})$, qqplot etc.
- Bivariate : Scatterplot , chisquare plot
- Multivariate :

 - Univariate plots as mentioned above
 - Scatterplot of pairs of variables
 - Obtain standardized values $z_{jk} = \frac{(x_{jk} - \bar{x}_k)}{\sqrt{s_{kk}}}$ for $j=1,2,\dots,n$
 and each column $k=1,2,\dots,p$
 If $|z_{jk}| > 3.5$ observation is an outlier
 - Obtain Generalized squared distance $d_j^2 = (x_j - \bar{x})^T S^{-1} (x_j - \bar{x})$ and observe the chi square plot. or compare each value to $\chi_p^2(0.005)$
 if $d_j^2 > \chi_p^2(0.005)$ observation is an outlier.



Transformations to near normality

- If normality assumption is violated we may need to transform our data and consider performing our analysis based on the transformed data (after it follows normality of course!)

Helpful transformations to near normality

Original Scale	Transformed Scale
① Count, y	\sqrt{y}
② Proportions, \hat{p}	$\text{logit}(\hat{p}) = \frac{1}{2} \log\left(\frac{\hat{p}}{1-\hat{p}}\right)$
③ Correlations, r	Fisher's z -transformation $z(r) = \frac{1}{2} \log\left(\frac{1+r}{1-r}\right)$

- In many instances, the choice of transformation to improve the approximation to normality is not obvious and so we could make use of a family of transformations known as **Power transformations**

For univariate data which deviate from the normality assumption, our goal is to find a value λ such that x^λ would be near normal. Some commonly used choices of λ are:

x	transformation	
-1	$x^{-1} = \frac{1}{x}$	
0	$\ln x$	
1	x	
$\frac{1}{4}$	$x^{\frac{1}{4}}$	

} shrink large values of x

$\frac{1}{2}$	$x^{\frac{1}{2}}$	increases large values of x .
2	x^2	
3	x^3	

We can use the Box Cox technique to analytically choose λ . In this technique, λ is chosen such that

$$x^{(\lambda)} = \begin{cases} \frac{x^{\lambda}-1}{\lambda} & \lambda \neq 0 \\ \ln x & \lambda = 0 \end{cases} ; x \geq 0$$

Given the observations x_1, x_2, \dots, x_n , the Box-Cox solution for the choice of an appropriate power λ is the λ value that maximizes the expression

$$\ell(\lambda) = -\frac{n}{2} \ln \left[\frac{1}{n} \sum_{j=1}^n (x_j^{(\lambda)} - \bar{x}^{(\lambda)})^2 \right] + (\lambda-1) \sum_{j=1}^n \ln x_j$$

$$\text{where } \bar{x}^{(\lambda)} = \frac{1}{n} \sum_{j=1}^n x_j^{(\lambda)} = \frac{1}{n} \sum_{j=1}^n \left(\frac{x_j^\lambda - 1}{\lambda} \right)$$

In the multivariate case, let $\boldsymbol{\lambda} = (\lambda_1, \lambda_2, \dots, \lambda_p)$ be the power transformations for the p variables. Each λ_k ; $k=1, \dots, p$ is obtained by maximizing

$$\ell_k(\lambda) = -\frac{n}{2} \ln \left[\frac{1}{n} \sum_{j=1}^n (x_{jk}^{(\lambda_k)} - \bar{x}_{jk}^{(\lambda_k)})^2 \right] + (\lambda_k - 1) \sum_{j=1}^n \ln x_{jk}$$

Check book examples 4.16 & 4.17 in the examples document and the R code.

Note: After transforming, one needs to check the transformed data to be sure that normality assumption is satisfied

Chapter 5: Inferences about the mean vector

Statistical inference involves reaching valid conclusions concerning a population based on information from a sample. This may involve hypothesis testing or confidence intervals.

Recall: In univariate data analysis; suppose x_1, \dots, x_n is a random sample from a univariate $N(\mu, \sigma^2)$. Also suppose we want to test a hypothesis concerning the population mean such that

$$H_0: \mu = \mu_0 \quad \text{and} \quad H_a: \mu \neq \mu_0 \quad (\text{2-sided})$$

Test statistic :

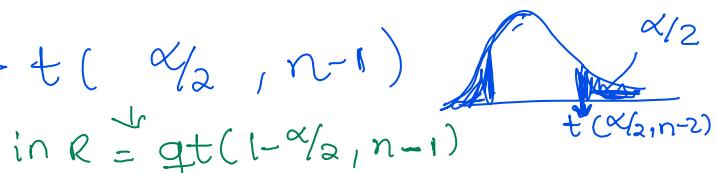
$$\text{Case 1: } t = \frac{(\bar{x} - \mu_0)}{\frac{s}{\sqrt{n}}} \sim t(n-1) \text{ under } H_0$$

This test statistic has an equivalent alternative given by

$$\text{Case 2: } t^2 = \frac{(\bar{x} - \mu_0)^2}{\frac{s^2}{n}} \sim F(1, n-1) \text{ under } H_0$$

Decision: Reject H_0 in favor of H_a if at a level of significance

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



$$\text{in R} = qt(1 - \alpha/2, n-1)$$

Case 2: Reject H_0 if $t^2 > F(1 - \alpha, 1, n-1)$

$$\text{in R} = qf(1 - \alpha, 1, n-1) \quad F(1 - \alpha, 1, n-1)$$

Conclusion: There is enough evidence to conclude " H_a "

There is not enough evidence to conclude " H_a "

A $100(1-\alpha)\%$ confidence interval for μ is given by

$$\bar{x} - t(\alpha/2, n-1) \frac{s}{\sqrt{n}} \leq \mu \leq \bar{x} + t(\alpha/2, n-1) \frac{s}{\sqrt{n}}$$

Note that the confidence interval consists of all the values of μ_0 that will not be rejected by a level α test of $H_0: \mu = \mu_0$, that is if $\mu_0 \in \text{Confidence Interval}$, fail to reject H_0 .

Check R for t-test example for $\mu_{\text{sweat rate}} = 0$

Consider now the problem of determining whether a pxl vector μ_0 is plausible for the mean of the multivariate normal distribution.

Why consider them jointly and not separately?

a) Overall Type I error becomes inflated when you test the variables individually. Type I error is the probability of falsely rejecting H_0 . ie $P(\text{Reject } H_0 | H_0 \text{ is true}) = \alpha$. For example, consider a dataset with 2 variables. Suppose we want an overall $\alpha = 0.05$. Performing each individual test at $\alpha = 0.05$ yields an overall $\alpha = 1 - (0.95)(0.95) = 0.1$. Of course, we can divide α by the number of tests (variables) and perform the individual tests using α/p to yield an overall α (Dunn-Bonferroni adjustment). This may solve the problem of inflating α

but then there is the 2nd reason:

- b) The individual univariate t-tests ignores correlation and covariance among dependent variables.
- c) It is possible that one may not find a single univariate mean difference between groups but a mean difference may exist when you consider the set of dependent variables jointly.

Note:

Situations may arise where all p variable means may not be jointly significant but some subsets could be, so check for subsets of dependent variables as well.

— We may extend the univariate t-test to the multivariate case:

Given a random sample $\underline{X}_1, \underline{X}_2, \dots, \underline{X}_n$ from an $N_p(\underline{\mu}, \Sigma)$ population, where $\underline{X}_j = (X_{j1}, X_{j2}, \dots, X_{jp})$ and both $\underline{\mu}$ and Σ are unknown. Let

$$\bar{\underline{X}} = \frac{1}{n} \sum_{j=1}^n \underline{X}_j \quad \text{and} \quad \bar{\Sigma} = \frac{1}{n-1} \sum_{j=1}^n (\underline{X}_j - \bar{\underline{X}})(\underline{X}_j - \bar{\underline{X}})^T$$

be the sample mean vector and sample covariance matrix respectively. Suppose we wish to test:

$$H_0: \underline{\mu} = \underline{\mu}_0 \quad \text{vs} \quad H_a: \underline{\mu} \neq \underline{\mu}_0 \quad (\text{two-sided})$$

where $\underline{\mu}_0$ is known.

Test statistic which is known as Hotelling's T^2 test statistic is given by

$$T^2 = n(\bar{\underline{X}} - \underline{\mu}_0)^T \bar{\Sigma}^{-1} (\bar{\underline{X}} - \underline{\mu}_0)$$

Under H_0 ,

$$\bar{T}^{*2} = \frac{n-p}{p(n-1)} T^2 \sim F(p, n-p)$$

Decision: Reject H_0 if $\bar{T}^{*2} > F(1-\alpha, p, n-p)$ 

Note: Hotelling's T^2 statistic reduces to the univariate test statistic t^2 when $p=1$

Refer to chapter 5 examples document corresponding to Book

examples 5.1 and 5.2. Check R code too

A confidence region for the mean vector $\underline{\mu}$ is also an extension of the CI for the μ in the univariate case. Recall the relationship between confidence interval / regions and hypothesis tests : A $100(1-\alpha)\%$ confidence interval / region can be derived when we invert the acceptance region of a level α test.

So by inverting the acceptance region of the Hotelling's T^2 test, we can obtain the confidence region for the mean vector $\underline{\mu}$.

The $100(1-\alpha)\%$ confidence region for the mean of p -dimensional normal distribution is the ellipsoid determined by all $\underline{\mu}$ such that

$$\nabla(\bar{x} - \underline{\mu})^T S^{-1} (\bar{x} - \underline{\mu}) \leq \frac{(n-1)p}{n-p} F(\alpha, p, n-p)$$

* If the squared distance $\nabla(\bar{x} - \underline{\mu}_0)^T S^{-1} (\bar{x} - \underline{\mu}_0) > \frac{(n-1)p}{n-p} F(\alpha, p, n-p)$, $\underline{\mu}_0$ is not in the confidence region and this leads to rejecting H_0 for a hypothesis test for $H_0: \underline{\mu} = \underline{\mu}_0$ vs $H_a: \underline{\mu} \neq \underline{\mu}_0$.

Hypothesis test on q Linear Combinations

More generally, suppose $\underline{x}_1, \dots, \underline{x}_n$ is a random sample from $N_p(\underline{\mu}, \Sigma)$. To test the hypothesis :

$$H_0: \underbrace{\underline{C} \underline{\mu}}_{q \times p} = \underline{\mu}_0 \quad \text{vs} \quad H_a: \underline{C} \underline{\mu} \neq \underline{\mu}_0$$

$$\text{Let } Y_i = \underline{C} \underline{x}_i - \underline{\mu}_0$$

$$\text{Then } Y_i \sim N_p(\underline{\mu}_Y, \Sigma_Y)$$

$$\text{where } \underline{\mu}_Y = \underline{C} \underline{\mu} - \underline{\mu}_0 \quad \text{and} \quad S_Y = \underline{C} S_X \underline{C}^T$$

Test statistic :

$$T^{*2} = n (\underline{C} \bar{\underline{x}} - \underline{\mu}_0)^T (\underline{C} S_X \underline{C}^T)^{-1} (\underline{C} \bar{\underline{x}} - \underline{\mu}_0)$$

$$= n \bar{Y}^T (\bar{S}_Y)^{-1} (\bar{Y})$$

Decision: Reject H_0 if $T^{*2} > \frac{(n-1)q}{n-q} F(1-\alpha, q, n-q)$

= Back to joint confidence region for μ

$p=2$: confidence ellipsoid

$p=3$: something spherical (like a football)

$p \geq 4$: hard to graph

Beginning from the center \bar{Y} , the axes of the confidence ellipsoids (within each confidence region),

there are ellipsoids are given by

$$\pm \sqrt{\lambda_i} \underbrace{\begin{bmatrix} p(n-1) \\ n(n-p) \end{bmatrix}}_{\text{length of the axes}} f(\alpha, p, n-p) \sim e_i ; \sum_{i=1,2,\dots,p} e_i = \lambda e_i$$

e_i - eigen vector corresponding to eigen value λ_i derived from the covariance matrix.

Example 5.3 from example document

Simultaneous Confidence Intervals

Suppose $\bar{Y} \sim N_p(\bar{\mu}, \Sigma)$ and $\bar{a} \in \mathbb{R}^p$ (a vector), we can form the linear combination

$$\bar{Z} = a_1 \bar{X}_1 + a_2 \bar{X}_2 + \dots + a_p \bar{X}_p = \bar{a}^T \bar{X}$$

We know $\bar{M}_Z = E(\bar{Z}) = \bar{a}^T \bar{\mu}$ and $\sigma_Z^2 = \text{Var}(\bar{Z}) = \bar{a}^T \Sigma \bar{a}$

$$\text{So } \bar{Z} \sim N(\bar{a}^T \bar{\mu}, \bar{a}^T \Sigma \bar{a})$$

For a random sample $\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n$ from $N_p(\tilde{\mu}, \Sigma)$, Consider a linear combination of the sample from the j th trial:

$$z_j = a_1 x_{j1} + a_2 x_{j2} + \dots + a_p x_{jp} = \tilde{a}^T \tilde{x}_j \quad j=1, 2, \dots, n$$

Then for n trials, we have $\bar{z}_1, \bar{z}_2, \dots, \bar{z}_n$ with

$$\bar{z} = \tilde{a}^T \tilde{x} \quad \text{and} \quad S_z^2 = \tilde{a}^T S_x \tilde{a} \xrightarrow{\Sigma \text{ unknown}}$$

Then for a fixed \tilde{a} and σ_z^2 unknown, the $100(1-\alpha)\%$ CI for $\mu_z = \tilde{a}^T \tilde{\mu}$ will be given by

$$\bar{z} \pm t(\alpha/2, n-1) \sqrt{\frac{\tilde{a}^T S_a \tilde{a}}{n}} \equiv \pm \frac{1}{n} S_z$$

fact:

Let $\tilde{x}_1, \dots, \tilde{x}_n$ be a random sample from $N_p(\tilde{\mu}, \Sigma)$ with Σ being positive definite. The $100(1-\alpha)\%$ Simultaneous confidence interval for all \tilde{a} is given by

$$\tilde{a}^T \tilde{x} \pm \sqrt{\frac{p(n-1)}{n(n-p)} F(\alpha, p, n-p) \tilde{a}^T S_a \tilde{a}} \quad \text{and this}$$

interval will contain $a^T \mu$ with probability $1-\alpha$.

Conveniently, if we choose $\tilde{a}^T = (0 \ 0 \ 0 \dots a_i \ 0 \ 0 \ 0)$ where $a_i = 1$, then the T^2 intervals allow us to conclude that

$$\bar{x}_i \pm \sqrt{\frac{p(n-1)}{n(n-p)} F(\alpha, p, n-p) s_{ii}} ; \quad \tilde{a}^T \tilde{x} = \bar{x}_i$$

- With the same confidence level $1-\alpha$, we can also make statements about mean differences, say $\mu_i - \mu_k$. In this case, choose $\tilde{a}^T = (0 \ 0 \dots a_i \ 0 \ 0 \dots 0 \ a_k \ 0 \ 0 \dots 0)$ where $a_i = 1$ and $a_k = -1$.

Point estimator for $\mu_i - \mu_k = \bar{x}_i - \bar{x}_k$

$$\begin{aligned}\text{Var}(\bar{x}_i - \bar{x}_k) &= V(\bar{x}_i) + V(\bar{x}_k) - 2 \text{Cov}(\bar{x}_i, \bar{x}_k) \\ &= \frac{1}{n} (V(x_i) + V(x_k) - 2 \text{Cov}(x_i, x_k)) \\ &= \frac{1}{n} (S_{ii} + S_{kk} - 2 S_{ik})\end{aligned}$$

Hence a $100(1-\alpha)\%$ CI for $\mu_i - \mu_k$ will be given as

$$\bar{x}_i - \bar{x}_k \pm \sqrt{\frac{p(n-1)}{n(n-p)} f(\alpha, p, n-p) (S_{ii} + S_{kk} - 2 S_{ik})}$$

Check examples 5.4 and 5.5 in examples document.

- Simultaneous confidence intervals are useful when one is interested in making specific inferences about each parameter separately, allowing for straightforward interpretation
 - Easier to visualize
- Joint confidence regions are useful when interested in the relationships among multiple parameters. They allow one to collectively assess the parameters and incorporate how changes

in one might be associated with changes in others.

- Complex or unfeasible to visualize, in higher dimensions.

More Precise Simultaneous Confidence Intervals

If p and/or the number of linear combinations $\underline{q}^T \underline{\mu}$ is small, we may be able to form shorter simultaneous confidence intervals. One procedure is the Bonferroni correction method in which to control the familywise error rate, we need to adjust the error rates of the individual test (or confidence intervals).

- We set $\alpha_k = \alpha/q$ to control the familywise error rate at α when conducting q tests.
- The $100(1-\alpha)$ % familywise confidence interval for q linear combinations, $\underline{q}^T \underline{\mu}$ is given as

$$\underline{q}^T \bar{\underline{x}} \pm t\left(\frac{\alpha_k}{2}, n-1\right) \sqrt{\frac{1}{n} \underline{q}^T \underline{S} \underline{q}} ; \alpha_k = \frac{\alpha}{q}$$

This is equivalent to

$$\underline{\mu}_1 \in \bar{x}_1 \pm t\left(\frac{\alpha}{2q}, n-1\right) \sqrt{\frac{s_{11}}{n}}$$

;

$$\underline{\mu}_q \in \bar{x}_q \pm t\left(\frac{\alpha}{2q}, n-1\right) \sqrt{\frac{s_{qq}}{n}}$$

Note the Bonferroni procedure uses $t(\alpha/2q, n-1)$ quantiles

instead of $\sqrt{\frac{(n-1)p}{n-p} F(\alpha, p, n-p)}$ as in the previously discussed simultaneous CI.

Check Example 5.6 in examples document.

fact:

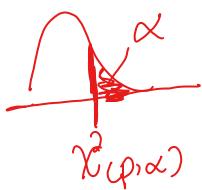
Let X_1, \dots, X_n be a random sample from a population with mean μ and positive definite covariance matrix Σ . When $n-p$ is large (regardless of whether data is normal or not),

$$\sqrt{n}(\bar{X} - \mu) \approx N(0, \Sigma) \text{ and}$$

$$n(\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu) \approx \chi^2(p)$$

Hence in testing $H_0: \mu = \mu_0$ against $H_1: \mu \neq \mu_0$ at a level of significance of α , we reject H_0 if

$$n(\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu) > \chi^2_{(p, \alpha)}$$



Note: Using this $\chi^2_{(p, \alpha)}$ or $\frac{(n-1)p}{n-p} F(\alpha, p, n-p)$ should yield same conclusions since these 2 are approximately equal when n is large relative to p .

The $100(1-\alpha)\%$ confidence region when $n-p$ is large is given by all μ such that:

$$n(\bar{X} - \mu)^T S^{-1} (\bar{X} - \mu) \leq \chi^2_{p(\alpha)}$$

Fact:

$\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n \sim N_p(\tilde{\mu}, \Sigma)$; Σ - positive definite. For

$n-p$ large,

$\tilde{a}^T \tilde{X} \pm \sqrt{\chi^2_{(p,\alpha)} \frac{\tilde{a}^T S \tilde{a}}{n}}$ will contain $\tilde{a}^T \tilde{\mu}$ for every

\tilde{a} .

Consequently, the $100(1-\alpha)\%$ simultaneous confidence intervals is given by

$$\tilde{\mu}_1 \in \bar{X}_1 \pm \sqrt{\chi^2_{(p,\alpha)} \frac{s_{11}}{n}}$$

$$\vdots$$

$$\tilde{\mu}_p \in \bar{X}_p \pm \sqrt{\chi^2_{(p,\alpha)} \frac{s_{pp}}{n}}$$

The Bonferroni simultaneous confidence intervals for q intervals is given by:

$$\tilde{a}^T \tilde{X} \pm Z\left(\frac{\alpha}{2q}\right) \sqrt{\frac{1}{n} \tilde{a}^T S \tilde{a}}$$

So for example for individual means, \bar{X}_i , we have

$$\bar{X}_i \pm Z\left(\frac{\alpha}{2p}\right) \sqrt{\frac{s_{ii}}{n}} ; i=1, 2, \dots, p$$

* Note: one at a time CI do not use the familywise adjustments so for Bonferroni, that is

$$\bar{X}_k \pm t\left(\frac{\alpha}{2}, n-1\right) \sqrt{\frac{s_{kk}}{n}} \quad \text{and} \quad \bar{X}_k \pm Z\left(\frac{\alpha}{2}\right) \sqrt{\frac{s_{kk}}{n}}, \quad \text{respectively}$$

large sample

Comparison - Univariate and Multivariate tests

Univariate

(1) σ^2 is known

(2) σ^2 unknown

a) n. large

b) n small

① Σ is known

② Σ is unknown

a) n-p large

b) n-p small

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 :

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

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

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

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:

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

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

$X_{1,j1} = \text{variable 1 under treatment 1}$ } Population 1
 $X_{1,j2} = \text{...}$ }
 \vdots

$X_{1,jp} = \text{...}$ } p under treatment 1

$$x_{2j1} = \text{Variable 1 under treatment 2} \quad \begin{matrix} & \\ & 2 \\ & \vdots \\ & p \end{matrix} \quad \begin{matrix} & \\ & 2 \\ & \vdots \\ & p \end{matrix} \quad \left. \begin{matrix} & \\ & 2 \\ & \vdots \\ & p \end{matrix} \right\} \text{Population 2}$$

The paired difference random variables becomes:

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

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

:

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

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_{cor} (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

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

