

The higher eigenvalue is more important variables.

Quantitative - Metrics, Ratio →
Qualitative - Non-metric ; Nominal
Ordinal

Course Description

- Real life examples of Multivariate Data
- The data matrix.
- Calculation of summary statistics, mean vectors, covariance & correlation matrices.
- Test for mean vector, one-sample and Hotelling T^2 tests based on union intersection approach.
- Simultaneous confidence intervals for detecting important components.
- Testing equality of two popn means.
- Basic assumptions and application of principal components, discriminant functions and canonical correlations.
- Distributions of linear functions of a random vector and of quadratic forms.
- Multivariate normal regression and correlation analysis.
- Elements of multivariate of variance.
- Use of computer packages.

$$z = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0.556 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

EXAM - Definitions

Modelling & Interpretation: Aims of multivariate adjustment

Introduction.

- The multivariate analysis methods are all statistical techniques that simultaneously analyse multiple measurements on individuals or objects under investigation, that is, any simultaneous analysis of more than two variables can be loosely considered multivariate analysis.

Multivariate Techniques and Real Life Multivariate Data Applications.

- The multivariate analysis techniques available for selection in various fields of application include the following techniques discussed below:

a) Principal Components.

This is a statistical approach that can be used to analyse interrelationships among a large no. of variables and to explain these variables in terms of their underlying factors.

- The objective is to find a way of condensing the information contained in a no. of original variables into a smaller set of variables (factors) with a minimal loss of info (Hair, Black, Babin, & Anderson, 2010).

Example 1.1

An investigator may use factor analysis to help understand relationships between customers' ratings of a fast-food restaurant by asking the customer to take the restaurant on the following 6 variables: food taste, food temperature, freshness, waiting time, cleanliness and friendliness of employees.

The investigator would like to combine these 6 variables into a smaller number, by analysing the customer responses and the investigator might find that the variables food taste, temperature, and freshness combine together to form a single factor of food quality, whereas the variables waiting time, cleanliness, and friendliness of employees combine to form another single variable, service quality.

b) Multiple Linear Regression Analysis

- The analysis method is used when the research problem involves a single metric dependent variable presumed to be related to two or more metric independent variables. The objective of multiple regression is to predict the changes in the dependent variable in response to

$$4x_1 = 2.414 \text{ either } x_2 = 1, x_1 = \frac{1}{2.414} = 0.414 \quad e_1 = \begin{bmatrix} 0.414 \\ 1.000 \end{bmatrix} \quad c_1 = \begin{bmatrix} 1.000 \\ \sqrt{0.414^2 + 1.000^2} \end{bmatrix} = 1.09240$$

$$x_2 = 0.414 \quad \text{or} \quad x_2 = 1, x_1 = 0.414$$

changes in the independent variables.

Example 1.2

The investigator may be interested in predicting the amount or size of the dependent variable for example, monthly expenditure on leisure (dependent variable) might be predicted from information regarding a household's income, its size, and the age of the head of the household (independent variables).

Similarly, the investigator might attempt to predict a company's sales from info on its expenditure for advertising, the no. of sales people, and the no. of stores carrying its products.

(c) Multiple Discriminant Analysis.

- The multiple discriminant " technique is an appropriate multivariate technique if the single dependent variable is dichotomous (for example male-female) or multichotomous (for example high-medium-low) and therefore non-metric. As with multiple regressions, the independent variables are assumed to be metric.
- Discriminant analysis is applicable in situations in which the total sample can be divided into groups based on a non-metric dependent variable characterizing several classes. The primary objectives of multiple discriminant analysis are to understand group differences and to predict the likelihood that an entity (individual or object) will belong to a particular class or group based on several metric independent variables.
- Logistic regression models, often referred to as logistic analysis, are a combination of multiple regression and multiple discriminant analysis.

Example 1.3

An investigator wishes to distinguish innovators from non-innovators according to their demographic & psychographic profiles or distinguish heavy product users from light users or distinguish males from females or distinguish national brand buyers from private-label buyers or distinguish good credit risks from poor credit risks.

Example 1.4

The revenue authority uses discriminant analysis to compare selected federal tax returns with a composite, hypothetical, normal tax payer's return (at diff. income levels) to identify the most promising returns and areas of audit.

d) Logistic Regression.

- Logistic P " models are distinguished from discriminant analysis primarily in that they accommodate all types of independent variables (metric & non-metric) & do not require the assumption of multivariate normality.
- However, in many instances particularly with more than two levels of dependent variables, discriminant analysis is the more appropriate analysis.

Example 1.5

- The logistic regression model is used to identify the financial & managerial data that best differentiated btwn \bar{e} successful & unsuccessful firms in order to select the best candidates for investment in the future by reviewing past records of financial management and placed firms data into one of two classes, successful over a five-yr period, and unsuccessful after 5 yrs.

e) Canonical Correlation.

- In " analysis the objective is to correlate simultaneously several metric dependent variables and several metric independent variables and is viewed as a logical extension of multiple regression analysis.
- Whereas multiple regressions involve a single dependent variable, canonical correlation involves multiple dependent variables. The underlying principle is to develop a linear combination of each set of variables (both independent and dependent) in a manner that maximizes the correlation between the two sets.

Example 1.6

- The canonical correlation could be used to compare the perception of the world-class companies over the 20 questions with the perception of the company such that the investigator could then conclude whether the perception of the company are correlated with those of world-class companies. The technique provides info on the overall correlation of perception as well as the correlation btwn each of \bar{e} 20 qstns.

f) Multivariate Analysis of Covariance and Variance.

- This is a statistical technique that can be used to simultaneously ^{explore} explain the relationship btwn several categorical independent variables (treatments) and two or more metric dependent variables.
- Multivariate analysis of covariance can be used in conjunction with multivariate analysis of variance

to remove (after experiment) the effect on any uncontrolled metric independent variables (covariates) on the dependent variables.

Example 1.7

An investigator may wish to know if a humorous ad will be more effective with its customers than a non-humorous ad. The investigator develops two ads - one humorous and one non-humorous ad and then shows a group of customers the two ads. After seeing the ads, the customers would be asked to rate the company and its products on several dimensions, such as modern vs traditional or high quality vs low quality.

The multivariate analysis of variance would be the technique to use to determine the extent of any statistical difference between the perception of customers who saw the humorous ad vs those who saw the non-humorous one.

g) Conjoint Analysis:

- This is a dependence technique that brings new sophistication to the evaluation of objects, such as new products, services or ideas. The most direct application is in new product or service dev., allowing for the evaluation of complex products while maintaining a realistic decision context for the respondent.
- The investigator is able to assess the importance of attributes as well as the levels of each attribute while consumers evaluate only a few product profiles, which are combination of product levels.

Example 1.8

Consider a developed product concept that has 3 attributes (price, quality, and colour) that have three, two & three levels respectively (for example colour has possible levels of for example red, yellow and blue). Instead of having to evaluate all 18 ($3 \times 2 \times 3$) possible combinations, a subset (9 or more) can be evaluated for 3 attractions to customers, and the investigator knows not only how important each attribute is but also the importance of each level.

h) Cluster Analysis

- This is an analytical technique for developing meaningful subgroups of individuals or objects. The objective is to classify a sample of entities (individuals or objects) into a small number of

- mutually exclusive groups based on similarities among entities.
- In cluster analysis, unlike discrimination analysis, the groups are not predefined. The technique is used to identify the groups.

Example 1.9

A restaurant owner wishes to know whether customers are patronizing the restaurant for diff. reasons. The data could be collected on perceptions of pricing and food quality. Cluster analysis could be used to determine whether some subgroups (clusters) are highly motivated by low prices vs those who are much less motivated to come to the restaurant based on the price considerations.

(i) Perception Mapping

- In " " (multidimensional scaling), the objective is to transform consumer judgements of similarity or preference (for example preference for stores or brands) into distances represented in multidimensional space. If objects A & B are judged by respondents as being the most similar compared with all other possible pairs of objects, perception mapping techniques will position objects A and B in such a way that the distance btwn them in multidimensional space is smaller than the distance btwn any other pairs of objects.

Example 1.10

The owner of a Burger King franchise wants to know whether the strongest competitor is MacDonald's or Wendy's. A sample of customers is given a survey and asked to rate the pairs of restaurants from most similar to least similar. The results show that the Burger King is most similar to Wendy's, so the owners know that the strongest competitor is the Wendy's restaurant because it is thought to be the most similar.

(j) Correspondence Analysis.

- This is the technique that facilitates the perceptual mapping of objects (for example products, persons) on a set of non-metric attributes. Investigators are constantly faced with the need to "qualify the qualitative data" found in normal variables. Correspondence analysis differs from the interdependence technique discussed earlier in its ability to accomodate both non-metric data and non-linear relationships.

Example 1-11

The respondents' brand preferences can be cross tabulated on demographic variables (for example gender, income, categories, occupation) by indicating how many people preferring each brand fall into each category of the demographic variables.

Through correspondence analysis, the association, or "correspondence", of brands and the distinguishing xtics of those preferring each brand are then shown in two- or three dimensional map of both brands and respondent xtics.

Brands perceived as similar are located close to one another. Likewise, the distinguishing xtic of respondents preferring each brand is also determined by the proximity of the demographic variable categories to the brand's position.

k) Structured Equation Modelling and Confirmatory Factor Analysis.

- This is a technique that allows separate relationships for each of a set of dependent variables. It provides the appropriate and most efficient estimation technique for a series of separate multiple regression equations estimated simultaneously.
- The multivariate analysis of variance and canonical correlation are not applicable in this situation because they allow only a single relationship btwn dependent & independent variables. In confirmatory factor analysis the researcher can assess the contribution of each scale item as well as incorporate how well the scale measures the concept (reliability).

Example 1-12

A study by management consultants identifies several factors that affect worker satisfaction: supervisor support, work environment, and job performance.

2. DESCRIPTIVE STATISTICS.

The data matrix.

- The measurements collected on several variables (characteristics) are frequently arranged and displayed in various ways that may be, for example, graphs or tabular arrangements

Arrays:

- The multivariate data arise whenever an investigator seeks to understand a social or physical phenomenon and selects a number $p \geq 1$ of variables or characteristics to record. The values of the variable recorded for each distinct item, individual or experiment unit for the p variables may be recorded in tabular form as shown below:

| 1 | 2 | ... | i | ... | p | |
|----------|----------|-----|----------|-----|----------|----------|
| X_{11} | X_{21} | ... | X_{i1} | ... | X_{p1} | X_{n1} |
| X_{12} | X_{22} | ... | X_{i2} | ... | X_{p2} | X_{n2} |
| : | : | ... | : | ... | : | : |
| X_{ij} | X_{2j} | ... | X_{ij} | ... | X_{pj} | X_{nj} |
| : | : | ... | : | ... | : | : |
| X_{1n} | X_{2n} | ... | X_{in} | ... | X_{pn} | X_{nn} |

Where X_{ij} is the value of the i th variable that is observed on the j th item or trial.

The data can alternatively be displayed as rectangular array of n rows and p columns denoted by X as shown below which is commonly referred to as data matrix.

$$X = \begin{bmatrix} X_{11} & X_{21} & \dots & X_{p1} \\ X_{12} & X_{22} & \dots & X_{p2} \\ \vdots & \vdots & \ddots & \vdots \\ X_{1n} & X_{2n} & \dots & X_{pn} \end{bmatrix} \quad X = [X_{11} \quad X_{12} \quad \dots \quad X_{1n}] \quad X = [X_{11} \quad X_{21} \quad \dots \quad X_{p1} \quad X_{12} \quad X_{22} \quad \dots \quad X_{p2} \quad \dots \quad X_{1n} \quad X_{2n} \quad \dots \quad X_{pn}]$$

Mean Vector, Variance - Covariance and Correlation Matrices.

- The info contained in the data can be assessed by calculating certain summary numbers commonly known as descriptive statistics that may be the measure of location or dispersion as described below:

(a) Arithmetic Mean.

- The measure of central location (the arithmetic mean) of the n measurements $X_{11}, X_{12}, \dots, X_{1n}$ of the first variable is given as:

$$\bar{X}_1 = \frac{1}{n} \sum_{j=1}^n X_{1j}$$

The general form of the arithmetic mean of the k^{th} variable is written as:

$$\bar{X}_k = \frac{1}{n} \sum_{j=1}^n X_{kj} \text{ where } k = 1, 2, \dots, p$$

If the n measurements represent a subset of the full set of measurements that might have been obtained, then, \bar{X}_k , is the sample mean of the k^{th} variable.

(b) Variance.

The measure of spread is provided by the sample variance defined for n measurements in the first variable given as:

$$S_1^2 = S_{11} = \frac{1}{n-1} \sum_{j=1}^n (X_{1j} - \bar{X}_1)^2$$

Where \bar{X}_1 is the sample mean of the X_{1j} 's, generally, for p variables, we have

$$S_k^2 = S_{kk} = \frac{1}{n-1} \sum_{j=1}^n (X_{kj} - \bar{X}_k)^2 \text{ where } k = 1, 2, \dots, p$$

1201499976

c) Standard deviation.

The square root of the sample variance is known as the sample standard deviation.

$$S_k = \sqrt{S_k^2} = \sqrt{S_{kk}}$$

d) Covariance.

The measure of the linear association btwn measurements of 2 variables X_1 & X_2 is given as

$$S_{12} = \frac{1}{n-1} \sum_{j=1}^n (X_{1j} - \bar{X}_1)(X_{2j} - \bar{X}_2)$$

The general association of the i^{th} variable and the k^{th} variable is given as:

$$S_{ik} = \frac{1}{n-1} \sum_{j=1}^n (X_{ij} - \bar{X}_i)(X_{kj} - \bar{X}_k) \quad \begin{cases} i = 1, 2, \dots, p \\ j = 1, 2, \dots, n \end{cases}$$

e) Correlation Coefficient (Pearson's Product moment correlation coefficient)

- This is the measure of the linear association btwn 2 variables that does not depend on the units of measurement. The sample correlation coefficient for the i^{th} & k^{th} variables is defined as:

$$r_{ik} = \frac{\sum_{j=1}^n (x_{ij} - \bar{x}_i)(x_{kj} - \bar{x}_k)}{\sqrt{\sum_{j=1}^n (x_{ij} - \bar{x}_i)^2} \sqrt{\sum_{j=1}^n (x_{kj} - \bar{x}_k)^2}} \quad \left\{ \begin{array}{l} i = 1, 2, \dots, p \\ k = 1, 2, \dots, p \end{array} \right. = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}}$$

- The sample correlation coefficient r has the following properties:

- The value of r must be between $-1 \leq r \leq 1$ inclusive, i.e. $-1 \leq r \leq 1$
- The r measures the strength of linear association.

> If $r = 0$, this indicates a lack of linear association b/w variables

> If $r < 0$, this indicates that there is a tendency for one value in the pair to be larger than its average when the other is smaller than its average.

> If $r > 0$, this indicates that there is a tendency for one value in the pair to be larger than the average when the other value is larger than the average and also for both values to be smaller than the average together.

- The value of r_{ik} remains unchanged if the measurements of the i^{th} variable are changed to $y_{ij} = ax_{ij} + b$, $j = 1, 2, \dots, n$ and the values of the k^{th} variable are changed to $y_{kj} = cx_{kj} + d$, $j = 1, 2, \dots, n$ provided that the constants a and c have the same sign.
- The quantities s_{ik} & r_{ik} do not, in general, convey all there is to know about the association b/w 2 variables. Non linear associations can exist that are not revealed by these descriptive statistics as they provide measures of linear association (association along a line).

Example 2-1

- The data of variables X_1, X_2 and X_3 were collected and recorded as shown in the table below. Determine the sample mean vector \bar{x} , variance-covariance matrix, S_n , and correlation coefficient, r_n for X_1, X_2 and X_3 variables.

| | | | | |
|-------|-----|-----|-----|-----|
| X_1 | 4.0 | 5.0 | 4.0 | 3.0 |
| X_2 | 4.2 | 5.2 | 4.5 | 5.8 |
| X_3 | 2.0 | 4.0 | 3.0 | 5.0 |

Solution-

a) Mean Vector: $\bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}$

$$\bar{x}_1 = \frac{1}{4} (4+5+4+3) = 4$$

$$\bar{x}_2 = \frac{1}{4} (4.2+5.2+4.5+5.8) = 5.0$$

$$\bar{x}_3 = \frac{1}{4} (2+4+3+5) = 3.5$$

$$\bar{x} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \bar{x}_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 5.0 \\ 3.5 \end{bmatrix}$$

b) Variance - Covariance Matrix:

$$S_{ik} = \frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x}_i)(x_{kj} - \bar{x}_k) = \frac{1}{n-1} \left\{ \sum_{j=1}^n x_{ij} x_{kj} - n \bar{x}_i \bar{x}_j \right\}$$

$$S_{11} = \frac{1}{3} [(4-4)^2 + \dots + (3-4)^2] = 0.67$$

$$S_{12} = \frac{1}{3} [(4-4)(4-2.5) + \dots + (3-4)(5-3.5)] = -0.20$$

$$S_{22} = \frac{1}{3} [(4-2.5)^2 + \dots + (5-3.5)^2] = 0.45$$

$$S_{13} = \frac{1}{3} [(4-4)(2-3.5) + \dots + (3-4)(5-3.5)] = -0.33$$

$$S_{23} = \frac{1}{3} [(4-2.5)(2-3.5) + \dots + (5-3.5)(5-3.5)] = 0.87$$

$$S_{33} = \frac{1}{3} [(2-3.5)^2 + \dots + (5-3.5)^2] = 1.67$$

$$S_n = \begin{bmatrix} S_{11} & S_{12} & S_{13} \\ S_{21} & S_{22} & S_{23} \\ S_{31} & S_{32} & S_{33} \end{bmatrix} = \begin{bmatrix} 0.667 & -0.200 & -0.333 \\ -0.200 & 0.453 & 0.867 \\ -0.333 & 0.867 & 1.667 \end{bmatrix}$$

c) Correlation matrix.

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}}$$

$$r_{11} = \frac{s_{11}}{\sqrt{s_{11}} \sqrt{s_{11}}} = \frac{0.67}{\sqrt{0.67} \sqrt{0.67}} = 1.000, r_{22} = \frac{s_{22}}{\sqrt{s_{22}} \sqrt{s_{22}}} = 1.000, r_{33} = \frac{1.67}{\sqrt{1.67} \sqrt{1.67}} = 1.000$$

$$r_{12} = \frac{s_{12}}{\sqrt{s_{11}} \sqrt{s_{22}}} = \frac{-0.20}{\sqrt{0.67} \sqrt{0.45}} = -0.364, r_{13} = \frac{s_{13}}{\sqrt{s_{11}} \sqrt{s_{33}}} = \frac{-0.33}{\sqrt{0.67} \sqrt{1.67}} = -0.316$$

$$r_{23} = \frac{s_{23}}{\sqrt{s_{22}} \sqrt{s_{33}}} = \frac{0.87}{\sqrt{0.45} \sqrt{1.67}} = 0.997$$

$$R_n = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} = \begin{bmatrix} 1.000 & -0.364 & -0.316 \\ -0.364 & 1.000 & 0.997 \\ -0.316 & 0.997 & 1.000 \end{bmatrix}$$

3x3

Alternative - 1.

to determine

Q: calculate

| x_1 | x_2 | x_3 | x_1^2 | x_2^2 | x_3^2 | $x_1 x_2$ | $x_1 x_3$ | $x_2 x_3$ | $(x_1 - \bar{x}_1)$ | $(x_2 - \bar{x}_2)$ |
|-------|-------|-------|---------|---------|---------|-----------|-----------|-----------|---------------------|---------------------|
| 4 | 4.2 | 2 | 16 | 17.64 | 4 | 16.8 | 8 | 8.4 | 0 | - |
| 5 | 5.2 | 4 | 25 | 27.04 | 16 | 26.0 | 20 | 20.8 | 1. | - |
| 4 | 4.8 | 3 | 16 | 23.04 | 9 | 19.2 | 12 | 14.4 | 0 | - |
| 3 | 5.8 | 5 | 9 | 33.64 | 25 | 17.4 | 15 | 29.0 | -1 | - |
| 16 | 20 | 14 | 66 | 101.36 | 54 | 79.4 | 55 | 72.6 | 0.3 | 6. |

$$\Sigma (x_1 - \bar{x}_1)(x_2 - \bar{x}_2)$$

$$\frac{\sqrt{(x_1 - \bar{x}_1)^2 + \Sigma (x_2 - \bar{x}_2)^2}}{\sqrt{1}}$$

a) Mean vector.

$$\bar{x}_i = \frac{1}{n} \sum_{j=1}^n x_{ij}$$

$$\bar{x}_1 = \frac{1}{4} \cdot 16 = 4, \quad \bar{x}_2 = \frac{20}{4} = 5, \quad \bar{x}_3 = \frac{14}{4} = 3.5$$

$$\bar{\mathbf{x}} = \begin{bmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \bar{x}_3 \end{bmatrix} = \begin{bmatrix} 4 \\ 5.0 \\ 3.5 \end{bmatrix}$$

b) Variance - covariance matrix.

$$S_n = \frac{1}{n-1} \left\{ \sum_{j=1}^n x_{ij} x_{kj} - n \bar{x}_i \bar{x}_k \right\}$$

$$S_{11} = \frac{1}{3} [66.0 - 4(\frac{16}{4} \times \frac{16}{4})] = \frac{1}{3} [66.0 - 64] = 0.667$$

$$S_{12} = \frac{1}{3} [79.4 - 4(\frac{16}{4} \times \frac{20}{4})] = \frac{1}{3} [79.4 - 80] = -0.200$$

$$S_{22} = \frac{1}{3} [101.36 - 4(\frac{20}{4} \times \frac{20}{4})] = \frac{1}{3} [101.36 - 100] = 0.453$$

$$S_{13} = \frac{1}{3} [55.0 - 4(\frac{16}{4} \times \frac{14}{4})] = \frac{1}{3} [55.0 - 56] = -0.333$$

$$S_{23} = \frac{1}{3} [72.6 - 4(\frac{20}{4} \times \frac{14}{4})] = \frac{1}{3} [72.6 - 70] = 0.867$$

$$S_{33} = \frac{1}{3} [54.0 - 4(\frac{14}{4} \times \frac{14}{4})] = \frac{1}{3} [54.0 - 49] = 1.667$$

$$S_n = \begin{bmatrix} S_{11} & S_{12} & S_{13} \\ S_{21} & S_{22} & S_{23} \\ S_{31} & S_{32} & S_{33} \end{bmatrix} = \begin{bmatrix} 0.667 & -0.200 & -0.333 \\ -0.200 & 0.453 & 0.867 \\ -0.333 & 0.867 & 1.667 \end{bmatrix}$$

c) Correlation Matrix.

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}}$$

As done in the previous solution

$$R_n = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} = \begin{bmatrix} 1.000 & -0.364 & -0.316 \\ -0.364 & 1.000 & 0.997 \\ -0.316 & 0.997 & 1.000 \end{bmatrix}$$

Example 2-2

The data for xtics x_1, x_2 & x_3 were collected and recorded in the table as shown below

Determine sample mean vector $\bar{\mathbf{x}}$, variance covariance matrix, S_n , and correlation coefficient matrix R_n for x_1, x_2 and x_3 variables.

Solution:

| X_1 | X_2 | X_3 | X_1^2 | X_2^2 | X_3^2 | $X_1 X_2$ | $X_1 X_3$ | $X_2 X_3$ |
|-------|-------|-------|---------|---------|---------|-----------|-----------|-----------|
| 0.80 | 21.4 | 7 | 0.64 | 457.96 | 49 | 17.120 | 5.6 | 149.8 |
| 0.82 | 26.8 | 7 | 0.6724 | 718.24 | 49 | 21.976 | 5.74 | 187.6 |
| 0.84 | 29.2 | 7.8 | 0.7056 | 852.64 | 60.84 | 24.528 | 6.552 | 227.76 |
| 0.82 | 31.8 | 7.5 | 0.6724 | 1011.24 | 56.25 | 26.076 | 6.15 | 238.5 |
| 0.80 | 36.0 | 7 | 0.64 | 1296 | 49 | 28.600 | 5.6 | 252 |
| 0.70 | 31.5 | 3.8 | 0.49 | 992.25 | 14.44 | 22.050 | 2.66 | 119.7 |
| 0.82 | 26.8 | 7 | 0.6724 | 718.24 | 49 | 21.976 | 5.74 | 187.6 |
| 0.80 | 15.2 | 7.2 | 0.64 | 231.04 | 51.84 | 12.160 | 5.76 | 109.44 |
| 0.82 | 28.0 | 8 | 0.6724 | 784 | 64 | 22.960 | 6.56 | 224 |
| 0.80 | 26.0 | 7.6 | 0.64 | 676 | 57.76 | 20.80 | 6.08 | 197.6 |
| 0.82 | 18.4 | 7 | 0.6724 | 338.56 | 49 | 15.088 | 5.74 | 128.8 |
| 0.80 | 14.2 | 7.2 | 0.64 | 201.64 | 51.84 | 11.360 | 5.76 | 102.24 |
| 0.80 | 20.4 | 6.8 | 0.64 | 416.16 | 46.24 | 16.320 | 5.44 | 138.72 |
| 0.80 | 15.8 | 6.8 | 0.64 | 249.64 | 46.24 | 12.640 | 5.44 | 107.44 |
| 0.82 | 18.5 | 7.2 | 0.6724 | 342.25 | 51.84 | 15.170 | 5.904 | 133.2 |
| 12.06 | 360.0 | 104.9 | 9.710 | 9285.86 | 746.29 | 289.02 | 84.726 | 2504.40 |

a) Mean Vector

$$\bar{X}_1 = \frac{12.06}{15} = 0.80, \quad \bar{X}_2 = \frac{360}{15} = 24.00, \quad \bar{X}_3 = \frac{104.9}{15} = 6.99$$

$$\bar{X} = \begin{bmatrix} 0.80 \\ 24.00 \\ 6.99 \end{bmatrix}$$

b) Variance - Covariance Matrix: $X_{ik} = \frac{1}{n-1} \left\{ \sum_{j=1}^n X_{ij} X_{kj} - n \bar{X}_i \bar{X}_k \right\}$

$$S_{11} = \frac{1}{14} \left[9.71 - 15 \left(\frac{12.06}{15} \cdot \frac{12.06}{15} \right) \right] = \frac{1}{14} [9.71 - 9.70] = 0.001$$

$$S_{12} = \frac{1}{14} [289.02 - 15 \times 0.8 \times 24] = \frac{1}{14} [289.02 - 289.44] = -0.030$$

$$S_{22} = \frac{1}{14} [9285.86 - 15 \times 24 \times 24] = \frac{1}{14} [9285.86 - 8640] = 46.133$$

$$S_{13} = \frac{1}{14} [84.726 - 15 \times 0.8 \times 6.99] = \frac{1}{14} [84.726 - 84.34] = 0.028$$

$$S_{23} = \frac{1}{14} [2504.4 - 15 \times 24 \times 6.99] = \frac{1}{14} [2504.4 - 2517.60] = -0.943$$

$$S_{33} = \frac{1}{14} [746.29 - 15 \times 6.99 \times 6.99] = \frac{1}{14} [746.29 - 733.60] = 0.906$$

$$S_n = \begin{bmatrix} S_{11} & S_{12} & S_{13} \\ S_{21} & S_{22} & S_{23} \\ S_{31} & S_{32} & S_{33} \end{bmatrix} = \begin{bmatrix} 0.001 & -0.030 & 0.028 \\ -0.030 & 46.133 & -0.943 \\ 0.028 & -0.943 & 0.906 \end{bmatrix}$$

c) Correlation Matrix.

$$r_{11} = \frac{s_{11}}{\sqrt{s_{11}} \sqrt{s_{11}}} = \frac{0.001}{\sqrt{0.001} \sqrt{0.001}} = 1.00, r_{22} = 1.00, r_{33} = 1.00$$

$$r_{12} = \frac{s_{12}}{\sqrt{s_{11}} \sqrt{s_{22}}} = \frac{-0.030}{\sqrt{0.001} \sqrt{46.133}} = -0.140$$

$$r_{13} = \frac{s_{13}}{\sqrt{s_{11}} \sqrt{s_{33}}} = \frac{0.028}{\sqrt{0.001} \sqrt{0.906}} = 0.925$$

$$r_{23} = \frac{s_{23}}{\sqrt{s_{22}} \sqrt{s_{33}}} = \frac{-0.9043}{\sqrt{46.133} \sqrt{0.906}} = -0.146$$

$$r = \begin{bmatrix} 1.000 & -0.140 & 0.925 \\ -0.140 & 1.000 & -0.146 \\ 0.925 & -0.146 & 1.000 \end{bmatrix}$$

3. TEST FOR MEAN VECTOR AND HOTELING T² TESTS.

Test for Mean Vector.

Consider testing the hypothesis in the univariate case given as:

Hypothesis (two-sided) $H_0: \mu = \mu_0$ against $H_1: \mu \neq \mu_0$

$H_a: \mu \neq \mu_0$ at α l.o.s $H_0: \mu = \mu_0$

If x_1, x_2, \dots, x_n are observations from a normal distribution, the appropriate test statistic is given as:

$$t = \bar{x} - \mu_0 \text{ where } \bar{x} = \frac{1}{n} \sum_{j=1}^n x_j \text{ and } s^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2 \quad (3.1)$$

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

The test statistic has a student's t distribution with $n-1$ degrees of freedom. When testing the hypothesis, we make a decision of rejecting H_0 , if the observed $|t|$ exceeds a specified point of a t-distribution with $n-1$ d.f. Rejecting H_0 when $|t|$ is large is equivalent to rejecting H_0 for a large square given as

$$\frac{(x - \mu_0)^2}{s^2/n} t^2 = \frac{(\bar{x} - \mu_0)^2}{s^2/n} = n(\bar{x} - \mu_0)^2 (s^2)^{-1} (\bar{x} - \mu_0) \quad (3.2)$$

The variance s^2 is the square of the distance from the sample mean \bar{x} to the test value μ_0 . The units of distance are expressed in terms of s/\sqrt{n} or estimated s.d. of \bar{x} , once $\bar{x} \notin S^2$ are observed. The test becomes, reject H_0 in favour of H_a at α l.o.s if:

$$n(\bar{x} - \mu_0)(s^2)^{-1} (\bar{x} - \mu_0) > t^2_{\alpha/2} (n-1) \quad (3.3)$$

Where $t_{\alpha/2}^2(n-1)$ denotes the upper $100(1-\alpha/2)^{\text{th}}$ percentile of the t dist with $n-1$ degrees d.f.

If H_0 is not rejected, i.e. $H_0: \mu = \mu_0$ at $\alpha 1\text{-o.s.}$, then,

$$\frac{|\bar{X} - \mu_0|}{s/\sqrt{n}} \leq t_{\alpha/2}(n-1) \quad \dots \quad (3.4)$$

(μ_0 lies in the $100(1-\alpha)$ percent confidence interval)

$$\bar{X} \pm t_{\alpha/2}(n-1) s/\sqrt{n} \text{ or } \bar{X} - t_{\alpha/2}(n-1) s/\sqrt{n} \leq \mu_0 \leq \bar{X} + t_{\alpha/2}(n-1) s/\sqrt{n}$$

In multivariate analysis we consider the problem of determining whether a given $p \times 1$ vector where, μ_0 is the mean of a multivariate normal dist, i.e., test the hypothesis.

Hypothesis: $H_0: \mu_1 = \mu_{10}, \mu_2 = \mu_{20}, \dots, \mu_p = \mu_{p0}$ against

$H_a: \mu_i \neq \mu_{i0}$ for at least one value of i at $\alpha 1\text{-o.s.}$

The hypothesis is rewritten as:

$$H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \\ \vdots \\ \mu_{p0} \end{bmatrix} \text{ against } H_a: \mu_i \neq \mu_{i0} \text{ for at least one value of } i \quad (i=1, 2, \dots, p) \text{ at } \alpha 1\text{-o.s.}$$

$$\dots \quad (3.6)$$

A natural generalization of the squared distance in eqtn (3.2) is its multivariate analog

$$T^2 = (\bar{X} - \mu_0)' (\frac{1}{n} \Sigma)^{-1} (\bar{X} - \mu_0) = n (\bar{X} - \mu_0)' \Sigma^{-1} (\bar{X} - \mu_0)$$

Where;

Exam
Derive.
How to
calculate
it

$$\bar{X} = \begin{bmatrix} \bar{X}_1 \\ \bar{X}_2 \\ \vdots \\ \bar{X}_p \end{bmatrix} = \bar{X}_i = \frac{1}{n} \sum_{i=1}^n X_{ij} (p \times p) = \frac{1}{n-1} \sum_{j=1}^n (X_{ij} - \bar{X}_i)' (X_{ij} - \bar{X}_i) \text{ and } (\mu_{10}, \mu_{20}, \dots, \mu_{p0})'$$

The T^2 is called Hotelling's T^2 in honour of Harold Hotelling, a pioneer in multivariate analysis who first obtained its sampling dist. If the observed statistical distance T^2 is too large, that is, if \bar{X} is "too far" from μ_0 , the hypothesis

$H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \\ \vdots \\ \mu_{p0} \end{bmatrix}$ is rejected. There are no special tables of T^2 percentage points

required for formal tests of hypothesis because T^2 is distributed as $\frac{(n-p)p}{n-p} F_{p, n-p}$ where $F(p, n-p)$ denote a r.v with $p \leq n-p$ d.f. from an F -dist.

Exq

$$\frac{(n-p)p}{n-p} F_{p, n-p}$$

Example 6.1

The data matrix for a random sample of size $n=3$ from a bivariate normal popn is given as $X = \begin{bmatrix} 7 & 9 \\ 9 & 6 \\ 8 & 3 \end{bmatrix}$. State the sampling dist of T^2 and determine whether the sample is drawn from a multivariate normal dist with mean, $\mu = \begin{bmatrix} 9 \\ 5 \end{bmatrix}$ at $\alpha = 0.01 - 1.0 \cdot 5$

Solution

Hypothesis: $H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \end{bmatrix}$ against $H_a: \mu_i \neq \mu_{i0}$ for at least 1 value of i ($i=1, 2, \dots, p$) at $\alpha = 0.01 - 1.0 \cdot 5$

Test statistic;

| X_1 | X_2 | X_1^2 | X_2^2 | $X_1 X_2$ |
|-------|-------|---------|---------|-----------|
| 7 | 9 | 49 | 81 | 63 |
| 9 | 6 | 81 | 36 | 54 |
| 8 | 3 | 64 | 9 | 24 |
| 24 | 18 | 194 | 126 | 141 |

$$\bar{X}_1 = \frac{24}{3} = 8, \bar{X}_2 = \frac{18}{3} = 6$$

$$\bar{X} = \begin{bmatrix} 8 \\ 6 \end{bmatrix}$$

$$S_{ik} = \frac{1}{n-1} \left\{ \sum_{j=1}^n X_{ij} X_{kj} - n \bar{X}_i \bar{X}_k \right\}$$

$$S_{11} = \frac{1}{2} [194 - 3(8)^2] = 1.0$$

$$S_{12} = \frac{1}{2} [141 - 3(8)(6)] = -1.5$$

$$S_{22} = \frac{1}{2} [126 - 3(6)^2] = 9.0$$

$$S_n = \begin{bmatrix} 1.0 & -1.5 \\ -1.5 & 9.0 \end{bmatrix}$$

$$T^2 = n(\bar{X} - \mu_0)' \Sigma^{-1} (\bar{X} - \mu_0) = 3[-1 \ 1] \frac{1}{6.75} \begin{bmatrix} 9.0 & 1.5 \\ 1.5 & 1.0 \end{bmatrix} [-1 \ 1]$$

$$= 3[-1 \ 1] \frac{1}{6.75} \begin{bmatrix} -7.5 \\ -0.5 \end{bmatrix} = \frac{28}{9}$$

= 3.111, that has approximately F distribution.

$$-9 + 1.5 = -7.5$$

$$-0.5 + 1 = 0.5$$

Tabulated

$$T^2 = (n-1)p F_{\alpha}(p, n-p) = 2(2)p F_{0.01}(2, 1)$$

$$= 3-2$$

$$= 4 \frac{4(4999.5)}{4(4999)} = \frac{1.9996}{19996} 19998$$

Decision:

Since $T^2 < 4 F_{0.01}(2, 1)$, fail to reject $H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \end{bmatrix}$ and conclude that the sample is from a multivariate normal dist. with mean vector.

Exq. 3.2

The data of the variables X_1, X_2 & X_3 were collected & recorded. Determine whether the sample is drawn from multivariate normal dist with mean vector $\mu_0 = [3 \ 6 \ 3]$ at $\alpha = 0.01$ l.o.s

Solution

Hypothesis: $H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \\ \mu_{30} \end{bmatrix}$ against $H_a: \mu_i \neq \mu_{i0}$ at $\alpha = 0.01$ l.o.s

Test statistic

| X_1 | X_2 | X_3 | X_1^2 | $X_1 X_2$ | $X_1 X_3$ | X_2^2 | $X_2 X_3$ | X_3^2 |
|-------|-------|-------|---------|-----------|-----------|---------|-----------|---------|
| 4 | 4.2 | 2 | 16 | 16.8 | 8 | 17.64 | 8.4 | 4 |
| 5 | 5.2 | 4 | 25 | 26 | 20 | 27.04 | 20.8 | 16 |
| 4 | 4.8 | 3 | 16 | 19.2 | 12 | 23.04 | 14.4 | 9 |
| 3 | 5.8 | 7 | 9 | 17.4 | 21 | 33.64 | 40.6 | 49 |
| 16 | 20 | 16 | 66 | 79.4 | 61 | 101.36 | 84.2 | 78 |

$$\bar{X}_1 = \frac{16}{4} = 4, \bar{X}_2 = \frac{20}{4} = 5, \bar{X}_3 = \frac{16}{4} = 4$$

$$\bar{X} = \begin{bmatrix} 4 \\ 5 \\ 4 \end{bmatrix}$$

$$S_{IK} =$$

$$S_{11} = \frac{1}{n-1} \left\{ \sum_{j=1}^n X_{1j}^2 - n \bar{X}_1^2 \right\} = \frac{1}{3} [66 - 4(4)^2] = 0.67$$

$$S_{12} = \frac{1}{3} [79.4 - 4(4 \times 5)] = -0.20$$

$$S_{22} = \frac{1}{3} [101.36 - 4(5)^2] = 0.45$$

$$S_{13} = \frac{1}{3} [61 - 4(4 \times 4)] = -1.00$$

$$S_{23} = \frac{1}{3} [84.2 - 4(5 \times 4)] = 1.40$$

$$S_{33} = \frac{1}{3} [78 - 4(4)^2] = 4.67$$

$$S_n = \begin{bmatrix} 0.67 & -0.20 & -1.00 \\ -0.20 & 0.45 & 1.40 \\ -1.00 & 1.40 & 4.67 \end{bmatrix}$$

Determine the inverse of S_n

Determinant,

$$D = 0.667 \{ 0.453(4.667) - 1.40(1.40) \} + 0.20 \{ -0.20(4.667) + 1.00(1.40) \} \\ - 1.00 \{ -0.20(1.40) + 1.00(0.453) \} \\ = 0.024$$

Using calc to get cofactors. Input a matrix.
Shift + Det (Mat A) x Trsp (Mat A)⁻¹

cofactor of element a_{ij}
(-1)^{i+j} times minor of a_{ij}

$$\text{Cofactors} = \begin{bmatrix} 0.156 & -0.467 & 0.173 \\ -0.467 & 2.111 & -0.733 \\ 0.173 & -0.733 & 0.262 \end{bmatrix}$$

$$S_n^{-1} = \frac{1}{0.024} \begin{bmatrix} 0.156 & -0.467 & 0.173 \\ -0.467 & 2.111 & -0.733 \\ 0.173 & -0.733 & 0.262 \end{bmatrix} = \begin{bmatrix} 6.6 & -19.7 & 7.3 \\ -19.7 & 89.1 & -30.9 \\ 7.3 & -30.9 & 11.1 \end{bmatrix}$$

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

$$= 4 \begin{bmatrix} 1 & -1 & 1 \end{bmatrix} \begin{bmatrix} 6.6 & -19.7 & 7.3 \\ -19.7 & 89.1 & -30.9 \\ 7.3 & -30.9 & 11.1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix} = 4 \begin{bmatrix} 33.6 & -139.7 & 49.3 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$

$$= 890.25$$

Tabulated,

$$T^2 = \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) = \frac{3(3)}{4-3} F_{0.01}(3, 1)$$

$$= 9(5403.4) = 48627^{48630.6}$$

Decision;

Since $T^2 < 9.0$ (~~reject~~) $F_{0.01}(3, 1)$, fail to reject H_0 : $\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \\ \mu_{30} \end{bmatrix}$ and conclude that the sample is from a multivariate normal dist with $\mu_0 = [3 \ 6 \ 3]$ at $\alpha = 0.01$ L.O.S

Example 3.3

The data of the variable $x_1, x_2 \in x_3$ were collected and recorded in a table as shown below. Determine whether the sample is drawn from a multivariate normal dist with mean vector $\mu_0 = [4 \ 50 \ 10]$ at $\alpha = 0.01$ L.O.S

Solution:

Hypothesis: $H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \\ \mu_{30} \end{bmatrix}$ against $H_a: \mu_i \neq \mu_{10}$ at $\alpha = 0.01$ L.O.S

Test statistics:

| X_1 | X_2 | X_3 | X_1^2 | $X_1 X_2$ | $X_1 X_3$ | X_2^2 | $X_2 X_3$ | X_3^2 |
|-------|-------|-------|---------|-----------|-----------|---------|-----------|---------|
| 4 | 36 | 12 | 16 | 144 | 48 | 1296 | 432 | 144 |
| 5 | 54 | 12 | 25 | 270 | 60 | 2916 | 648 | 144 |
| 4 | 26 | 10 | 16 | 104 | 40 | 676 | 260 | 100 |
| 5 | 40 | 10 | 25 | 200 | 50 | 1600 | 400 | 100 |
| 2 | 20 | 10 | 4 | 40 | 20 | 400 | 200 | 100 |
| 8 | 54 | 14 | 64 | 432 | 112 | 2916 | 756 | 196 |
| 5 | 48 | 16 | 25 | 240 | 80 | 2304 | 768 | 256 |
| 7 | 46 | 12 | 49 | 322 | 84 | 2116 | 552 | 144 |
| 4 | 36 | 12 | 16 | 144 | 48 | 1296 | 432 | 144 |
| 6 | 40 | 12 | 36 | 240 | 72 | 1600 | 480 | 144 |
| 50 | 400 | 120 | 276 | 2136 | 614 | 17120 | 4928 | 1472 |

$$\bar{X}_1 = \frac{50}{10} = 5.0, \bar{X}_2 = \frac{400}{10} = 40, \bar{X}_3 = \frac{120}{10} = 12.0$$

$$\bar{X} = \begin{bmatrix} 5 \\ 40 \\ 12 \end{bmatrix}$$

$$S_{11} = \frac{1}{n-1} \left[\sum_{j=1}^n X_{1j}^2 - n\bar{X}_1^2 \right] = \frac{1}{9} [276 - 10(5)^2] = 2.9$$

$$S_{12} = \frac{1}{n-1} \left[\sum_{j=1}^n X_{1j} X_{2j} - n\bar{X}_1 \bar{X}_2 \right] = \frac{1}{9} [2136 - 10(5)(40)] = 15.1$$

$$S_{13} = \frac{1}{9} [614 - 10(5)(12)] = 1.6$$

$$S_{22} = \frac{1}{9} [17120 - 10(40)^2] = 124.4$$

$$S_{23} = \frac{1}{9} [4928 - 10(40)(12)] = 14.2$$

$$S_{33} = \frac{1}{9} [1472 - 10(12)^2] = 3.6$$

Determine the inverse of S_n

$$\text{Determinant, } D = 2.89 \{ 124.44(3.56) - 14.22(14.22) \} - 15.1 \{ (15.1(3.56) - 1.56(14.22)) \} - 1.56 \{ 15.11(14.22) - 124.44(1.56) \} = 249.5$$

Cofactors, $C = \begin{bmatrix} 240.2 & -31.6 & 21.33 \\ -31.6 & 7.85 & -17.6 \\ 21.33 & -17.6 & 131.2 \end{bmatrix}$

$$S_n^{-1} = \frac{1}{249.5} \begin{bmatrix} 240.2 & -31.6 & 21.33 \\ -31.6 & 7.85 & -17.6 \\ 21.33 & -17.6 & 131.2 \end{bmatrix} = \begin{bmatrix} 0.963 & -0.127 & 0.086 \\ -0.127 & 0.031 & -0.070 \\ 0.086 & -0.070 & 0.526 \end{bmatrix}$$

$$\bar{X} - \mu_0 = \begin{bmatrix} 40 \\ 12 \end{bmatrix} - \begin{bmatrix} 50 \\ 10 \end{bmatrix} = \begin{bmatrix} -10 \\ 2 \end{bmatrix}$$

$$T^2 = n(\bar{X} - \mu_0)' \Sigma^{-1} (\bar{X} - \mu_0)$$

$$= 10 \begin{bmatrix} 1 & -10 & 2 \end{bmatrix} \begin{bmatrix} 0.963 & -0.127 & 0.086 \\ -0.127 & 0.031 & -0.070 \\ 0.086 & -0.070 & 0.526 \end{bmatrix} \begin{bmatrix} 1 \\ -10 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 24 & -5.8 & 18.4 \end{bmatrix} \begin{bmatrix} 1 \\ -10 \\ 2 \end{bmatrix}$$

$$= 119.06$$

Tabulated,

$$T^2 = \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) = \frac{9(3)}{10-3} F_{0.01}(3, 7)$$

$$= 3.86(8.45) = 32.62$$

Decision :

Since $T^2 > 3.86 F_{0.01}(3, 7)$, reject H_0 : $\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \\ \mu_{30} \end{bmatrix}$ and conclude that the sample is not from a multivariate normal dist with $\mu_0 = [4 50 10]$ at $\alpha = 0.01$ l.o.s

NOTE : The H_0 will be rejected if one or some combination of means differs too much from the hypothesized values $[4 50 10]$. At this stage of analysis if we reject H_0 we would not know which of the hypothesized values would be supported by the data.

Simultaneous Confidence Intervals

- In order to obtain the method of making inferences from sample, there is need to extend the concept of univariate interval to multivariate confidence region. The rejection region for the mean of μ of p -dimensional normal population before the sample is selected from the previous section is given as:

$$p \left[n(\bar{X} - \mu_0)' S^{-1} (\bar{X} - \mu_0) > \frac{(n-1)p}{n-p} F_{\alpha/2}(p, n-p) \right] = \alpha/2$$

- The confidence region for the mean μ of a p -dimensional normal popn is then given as:

$$p \left[n(\bar{X} - \mu_0)' S^{-1} (\bar{X} - \mu_0) \leq \frac{(n-1)p}{n-p} F_{\alpha/2}(p, n-p) \right] = 1 - \alpha/2$$

- Whatever the values of the unknown μ and Σ in other words \bar{X} will be within

$$\left[\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \right]$$

- We write the intervals as:

$$\bar{X}_1 = \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)} \sqrt{\frac{s_{11}}{n}} \leq \mu_1 \leq \bar{X}_1 + \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)} \sqrt{\frac{s_{11}}{n}}$$

$$\bar{X}_2 = \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)} \sqrt{\frac{s_{22}}{n}} \leq \mu_2 \leq \bar{X}_2 + \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)} \sqrt{\frac{s_{22}}{n}}$$

$$\bar{X}_p = \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)} \sqrt{\frac{s_{pp}}{n}} \leq \mu_p \leq \bar{X}_p + \sqrt{\frac{(n-1)p}{n-p} F_{\alpha}(p, n-p)} \sqrt{\frac{s_{pp}}{n}}$$

All hold simultaneously with confidence coefficient $1-\alpha$.

Ex 3.4

The data of the variables X_1 & X_2 for $n=42$ were collected and statistics recorded;

$$\bar{X} = \begin{bmatrix} 0.564 \\ 0.603 \end{bmatrix}, S = \begin{bmatrix} 0.0144 & 0.0117 \\ 0.0117 & 0.0146 \end{bmatrix}, S^{-1} = \begin{bmatrix} 203.018 & -163.391 \\ -163.391 & 200.225 \end{bmatrix}$$

Determine the following.

a) Whether the sample is drawn from multivariate normal dist with mean vector $\mu = \begin{bmatrix} 0.562 \\ 0.589 \end{bmatrix}$ at $\alpha = 0.05$ l.o.s

b) The 95% confidence interval for the mean.

Solution:

(a) Hypothesis: $H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \end{bmatrix}$ against $H_a: \mu_i \neq \mu_{i0}$ for at least 1 value of i at $\alpha = 0.05$ l.o.s

Test statistic;

$$T^2 = n(\bar{X} - \mu_0)' \Sigma^{-1} (\bar{X} - \mu_0)$$
$$= 42 [0.564 - 0.562, 0.603 - 0.589] \begin{bmatrix} 203.018 & -163.391 \\ -163.391 & 200.225 \end{bmatrix} \begin{bmatrix} 0.002 \\ 0.014 \end{bmatrix}$$
$$= 1.21$$

Tabulated

$$T^2 = \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) = \frac{41(2)}{40} F_{0.05}(2, 40)$$

$$= 8.2(3.23) = 6.62$$

Conclusion: Since test statistic $T^2 <$ tabulated value, then $\bar{X}' = [0.564 \ 0.603]$ is in the acceptance region, therefore we fail to reject H_0 and conclude that the sample data is from the multivariate normal dist with mean, $\mu = [0.562 \ 0.589]$ at $\alpha = 0.05$ l.o.s

95% C.I of the mean

$$b) \bar{X}_1 - \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \sqrt{\frac{S_{11}}{n}} \leq \mu_1 \leq \bar{X}_1 + \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \sqrt{\frac{S_{11}}{n}}$$

$$= 0.564 - \frac{41(2)(3.23)}{40} \sqrt{\frac{0.0144}{4^2}} \leq \mu_1 \leq 0.564 + \frac{41(2)(3.23)}{40} \sqrt{\frac{0.0144}{4^2}}$$

$$= (0.516, 0.612)$$

$$\bar{X}_2 - \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \sqrt{\frac{S_{22}}{n}} \leq \mu_2 \leq \bar{X}_2 + \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \sqrt{\frac{S_{22}}{n}}$$

$$= 0.603 - \frac{41(2)(3.23)}{40} \sqrt{\frac{0.0146}{4^2}} \leq \mu_2 \leq 0.555 + 0.603 - \frac{41(2)(3.23)}{40} \sqrt{\frac{0.0146}{4^2}}$$

$$= (0.555, 0.651)$$

Exa 3.5

The data of the variables $X_1, X_2, \& X_3$ for $n=40$ were collected and recorded as below

$$\bar{X} = [5.4, 5.5, 2.5], S = [5.6, 6.0, 2.2; 6.0, 7.0, 2.4; 2.2, 2.4, 2.4], S^{-1} = [2.251, -1.860, -0.204; -1.860, 1.754, -0.049; -0.204, -0.049, 0.653]$$

Determine the following:

- a) Whether the sample is drawn from multivariate normal dist. with mean vector $\mu = \begin{bmatrix} 5 \\ 5 \\ 2 \end{bmatrix}$ at $\alpha=0.05$ b) The 95% C.I for the mean.

Solution

(a) Hypothesis: $H_0: \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix} = \begin{bmatrix} \mu_{10} \\ \mu_{20} \\ \mu_{30} \end{bmatrix}$ against $H_a: \mu_i \neq \mu_{i0}$ for at least 1 value of i ($i=1, 2, 3$) at $\alpha=0.05$

Test-statistic

$$T^2 = \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) = \frac{39(3)}{37} F_{0.05}(3, 37)$$

$$= \frac{117(2.86)}{37} = 9.04$$

37

Conclusion:

The mean vector $\bar{X}' = [5.4, 5.5, 2.5]$ is in the acceptance region, therefore the sample data is from multivariate normal dist. with mean vector, $\mu' = [5, 5, 2]$ at $\alpha=0.05$

(b) 95% C.I of the mean

$$\bar{X}_1 - \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \sqrt{\frac{S_{11}}{n}} \leq \mu_1 \leq \bar{X}_1 + \frac{(n-1)p}{n-p} F_{\alpha}(p, n-p) \sqrt{\frac{S_{11}}{n}}$$

$$= 5.4 - \sqrt{\frac{39(3)(2.86)}{37}} \sqrt{\frac{5.6}{40}} \leq \mu_1 \leq 5.4 + \sqrt{\frac{39(3)(2.86)}{37}} \sqrt{\frac{5.6}{40}}$$

$$= (4.27, 6.63)$$

$$\bar{X}_2 = 5.5 - \sqrt{\frac{39(3)(2.86)}{37}} \sqrt{\frac{7.0}{40}} \leq \mu_2 \leq 5.5 + \sqrt{\frac{39(3)(2.86)}{37}} \sqrt{\frac{7.0}{40}}$$

$$= (4.24, 6.76)$$

$$\bar{X}_3 = 2.5 - \sqrt{\frac{39(3)(2.86)}{37}} \sqrt{\frac{2.4}{40}} \leq \mu_3 \leq 2.5 + \sqrt{\frac{39(3)(2.86)}{37}} \sqrt{\frac{2.4}{40}}$$

$$= (1.76, 3.24)$$

Testing Equality of the Population Means

A T^2 -statistic for testing the equality of vector means from 2 multivariate populations can be developed by analogy with the univariate procedures. The T^2 -statistic is appropriate for comparing responses from one set of experimental settings (popn 1) with independent responses from another set of experimental settings (popn 2). If possible, the experimental units should be randomly assigned to the sets of experimental conditions. Randomization will, to some extent, mitigate the effect of unit-to-unit variability in a subsequent comparison of treatments. Although, some precision is lost relative to paired comparisons, the inferences in the two-popn cases are, ordinarily, applicable to a more general collection of experimental units simply because unit homogeneity is not required.

Table : The arrangement of observations on p ($p = 1, 2, \dots$) variables.

| Population | Sample | Mean | Variance |
|------------|-----------------------------------|---|--|
| 1 | $X_{11}, X_{12}, \dots, X_{1n_1}$ | $\bar{X}_1 = \frac{1}{n_1} \sum_{j=1}^{n_1} X_{1j}$ | $S_1 = \frac{1}{n_1-1} (X_{1j} - \bar{X}_1)' (X_{1j} - \bar{X}_1)$ |
| 2 | $X_{21}, X_{22}, \dots, X_{2n_2}$ | $\bar{X}_2 = \frac{1}{n_2} \sum_{j=1}^{n_2} X_{2j}$ | $S_2 = \frac{1}{n_2-1} (X_{2j} - \bar{X}_2)' (X_{2j} - \bar{X}_2)$ |
| q | $X_{q1}, X_{q2}, \dots, X_{qn_q}$ | $\bar{X}_q = \frac{1}{n_q} \sum_{j=1}^{n_q} X_{qj}$ | $S_q = \frac{1}{n_q-1} (X_{qj} - \bar{X}_q)' (X_{qj} - \bar{X}_q)$ |

We wish to make inference about $(\text{mean vector of popn 1}) - (\text{mean vector of popn 2}) = \mu_1 - \mu_2$.

Hypothesis: $H_0: \mu_1 = \mu_2$ against $H_a: \mu_1 \neq \mu_2$ or $H_0: \mu_1 - \mu_2 = 0$ against $H_a: \mu_1 - \mu_2 \neq 0$

Assumptions

- The sample $X_{11}, X_{12}, \dots, X_{1n_1}$ is a random sample of size n_1 from a p -variate popn with mean μ_1 and covariance matrix Σ_1 .
 - The sample $X_{21}, X_{22}, \dots, X_{2n_2}$ is a random sample of size n_2 from a p -variate popn with mean μ_2 and covariance matrix Σ_2 .
 - $X_{11}, X_{12}, \dots, X_{1n_1}$ are independent $X_{21}, X_{22}, \dots, X_{2n_2}$ when $n_1 \neq n_2$ are small.
 - Both popns are multivariate normal.
 - $\Sigma_1 = \Sigma_2$ (same variance-covariance matrix).
- The assumption that $\Sigma_1 = \Sigma_2$, is much stronger than its univariate counterpart. Here we are assuming that several pairs of variance & covariance are nearly equal. When $\Sigma_1 = \Sigma_2 = \Sigma$, $(X_{ij} - \bar{X}_i)'(X_{ij} - \bar{X}_i)$ is an estimate of $(n_1 - 1)\Sigma$ for popn 1 and $(X_{2j} - \bar{X}_2)'(X_{2j} - \bar{X}_2)$ is an estimate $(n_2 - 1)\Sigma$ for popn 2. Consequently we pool the information in both samples in order to estimate the common variance, Σ .

- We set,

$$\text{Spooled} = \frac{1}{n_1 + n_2 - 2} \left\{ (X_{ij} - \bar{X}_i)'(X_{ij} - \bar{X}_i) + (X_{2j} - \bar{X}_2)'(X_{2j} - \bar{X}_2) \right\}$$

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

- Since $(X_{ij} - \bar{X}_i)'(X_{ij} - \bar{X}_i)$ from popn 1 has $n_1 - 1$ d.f and $(X_{2j} - \bar{X}_2)'(X_{2j} - \bar{X}_2)$ from popn 2 has $n_2 - 1$ d.f, the divisor $n_1 + n_2 - 2$ is obtained by combining the 2 components.

- If $X_{11}, X_{12}, \dots, X_{1n_1}$ is a random sample of size n_1 from $N_p(\mu_1, \Sigma)$ and

$X_{21}, X_{22}, \dots, X_{2n_2}$ is an independent random sample of size n_2 from $N_p(\mu_2, \Sigma)$, then

$T^2 = [\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)]' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) \text{Spooled} \right]^{-1} [\bar{X}_1 - \bar{X}_2 - (\mu_1 - \mu_2)]$ is distributed as

$$\frac{(n_1 + n_2 - 2)p}{(n_1 + n_2 - p - 1)} F(p, n_1 + n_2 - p - 1) \quad (\text{Johnson \& Wichern, 2005})$$

$$(n_1 + n_2 - p - 1)$$

Confidence Interval for Difference in Means

- The C.I. for difference in means is given as:

$$\mu_1 - \mu_2 = \bar{X}_1 - \bar{X}_2 = \frac{(n_1 - n_2 - 2)p}{(n_1 + n_2 - p - 1)} F(p, n_1 + n_2 - p - 1) \left[\sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_{11 \text{ pooled}}} \right]$$

$$\left[\sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_{22 \text{ pooled}}} \right]$$

Where,

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

$$\left(\frac{1}{n_1} + \frac{1}{n_2} \right) = \frac{2}{15} = \frac{2}{9}$$

Exa 3.6

The data of samples from 2 popltns were collected and recorded, below for xtics $x_1 \in x_2$.

~~Exam~~

$$n_1 = 50, \bar{x}_1 = \begin{bmatrix} 8.3 \\ 4.1 \end{bmatrix}, s_1 = \begin{bmatrix} 2 & 1 \\ 1 & 6 \end{bmatrix}$$

$$n_2 = 50, \bar{x}_2 = \begin{bmatrix} 10.2 \\ 3.9 \end{bmatrix}, s_2 = \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix}$$

Determine the following:

a) Whether the popltn have the same means at $\alpha = 0.05 - 1.0.5$

b) The 95% C.I. for difference in means.

Solution:

Hypothesis: $H_0: \mu_1 = \mu_2$ against $H_a: \mu_1 \neq \mu_2$ or $H_0: \mu_1 - \mu_2 = 0$ against $H_a: \mu_1 - \mu_2 \neq 0$ at $\alpha = 0.05$

Test statistic:

$$\bar{x}_1 - \bar{x}_2 = \begin{bmatrix} -1.9 \\ 0.2 \end{bmatrix}$$

$$\left[\frac{1}{n_1} + \frac{1}{n_2} \right] S_{\text{pooled}}^2 = \begin{bmatrix} 0.5 & 3 \\ 0.5 & 2 \end{bmatrix}$$

$$\text{Spooled} = \frac{n_1 - 1}{n_1 + n_2 - 2} S_1^2 + \frac{n_2 - 1}{n_1 + n_2 - 2} S_2^2 = \frac{49}{98} \begin{bmatrix} 2 & 1 \\ 1 & 6 \end{bmatrix} + \frac{49}{98} \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 1 & 5 \end{bmatrix} S_{12}$$

$$\begin{aligned} T^2 &= [\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)]' \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) S_{\text{pooled}} \right]^{-1} [\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)] \\ &= [-1.9 \ 0.2] \left[\left(\frac{1}{n_1} + \frac{1}{n_2} \right) \begin{bmatrix} 2 & 1 \\ 1 & 5 \end{bmatrix} \right]^{-1} \begin{bmatrix} -1.9 \\ 0.2 \end{bmatrix} \\ &= \frac{25}{9} \begin{bmatrix} -1.9 & 0.2 \\ 0.2 & 1 \end{bmatrix}^{-1} \begin{bmatrix} -1.9 \\ 0.2 \end{bmatrix} \\ &= \frac{25}{9} [-9.7 \ 12.3] \begin{bmatrix} -1.9 \\ 0.2 \end{bmatrix} \\ &= 52.47 \end{aligned}$$

Tabulated value

$$\begin{aligned} T^2 &= (n_1 + n_2 - 2) p F_{\alpha}(p, n_1 + n_2 - p - 1) \\ &= (n_1 + n_2 - p - 1) \\ &= 98(2) F_{0.05}(2, 97) = 98(2)(3.09) = 6.24 \end{aligned}$$

Conclusion:

Since $T^2 > 6.24$, reject H_0 and conclude that the popltn do not have the same means at $\alpha = 0.05 - 1.0.5$

b) 95% C.I

$$\begin{aligned} \mu_1 - \mu_2 &= \bar{x}_1 - \bar{x}_2 = \frac{(n_1 + n_2 - 2)p}{\sqrt{(n_1 + n_2 - p - 1)}} f_{\alpha}(p, n_1 + n_2 - p - 1) \left[\sqrt{\frac{1}{n_1} + \frac{1}{n_2}} S_{11 \text{ pooled}} \right] \\ &= \left[\frac{-1.9}{0.2} \right] \pm \sqrt{\frac{98(2)}{97}} f_{0.05}(2, 97) \left[\sqrt{2 \left(\frac{1}{50} + \frac{1}{50} \right)} \right] = \left[\frac{-1.9}{0.2} \right] \pm \left[\sqrt{6.24(2) \left(\frac{1}{50} \right)} \right] 0.707 \\ &= \left[-1.9 - 0.707 \right] \pm \left[-1.9 + \left[\sqrt{6.24(5) \left(\frac{1}{50} \right)} \right] 1.117 \right] \\ &= \left[-2.606 \leq \mu_1 - \mu_2 \leq -1.193 \right] \\ &\quad \left[-0.917 \leq \mu_1 - \mu_2 \leq 1.317 \right] \end{aligned}$$

*Exercise.

4: PRINCIPAL COMPONENTS, DISCRIMINANT FUNCTIONS AND CANONICAL CORRELATION.

Principal Components.

A principal component analysis is concerned with explaining the variance-covariance structure of variables through a few linear combination of these variables. The general objectives include the following:

a) Data reduction

b) Interpretation.

Although p components are required to reproduce the total system variability, often much of this variability can be accounted for by a small number, k , of principal components. If so, there is almost as much information in the k components as there is in the original p variables. The k principal components can then replace the initial p variables, and the original data set, consisting of n measurements on p variables, is reduced to a data set containing of n measurements on k principal components.

Population Principal Components.

Algebraically, principal components are particular linear combinations of the p r.v.s x_1, x_2, \dots, x_p . Geometrically, these linear combinations represent a selection of a new coordinate system obtained by rotating the original system with x_1, x_2, \dots, x_p as the coordinate axes. The new axes represent the directions with maximum variability and provide a simpler and more parsimonious description of the covariate structure. The development does not require a multivariate normal assumption.

Let the random vector $X' = [x_1, x_2, \dots, x_p]$ have a variance-covariance matrix Σ with Eigen values $\lambda_1 \geq \lambda_2 \geq \dots, \lambda_p \geq 0$. Consider a linear combination

$$Y_1 = a_1' X = a_{11} x_1 + a_{12} x_2 + \dots + a_{1p} x_p$$

$$Y_2 = a_2' X = a_{21} x_1 + a_{22} x_2 + \dots + a_{2p} x_p$$

⋮

$$Y_p = a_p' X = a_{p1} x_1 + a_{p2} x_2 + \dots + a_{pp} x_p$$

We obtain;

$$V(Y_i) = a_i' \Sigma a_i, \quad i = 1, 2, \dots, p \quad V(Y_i) = q_i' \Sigma q_i$$

$$\text{Cov}(Y_i, Y_k) = a_i' \Sigma a_k, \quad i, k = 1, 2, \dots, p \quad \text{Cov}(Y_i, Y_k) = a_i' \Sigma a_k$$

- The principal components are uncorrelated linear combinations Y_1, Y_2, \dots, Y_p whose variances are large as possible.

~~First principal components = linear combination $a_1' X$ that maximizes $V(a_1' X)$ subject to $a_1' a_1 = 1$~~

~~Second " " = linear combination $a_2' X$ that maximizes $V(a_2' X)$ subject to $a_2' a_2 = 1$
and $\text{Cov}(a_1' X, a_2' X) = \text{Cov}(Y_1, Y_2) = 0$~~

~~ith " " = linear combination $a_i' X$ that maximizes $V(a_i' X)$ subject to $a_i' a_i = 1$
and $\text{Cov}(a_i' X, a_k' X) = \text{Cov}(Y_i, Y_k) = 0$ for $k < i$.~~

Ex 4.1

Determine the principal components if the random variables x_1, x_2, x_3 have a variance-covariance

Xam

$$\text{matrix } \Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

Solution:

Determine the eigen values using the relation $A = \lambda$, such that $|A - \lambda I| = 0$

$$\begin{vmatrix} 1-\lambda & -2 & 0 \\ -2 & 5-\lambda & 0 \\ 0 & 0 & 2-\lambda \end{vmatrix} = (1-\lambda)[(5-\lambda)(2-\lambda)] - 2(2)(2-\lambda) = 0$$

$$= 2\lambda^2 - 12\lambda + 2 - \lambda^3 + 6\lambda^2 - \lambda = -\lambda^3 + 8\lambda^2 - 12\lambda + 2$$

$$(2-\lambda)(5-6\lambda+\lambda^2-4) = (2-\lambda)(\lambda^2-6\lambda+1) = 0$$

$$2(\lambda^2-6\lambda+1) - \lambda(\lambda^2-6\lambda+1) = 2\lambda^2 - 12\lambda + 2 - \lambda^3 + 6\lambda^2 - \lambda$$

$$\lambda_1 = 2 \quad \text{and,} \quad \lambda_2 = \frac{-2\lambda^2 + 12\lambda + 2 - \lambda^3 + 6\lambda^2 - \lambda}{-\lambda^3 + 8\lambda^2 - 13\lambda + 2}$$

$$\lambda = \frac{6 \pm \sqrt{36-4}}{2} = 3 \pm \sqrt{8} = 3 \pm 2.828$$

such that,

$$\lambda_2 = 5.828, \lambda_3 = 0.172$$

Eigen Vectors

$$Ax = \lambda x$$

$$a) \lambda_1 = 5.828$$

$$\begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 5.828 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$1x_1 - 2x_2 = 5.828x_1 \text{ such that } 4.828x_1 = -2x_2 \quad \dots \text{(i)}$$

$$-2x_1 + 5x_2 = 5.828x_2 \text{ such that } 2x_1 = 0.828x_2 \quad \dots \text{(ii)}$$

$$2x_3 = 5.828x_3 \text{ such that } x_3 = 0 \quad \dots \text{(iii)}$$

consider eqtn (i)

$$1x_1 = -2x_2 \quad x_1 = -0.414x_2 \quad \dots \text{(iv)} \quad x_1 = -2x_2 \quad x_2 =$$

$\frac{4.828}{4.828}$

Let $x_2 = 1$ such that $x_1 = -0.414$ then,

$$e_1 = \begin{bmatrix} -0.414 \\ 1.000 \\ 0.000 \end{bmatrix}$$

$$e_1 = \begin{bmatrix} -0.383 \\ 0.924 \\ 0.000 \end{bmatrix}$$

$$e = \frac{x}{\sqrt{x^T x}} = \frac{0.414}{\sqrt{0.414^2 + 1^2}} = \frac{1}{\sqrt{0.414^2 + 1^2}}$$

$$b) \lambda_2 = 2$$

$$\begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 2 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$1x_1 - 2x_2 = 2x_1 \text{ such that } x_1 = -2x_2 \quad \dots \text{(i)}$$

$$-2x_1 + 5x_2 = 2x_2 \text{ such that } 2x_1 = 3x_2 \quad \dots \text{(ii)}$$

$$2x_3 = 2x_3 \text{ such that } x_3 = 0 \quad \dots \text{(iii)}$$

consider eqtn (ii)

Let $x_2 = 1$

consider eqtn (i)

$$x_1 = \frac{3}{2}x_2 \quad \dots \text{(iv)}$$

Let $x_2 = 0$ such that $x_1 = 0$ then, $e_2 =$

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

$$c) \lambda_3 = 0.172$$

$$\begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0.172 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$1x_1 - 2x_2 = 0.172x_1 \text{ such that } 0.828x_1 = 2x_2 \quad \dots \text{(i)}$$

$$-2x_1 + 5x_2 = 0.172x_2 \quad " " \quad 2x_1 = 4.828x_2 \quad \dots \text{(ii)}$$

$$2x_3 = 0.172x_3 \quad " " \quad x_3 = 0 \quad \dots \text{(iii)}$$

consider eqtn (i)

$$1x_2 = 0.828x_1 = 0.414x_1 \quad \dots \text{(iv)}$$

2

Let $x_1 = 1$ such that $x_2 = 0.414$ then,

$$C_3 = \begin{bmatrix} 1.000 \\ 0.414 \\ 0.000 \end{bmatrix} \text{ such that } e_3 = \begin{bmatrix} 0.924 \\ 0.383 \\ 0.000 \end{bmatrix} \text{ is unit vector.}$$

Therefore, the principal components for $\lambda_1 \geq \lambda_2 \geq \lambda_3$ are

$$Y_1 = e_1^T X = -0.383x_1 + 0.924x_2$$

$$Y_2 = e_2^T X = x_3$$

$$Y_3 = e_3^T X = 0.924x_1 + 0.383x_2$$

The variable x_3 is one of the principal components because it is uncorrelated with the other two variables.

$$\begin{aligned} V(Y_1) &= V(-0.383x_1 + 0.924x_2) = V\left\{[-0.383 \ 0.924] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right\} \\ &= [-0.383 \ 0.924] \begin{bmatrix} 1 & -2 \\ -2 & 5 \end{bmatrix} \begin{bmatrix} -0.383 \\ 0.924 \end{bmatrix} = [-0.383 \ 0.924] \begin{bmatrix} -2.231 \\ 5.386 \end{bmatrix} \\ &= 5.83 \end{aligned}$$

$$V(Y_2) = V(x_3) = 2.00$$

$$\begin{aligned} \text{Cov}(Y_1, Y_2) &= \text{Cov}(-0.383x_1 + 0.924x_2, x_3) = \text{Cov}\left\{[-0.383 \ 0.924] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, x_3\right\} \\ &= [-0.383 \ 0.924] \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0 \end{aligned}$$

$$\begin{aligned} V(Y_3) &= V(0.924x_1 + 0.383x_2) = V\left\{[0.924 \ 0.383] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right\} \\ &= [0.924 \ 0.383] \begin{bmatrix} 1 & -2 \\ -2 & 5 \end{bmatrix} \begin{bmatrix} 0.924 \\ 0.383 \end{bmatrix} = [0.924 \ 0.383] \begin{bmatrix} 0.158 \\ 0.061 \end{bmatrix} \\ &= 0.172. \end{aligned}$$

It is also readily apparent that,

$$\sigma_{11} + \sigma_{22} + \sigma_{33} = 1 + 5 + 2 = 5.828 + 2 + 0.172$$

$$= 8$$

The proportion of total variances accounted for by,

$$\text{The first principal component is } \frac{\lambda_1}{\lambda_1 + \lambda_2 + \lambda_3} = \frac{5.828}{5.828 + 2 + 0.172} = 0.73$$

$$\text{The first 2 principal components is } \frac{\lambda_1 + \lambda_2}{\lambda_1 + \lambda_2 + \lambda_3} = \frac{5.83 + 2}{5.828 + 2 + 0.172} = 0.98 \text{ of the population variances.}$$

In this case, the components γ_1 and γ_2 could replace the original 3 variables with little loss of information.

$$p\gamma_1 X_1 = \frac{c_{11}\sqrt{\lambda_1}}{\sqrt{\sigma_{11}}} = \frac{-0.383\sqrt{5.83}}{\sqrt{1}} = -0.925$$

$$p\gamma_1 X_2 = \frac{c_{12}\sqrt{\lambda_1}}{\sqrt{\sigma_{22}}} = \frac{0.924\sqrt{5.83}}{\sqrt{5}} = 0.998$$

Remarks:

- The variable X_2 , with coefficient 0.998, records the greatest weight in component γ_1 . It also has the largest correlations (in absolute value) with γ_1 .
- The correlations of X_1 with γ_1 , -0.925 is almost as large as that for X_2 , indicating that the variables are about equally important to the principal components. The relative sizes of the coefficients of X_1 & X_2 suggest, however, that X_2 contributes more to the determination of γ_1 than does X_1 . Since, in this case both coefficients are reasonably large and they have opposite signs, one would argue that both variables aid in the interpretation of γ_2 .

• Finally $p\gamma_2 X_1 = \frac{c_{21}\sqrt{\lambda_2}}{\sqrt{\sigma_{11}}} = 0$, $p\gamma_2 X_2 = \frac{c_{22}\sqrt{\lambda_2}}{\sqrt{\sigma_{22}}} = 0$ & $p\gamma_2 X_3 = \frac{c_{23}\sqrt{\lambda_2}}{\sqrt{\sigma_{33}}} = \frac{\sqrt{2}}{\sqrt{2}} = 1$ (as it should)

- The remaining correlations can be neglected since the 3rd component is unimportant.

- Variables should probably be standardized if they are measured on scales with widely different ranges or if the units of measurement are not commensurate. If variables are standardized their subsequent magnitudes will be of the same order.

- The principal components derived from Σ are different from those derived from R . Furthermore, one set of principal components is not a simple function of the other. This suggests that standardization is not inconsequential.

Ex 4.2

The data of the variables X_1 , X_2 & X_3 were collected & recorded.

Determine the principal components, loadings, scores and interpret the results.

Solution:

| x_1 | x_2 | x_3 | x_1^2 | x_1x_2 | x_1x_3 | x_2^2 | x_2x_3 | x_3^2 |
|-------|-------|-------|---------|----------|----------|---------|----------|---------|
| 7 | 4 | 2 | 49 | 28 | 14 | 16 | 8 | 4 |
| 5 | 6 | 1 | 25 | 30 | 5 | 36 | 6 | 1 |
| 4 | 4 | 2 | 16 | 16 | 8 | 16 | 8 | 4 |
| 4 | 4 | 2 | 16 | 16 | 8 | 16 | 8 | 4 |
| 5 | 2 | 3 | 25 | 10 | 15 | 4 | 6 | 9 |
| 25 | 20 | 10 | 131 | 100 | 50 | 88 | 36 | 22 |

$$\bar{x}_1 = \frac{25}{5} = 5.0, \bar{x}_2 = \frac{20}{5} = 4.0, \bar{x}_3 = \frac{10}{5} = 2.0$$

$$\bar{x} = \begin{bmatrix} 5 \\ 4 \\ 2 \end{bmatrix}$$

$$\text{Covariance; } S_{ik} = \frac{1}{n-1} \left\{ \sum_{j=1}^n x_{ij} x_{kj} - n \bar{x}_i \bar{x}_k \right\}$$

$$S_{11} = \frac{1}{4} \{ 131 - 5(5)^2 \} = 1.5, S_{12} = \frac{1}{4} \{ 100 - 5(5)(4) \} = 0$$

$$S_{13} = \frac{1}{4} \{ 50 - 5(5)(2) \} = 0, S_{22} = \frac{1}{4} \{ 88 - 5(4)^2 \} = 2$$

$$S_{23} = \frac{1}{4} \{ 36 - 5(4)(2) \} = -1, S_{33} = \frac{1}{4} \{ 22 - 5(2)^2 \} = 0.5$$

$$S_n = \begin{bmatrix} 1.5 & 0 & 0 \\ 0 & 2 & -1 \\ 0 & -1 & 0.5 \end{bmatrix}$$

Eigen Values

Solve for λ in the equation $|A - \lambda I| = 0$.

$$\begin{bmatrix} 1.5 - \lambda & 0 & 0 \\ 0 & 2 - \lambda & -1 \\ 0 & -1 & 0.5 - \lambda \end{bmatrix} = 0$$

$$0(0 - \lambda)(0 - 1) = 0$$

$$(1.5 - \lambda) \{ (2.0 - \lambda)(0.5 - \lambda) - 1.0 \} = 0$$

$$(1.5 - \lambda) \{ (\lambda^2 - 2.5\lambda) \} = 0$$

$$\lambda(1.5 - \lambda)(\lambda - 2.5) = 0$$

$$\lambda_1 = 2.5, \lambda_2 = 1.5, \lambda_3 = 0$$

Eigen Vectors

Solve for x in the relation $Ax = \lambda x$

$$\lambda_1 = 2.5$$

$$\begin{bmatrix} 1.5 & 0 & 0 \\ 0 & 2 & -1 \\ 0 & -1 & 0.5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 2.5 \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$1.5X_1 + 0X_2 + 0X_3 = 2.5X_1 \quad \text{or} \quad X_1 = 0 \quad \dots \text{(i)}$$

$$0X_1 + 2X_2 - 1X_3 = 2.5X_2 \quad " \quad 0.5X_2 = -1X_3 \quad \dots \text{(ii)}$$

$$0X_1 - 1X_2 + 0.5X_3 = 2.5X_3 \quad " \quad X_2 = 2X_3 \quad \dots \text{(iii)}$$

From eqtn (i), $X_1 = 0$

From eqtn (ii)

$$\text{Let } X_3 = -1 \text{ such that } X_2 = 2 \text{ then, } C_1 = \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix} = \begin{bmatrix} 0 \\ 2/\sqrt{5} \\ 1/\sqrt{5} \end{bmatrix} = \begin{bmatrix} 0 \\ 0.894 \\ -0.447 \end{bmatrix} \quad e = \frac{x}{\sqrt{x^T x}}$$

$$\lambda_2 = 1.5$$

$$\begin{bmatrix} 1.5 & 0 & 0 \\ 0 & 2 & -1 \\ 0 & -1 & 0.5 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = 1.5 \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$$

$$1.5X_1 + 0X_2 + 0X_3 = 1.5X_1 \quad \text{or} \quad 0X_1 = 0 \quad \dots \text{(i)}$$

$$0X_1 + 2X_2 - 1X_3 = 1.5X_2 \quad \text{or} \quad 0.5X_2 = 1X_3 \quad \dots \text{(ii)}$$

$$0X_1 - 1X_2 + 0.5X_3 = 1.5X_3 \quad \text{or} \quad 1X_2 = -1X_3 \quad \dots \text{(iii)}$$

From eqtn (i), Let $X_1 = 1$

From eqtn (ii)

$$\text{Let } X_2 = 0 \text{ then } X_3 = 0 \text{ such that } C_2 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$\lambda_3 = 0$$

$$\begin{bmatrix} 1.5 & 0 & 0 \\ 0 & 2 & -1 \\ 0 & -1 & 0.5 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix} = 0 \begin{bmatrix} X_1 \\ X_2 \\ X_3 \end{bmatrix}$$

$$1.5X_1 = 0 \quad \text{or} \quad 1.5X_1 = 0 \quad \dots \text{(i)}$$

$$2X_2 - X_3 = 0 \quad \text{or} \quad 2X_2 = X_3 \quad \dots \text{(ii)}$$

$$-X_2 + 0.5X_3 = 0 \quad \text{or} \quad X_2 = 0.5X_3 \quad \dots \text{(iii)}$$

From eqtn (i), $X_1 = 0$

From eqtn (ii)

$$\text{Let } X_2 = 1 \text{ such that } X_3 = 2 \text{ then } C_3 = \begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1/\sqrt{5} \\ 2/\sqrt{5} \end{bmatrix} = \begin{bmatrix} 0 \\ 0.447 \\ 0.894 \end{bmatrix}$$

Principal components

$$Y_1 = C_1^T X = 0.894X_2 - 0.447X_3$$

$$Y_2 = C_2^T X = X_1$$

$$Y_3 = C_3^T X = 0.447X_2 + 0.894X_3$$

components

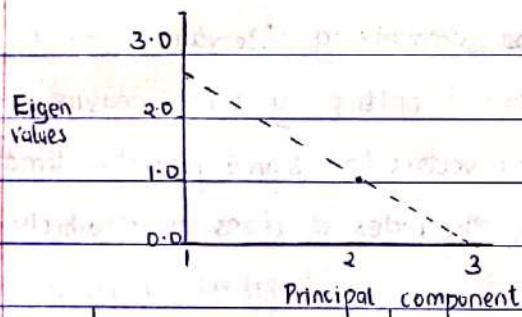
The principal loadings are the coefficients of X in the principal components or simply Eigen vectors that are given as:

| PC ₁ | PC ₂ | PC ₃ |
|-----------------|-----------------|-----------------|
| 0.000 | 1.000 | 0.000 |
| 0.894 | 0.000 | 0.447 |
| -0.447 | 0.000 | 0.894 |

- The principal component one indicates that there is a factor that affects X_2 positively while affecting X_3 negatively whereas principal component γ_2 indicates that there is a factor that is affecting X_1 positively and does not consider any measure of $X_2 \& X_3$ therefore γ_2 measures the scores of X_1 .

- Considering the scree diagram below the Eigen values of $\lambda_1 = 5/2$ and $\lambda_2 = 3/2$ have values more than 1, then the principal components of $\lambda_1 \& \lambda_2$ are considered sufficient enough in providing enough info. Alternatively the contribution of the variances can be used by considering the relative cumulative variances.

Scree Plot



| λ | Cumulative Variance | % Cumulative Variance |
|-----------|---------------------|-----------------------|
| 2.5 | 2.5 | 0.625 |
| 1.5 | 4.0 | 1.000 |
| 0 | 4.0 | 1.000 |

The first principal components contributes 62.5% of the variations and this is considered lower and therefore more than one principal components are required in the analysis in order to have at least 90% of the variance explained by the principal components. The first 2 principal components have 100% of the variance explained by the " " " that is more than minimum of 90% and therefore the first 2 principal components are considered sufficient.

- The principal components scores are computed by first transforming the initial data of the variables by subtracting the means from the initial values then multiplying by the principal component loadings, that is,

$$\text{Principal component score, } PC_{1i} = [X_{1j} - \bar{X}_1 \quad X_{2j} - \bar{X}_2 \quad X_{3j} - \bar{X}_3] [e_1]$$

$$\text{ " " " , } PC_{2i} = [X_{1j} - \bar{X}_1 \quad X_{2j} - \bar{X}_2 \quad X_{3j} - \bar{X}_3] [e_2]$$

$$\text{ " " " , } PC_{3i} = [X_{1j} - \bar{X}_1 \quad X_{2j} - \bar{X}_2 \quad X_{3j} - \bar{X}_3] [e_3]$$

Deviations from \bar{x} mean

Loadings

Scores

| X_1 | X_2 | X_3 | C_1 | C_2 | C_3 | PC_1 | PC_2 | PC_3 |
|-------|-------|-------|--------|-------|-------|--------|--------|--------|
| 2.0 | 0.0 | 0.0 | 0.000 | 1.000 | 0.000 | 0.000 | 2.000 | 0.000 |
| 0.0 | 2.0 | -1.0 | -0.894 | 0.000 | 0.447 | -2.236 | 0.000 | 0.000 |
| -1.0 | 0.0 | 0.0 | 0.447 | 0.000 | 0.894 | 0.000 | -1.000 | 0.000 |
| -1.0 | 0.0 | 0.0 | 0.447 | 0.000 | 0.894 | 0.000 | -1.000 | 0.000 |
| 0.0 | -2.0 | 1.0 | -0.894 | 0.000 | 0.447 | 2.236 | 0.000 | 0.000 |

Exam

Discriminant Functions.

Consider random samples of N_1 & N_2 observation vectors drawn independently from respective p -dimensional multi-normal popns with mean vectors μ_1 & μ_2 , and a common variance-covariance matrix Σ . Rather than test the usual hypothesis of equal mean vectors we wish to construct a linear compound ^{or} index of summarizing observations from \bar{e} groups on a one-dimensional scale that discriminates btwn \bar{e} popns by some measure of maximal separation. If \bar{X}_1 and \bar{X}_2 are \bar{e} sample mean vectors an S is the pooled estimate of Σ , we shall determine \bar{e} coefficient vector a of the index $a'x$ as that which gives \bar{e} greatest squared critical ratio.

$$t^2(a) = [a'(\bar{X}_1 - \bar{X}_2)]^2 / N_1 N_2 / (N_1 + N_2) \quad \dots (4.1)$$

$$a'Sa$$

Equivalently, that which maximizes \bar{e} absolute difference $|a'(\bar{X}_1 - \bar{X}_2)|$ in \bar{e} ave. values of the index of \bar{e} two groups subject to \bar{e} constraint $a'Sa = 1$. As in \bar{e} union intersection construction of \bar{e} single-sample T^2 , the coefficient vector a is given by \bar{e} homogeneous system of eqtns

$$N_1 N_2 [(\bar{X}_1 - \bar{X}_2)(\bar{X}_1 - \bar{X}_2)' - \lambda S] a = 0 \quad \dots (4.2)$$

$$N_1 + N_2$$

Such that

$$\lambda = \max t^2(a) = N_1 N_2 / (N_1 + N_2) (\bar{X}_1 - \bar{X}_2)' S^{-1} (\bar{X}_1 - \bar{X}_2) = T^2 \quad \dots (4.3)$$

The rank of \bar{e} coefficient matrix is $p-1$, so that \bar{e} system has only \bar{e} single solution

$$a = S^{-1}(\bar{X}_1 - \bar{X}_2) \quad \dots (4.4)$$

The linear discriminant function is given as:

$$y = (\bar{X}_1 - \bar{X}_2)' S^{-1} x \quad \dots (4.5)$$

If \bar{e} variances of \bar{e} responses are nearly equal, the elements give \bar{e} relative importance of \bar{e} contribution of each response to \bar{e} T^2 statistic.

In order to use a linear discriminant function for classifying an observation of an unknown popn we begin by comparing the mean values of the scores of the 2 samples.

$$\bar{y}_1 = (\bar{x}_1 - \bar{x}_2)' S^{-1} \bar{x}_1 \text{ and } \bar{y}_2 = (\bar{x}_1 - \bar{x}_2)' S^{-1} \bar{x}_2$$

The midpoint of these 2 means on the discriminant func scale is $\frac{1}{2}(\bar{x}_1 - \bar{x}_2)' S^{-1} (\bar{x}_1 + \bar{x}_2)$ and we might adopt the classification rule, assign an individual with an observation x to popn one (1) if $(\bar{x}_1 - \bar{x}_2)' S^{-1} x > \frac{1}{2}(\bar{x}_1 - \bar{x}_2)' S^{-1} (\bar{x}_1 + \bar{x}_2)$ and to popn two (2) if $(\bar{x}_1 - \bar{x}_2)' S^{-1} x \leq \frac{1}{2}(\bar{x}_1 - \bar{x}_2)' S^{-1} (\bar{x}_1 + \bar{x}_2)$

That is, sampling units are assigned to the group with the closer discriminant score. Because the mean point is a value of a random var. it would seem more appropriate to state that rule in terms of a single statistic.

$$W = x' S^{-1} (\bar{x}_1 - \bar{x}_2) - \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' S^{-1} (\bar{x}_1 + \bar{x}_2)$$

That is, assign x to popn one (1) if $W > 0$, otherwise assign x to popn two (2), W is called the Wald - Anderson classification statistic.

Exa 4.3

49 women participating in study of human aging were classified into diagnostic categories "civilized" and "no civilization" in the basis of an intensive psychiatric examination. Determine the classification criterion using the discriminant function for the results recorded in the table below.

| Subset | Group Mean | |
|---------------------|-----------------|--------------|
| | No civilization | Civilization |
| Information | 12.75 | 8.75 |
| Similarities | 9.75 | 5.33 |
| Arithmetic | 11.49 | 8.50 |
| Picture Competition | 7.97 | 4.75 |

The within-group covariance matrix and its inverse is given as

$$S = \begin{bmatrix} 11.26 & 9.41 & 7.16 & 3.38 \\ 9.41 & 13.53 & 7.38 & 2.50 \\ 7.16 & 7.38 & 11.58 & 2.62 \\ 3.38 & 2.50 & 2.62 & 5.81 \end{bmatrix} \quad S^{-1} = \begin{bmatrix} 0.260 & -0.137 & -0.059 & -0.066 \\ -0.137 & 0.187 & -0.038 & 0.016 \\ -0.059 & -0.038 & 0.151 & -0.017 \\ -0.066 & 0.016 & -0.017 & 0.211 \end{bmatrix}$$

Solution:

The vector of the mean differences for the two diagnostic groups is

$$(\bar{x}_1 - \bar{x}_2)' = [3.82 \ 4.24 \ 2.99 \ 3.22] \text{ and discriminant function is}$$

$$y = (\bar{x}_1 - \bar{x}_2) S_{\bar{x}}^{-1} = [3.82 \ 4.24 \ 2.99 \ 3.22] \begin{bmatrix} 0.260 & -0.137 & -0.059 & -0.066 \\ -0.137 & 0.187 & -0.038 & 0.016 \\ -0.059 & -0.038 & 0.151 & -0.017 \\ -0.066 & 0.016 & -0.017 & 0.211 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix}$$

$$= 0.0259x_1 + 0.2083x_2 + 0.0086x_3 + 0.4456x_4$$

The mean discriminant scores for non-civilization factor & civilization factor groups are:

$$\bar{Y}_1 = (\bar{x}_1 - \bar{x}_2) S_{\bar{x}}^{-1} = [3.82 \ 4.24 \ 2.99 \ 3.22] \begin{bmatrix} 0.260 & -0.137 & -0.059 & -0.066 \\ -0.137 & 0.187 & -0.038 & 0.016 \\ -0.059 & -0.038 & 0.151 & -0.017 \\ -0.066 & 0.016 & -0.017 & 0.211 \end{bmatrix} \begin{bmatrix} 12.57 \\ 9.57 \\ 11.49 \\ 5.97 \end{bmatrix}$$

$$\bar{Y}_2 = (\bar{x}_1 - \bar{x}_2) S_{\bar{x}}^{-1} = [3.82 \ 4.24 \ 2.99 \ 3.22] \begin{bmatrix} 0.260 & -0.137 & -0.059 & -0.066 \\ -0.137 & 0.187 & -0.038 & 0.016 \\ -0.059 & -0.038 & 0.151 & -0.017 \\ -0.066 & 0.016 & -0.017 & 0.211 \end{bmatrix} \begin{bmatrix} 8.57 \\ 5.33 \\ 8.50 \\ 4.75 \end{bmatrix}$$

The midpoint of these 2 means on the discriminant function scale is given as

$$\frac{1}{2}(\bar{x}_1 - \bar{x}_2) S_{\bar{x}}^{-1} (\bar{x}_1 + \bar{x}_2) = \frac{1}{2}(5.97 + 3.527) = \frac{9.497}{2} = 4.75$$

Decision:

Assign i^{th} individual to civilization factor diagnostic group if $y_i \leq 4.75$ and to no-civilization factor group if $y_i > 4.75$

Canonical Correlations.

- In this section we shall consider the use of the quantities of coefficient of linear compounds and discuss the test statistic in its squared correlation sense. Suppose that the $p+q$ variates $[x_1' \ x_2']'$ of some multi-dimensional population have been divided so that their covariance matrix Σ has the partitioned form. In this section we shall drop the requirement that the distribution is multi-normal and merely specify that:
 - The elements of Σ are finite.
 - Σ is of the full rank $p+q$.
 - The first $r \leq \min(p+q)$ xtic roots of $\Sigma_{11}' \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{12}'$ are distinct.

The N observation vectors have been drawn randomly from the population and their sample covariance matrix is partitioned as:

$$S = \begin{bmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{bmatrix} \quad \dots \quad (4.6)$$

In conformance with, Σ , in an attempt to reduce the no. of variates and their non-zero correlations to more parsimonious degrees we seek to answer the following qstns.

What are the linear compounds?

$$U_1 = a_1' x_1, V_1 = b_1' x_2 \dots U_s = a_s' x_1, V_s = b_s' x_2 \dots \quad (4.7)$$

- With the property that the sample correlation of $U_1 \& V_1$ is the greatest, the sample correlation of $U_2 \& V_2$ greatest among all linear compounds uncorrelated with $U_1 \& V_1$, and so on, for all $s = \min(p+q)$ possible pairs. Clearly, the first pairs have the coefficients determined implicitly in the union-intersection test of set independence, and its squared correlation is the test statistic. We introduce the constraints in the maximization problem

$$a_i' S_{11} a_j = 0, b_i' S_{22} b_j = 0, a_i' S_{12} b_j = 0, a_j' S_{12} b_i = 0, i \neq j \quad (4.8)$$

- It can be shown in the maximization problem that the coefficients of the i^{th} pair are given by the homogeneous linear eqtn

$$(S_{12} S_{22}^{-1} S_{12}' - c_i S_{11}) a_i = 0, (S_{12}' S_{11}^{-1} S_{12} - c_i S_{22}) b_i = 0$$

- Where c_i is the largest root of the determinant eqtn $|S_{12} S_{22}^{-1} S_{12}' - \lambda S_{11}| = 0$ or $|S_{12}' S_{11}^{-1} S_{12} - \lambda S_{22}| = 0$

- Rewrite c_i as the squared product-moment correlation of the i^{th} linear compounds:

$$c_i = r_{U_i, V_i}^2 = \frac{(a_i' S_{12} b_i)^2}{a_i' S_{11} a_i b_i' S_{22} b_i}$$

- If c_i is a distinct root, the coefficient vectors $a_i \& b_i$ will be unique, and their linear compounds will be uncorrelated with the other canonical variates. In order to demonstrate this property we assume that $c_i \neq c_j$ are distinct roots & write the homogeneous eqtns for their respective pairs of coefficients as:

$$(S_{12} S_{22}^{-1} S_{12}' - c_i S_{11}) a_i = 0, (S_{12}' S_{11}^{-1} S_{12} - c_i S_{22}) b_i = 0$$

$$(S_{12} S_{22}^{-1} S_{12}' - c_j S_{11}) a_j = 0, (S_{12}' S_{11}^{-1} S_{12} - c_j S_{22}) b_j = 0$$

- Pre-multiply these equations by a_j', b_j', a_i' and b_i' , respectively. Then,

$$(c_i - c_j) a_i' S_{11} a_j = 0 \text{ and } (c_i - c_j) b_i' S_{22} b_j = 0$$

- Since $c_i \neq c_j$ it must follow that the bilinear forms equal to the covariance of U_i, V_j are zero. Similarly, pre-multiplication of the 1st & 4th eqtns by $b_j' S_{12} S_{11}^{-1}$ and $a_i' S_{12} S_{22}^{-1}$, respectively, and subtraction of the resulting 1st scalar eqtn from the second lead to the condition $(c_i - c_j) b_j' S_{12} S_{11}^{-1} a_i = 0$ or zero correlation of

$u_i \& v_i$ and $c_i \& c_j$ are distinct. If, as in the population, the r largest c_i are distinct, r pairs of unique canonical variables are formed.

Summary:

- We begin with a sample whose observation vectors have the covariance matrix $\Sigma = \Sigma_0$. Through the canonical variate transformation we have passed to the new scores $u_1, \dots, u_s, v_1, \dots, v_s$ with conditional matrix

$$\begin{bmatrix} 1 & \cdots & 0 & c_1^{-1/2} & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & 1 & 0 & \cdots & c_s^{-1/2} \\ c_1^{-1/2} & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & \cdots & c_s^{-1/2} & 0 & \cdots & 1 \end{bmatrix}$$

- All the correlation btwn \bar{e} sets of \bar{e} original variables has been channelled thro' the s canonical correlations. Identical canonical correlation will be obtained for $S_{11}^{-1}S_{12}S_{22}^{-1}S_{12}$ or $R_{11}^{-1}R_{12}R_{22}^{-1}R_{12}$. In the first case the elements of a_i and b_i will be in units proportional to those of the respective responses in \bar{e} sets and \bar{e} dimensions of $u_i \& v_i$ will be meaningless. Canonical variates based on \bar{e} correlation matrix and dimensionless should be evaluated in terms of the standard score $Z_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j}$ of \bar{e} original observations.

Exq 4.4

Determine the canonical correlation and variates for the following correlation matrix which has been partitioned into four matrices of order 2×2 :

$$R = \begin{bmatrix} 1.00 & 0.45 & \cdots & -0.19 & 0.43 \\ 0.45 & 1.00 & \cdots & -0.02 & 0.62 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ -0.19 & -0.02 & 1.00 & -0.29 & \\ 0.43 & 0.62 & -0.29 & 1.00 & \end{bmatrix}$$

Solution

$$R_{11}^{-1} = \frac{1}{1 - 0.45^2} \begin{bmatrix} 1.00 & -0.45 \\ -0.45 & 1.00 \end{bmatrix} = \begin{bmatrix} 1.254 & -0.564 \\ -0.564 & 1.254 \end{bmatrix}$$

$$R_{22}^{-1} = \frac{1}{1 - 0.29^2} \begin{bmatrix} 1.00 & 0.29 \\ 0.29 & 1.00 \end{bmatrix} = \begin{bmatrix} 1.092 & 0.317 \\ 0.317 & 1.092 \end{bmatrix}$$

$$R_{11}^{-1} R_{12} R_{22}^{-1} R_{12}^T = \begin{bmatrix} 1.254 & -0.564 \\ -0.564 & 1.254 \end{bmatrix} \begin{bmatrix} -0.19 & 0.43 \\ -0.02 & 0.62 \end{bmatrix} \begin{bmatrix} 1.092 & 0.317 \\ 0.317 & 1.092 \end{bmatrix} \begin{bmatrix} -0.19 & -0.02 \\ 0.43 & 0.62 \end{bmatrix}$$

$$= \begin{bmatrix} 0.094 & 0.088 \\ 0.213 & 0.373 \end{bmatrix}$$

$$|A - \lambda I| = \begin{vmatrix} 0.094 - \lambda & 0.087 \\ 0.213 & 0.373 - \lambda \end{vmatrix} = 0$$

$$(0.094 - \lambda)(0.373 - \lambda) - (0.213)(0.087) = \lambda^2 - 0.467\lambda + 0.0165 = 0$$

$$\lambda = 0.467 \pm \sqrt{0.152}$$

2

$$\lambda_1 = 0.428, \lambda_2 = 0.038,$$

The first xtic, $\lambda_1 = 0.428$ and the largest canonical correlation is $\sqrt{0.428} = 0.655$

The second xtic, $\lambda_2 = 0.038$ and the second " " " $\sqrt{0.038} = 0.195$

The 1st coefficients of the canonical associated with 1st xtic root are given by the equations

$$R_{12} R_{22}^{-1} R_{12}^T = \begin{bmatrix} -0.19 & 0.43 \\ -0.02 & 0.62 \end{bmatrix} \begin{bmatrix} 1.092 & 0.317 \\ 0.317 & 1.092 \end{bmatrix} \begin{bmatrix} -0.19 & -0.02 \\ 0.43 & 0.62 \end{bmatrix}$$

$$= \begin{bmatrix} 0.190 & 0.255 \\ 0.255 & 0.412 \end{bmatrix}$$

$$\begin{bmatrix} 0.190 & 0.255 \\ 0.255 & 0.412 \end{bmatrix} \begin{bmatrix} q_{11} \\ q_{12} \end{bmatrix} = 0.428 \begin{bmatrix} 1.00 & 0.45 \\ 0.45 & 1.00 \end{bmatrix} \begin{bmatrix} q_{11} \\ q_{12} \end{bmatrix}$$

Then,

$$q_{11} = 0.261 q_{12}$$

$$0.190 q_{11} + 0.255 q_{12} = 0.428 q_{11} + 0.193 q_{12} \quad \text{or} \quad 0.238 q_{11} = 0.062 q_{12} = 0 \quad q_{11} = 0.261$$

$$0.255 q_{11} + 0.413 q_{12} = 0.193 q_{11} + 0.428 q_{12} \quad \text{or} \quad 0.062 q_{11} + 0.016 q_{12} = 0 \quad q_{12} = 1.$$

Let $q_{12} = 1$, then $q_{11} = 0.261$ or $q_{11} = 0.255$ and $q_{12} = 0.967$

The 2nd coefficients of canonical associated with the 1st xtic root are given by the eqns:

$$R_{12}^T R_{11}^{-1} R_{12} = \begin{bmatrix} -0.19 & -0.02 \\ 0.43 & 0.62 \end{bmatrix} \begin{bmatrix} 1.254 & -0.564 \\ -0.564 & 1.254 \end{bmatrix} \begin{bmatrix} -0.19 & 0.43 \\ -0.02 & 0.62 \end{bmatrix}$$

$$= \begin{bmatrix} 0.041 & -0.047 \\ -0.047 & 0.413 \end{bmatrix}$$

$$\begin{bmatrix} 0.041 & -0.047 \\ -0.047 & 0.413 \end{bmatrix} \begin{bmatrix} b_{11} \\ b_{12} \end{bmatrix} = 0.429 \begin{bmatrix} 1.00 & -0.29 \\ -0.29 & 1.00 \end{bmatrix} \begin{bmatrix} b_{11} \\ b_{12} \end{bmatrix}$$

Then

$$0.041b_{11} - 0.047b_{12} = 0.429b_{11} - 0.124b_{12} \text{ or } 0.388b_{11} - 0.077b_{12} = 0$$

$$-0.047b_{11} + 0.413b_{12} = -0.124b_{11} + 0.428b_{12} \text{ or } 0.077b_{11} - 0.015b_{12} = 0$$

$$\text{Let } b_{12} = 1 \text{ then } b_{11} = 0.198 \text{ or } b_{11} = 0.194 \text{ and } b_{12} = 0.981$$

Replacement of 0.428 by the smaller root 0.388 in the preceding systems of equations to solve for coefficients for the 2nd canonical variates

The 1st coefficients of canonical associated with 2nd characteristic root are given by the equations

$$\begin{aligned} R_{12}R_{22}^{-1}R_{12}' &= \begin{bmatrix} -0.19 & 0.43 \\ -0.02 & 0.62 \end{bmatrix} \begin{bmatrix} 1.092 & 0.317 \\ 0.317 & 1.092 \end{bmatrix} \begin{bmatrix} -0.19 & -0.02 \\ 0.43 & 0.62 \end{bmatrix} \\ &= \begin{bmatrix} 0.190 & 0.255 \\ 0.255 & 0.412 \end{bmatrix} \end{aligned}$$

$$\begin{bmatrix} 0.190 & 0.255 \\ 0.255 & 0.412 \end{bmatrix} \begin{bmatrix} q_{21} \\ q_{22} \end{bmatrix} = 0.038 \begin{bmatrix} 1.00 & 0.45 \\ 0.45 & 1.00 \end{bmatrix} \begin{bmatrix} q_{21} \\ q_{22} \end{bmatrix}$$

Then

$$0.189q_{21} + 0.255q_{22} = 0.038q_{21} + 0.017q_{22} \text{ or } 0.151q_{21} + 0.238q_{22} = 0$$

$$0.255q_{21} + 0.413q_{22} = 0.017q_{21} + 0.038q_{22} \text{ or } 0.238q_{21} + 0.375q_{22} = 0$$

$$\text{Let } q_{21} = -1.00 \text{ then } q_{22} = 0.63 \text{ or } q_{21} = -0.844 \text{ then } q_{22} = 0.93$$

The second coefficients of canonical associated with the 2nd xtic root are given by the equations

$$\begin{aligned} R_{12}'R_{11}^{-1}R_{12} &= \begin{bmatrix} -0.19 & -0.02 \\ 0.43 & 0.62 \end{bmatrix} \begin{bmatrix} 1.254 & -0.564 \\ -0.564 & 1.254 \end{bmatrix} \begin{bmatrix} -0.19 & 0.43 \\ -0.02 & 0.62 \end{bmatrix} \\ &= \begin{bmatrix} 0.041 & -0.047 \\ -0.047 & 0.413 \end{bmatrix} \end{aligned}$$

$$\begin{bmatrix} 0.041 & -0.047 \\ -0.047 & 0.413 \end{bmatrix} \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix} = 0.038 \begin{bmatrix} 1.00 & -0.29 \\ -0.29 & 1.00 \end{bmatrix} \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}$$

Then

$$0.041b_{21} - 0.047b_{22} = 0.038b_{21} - 0.011b_{22} \text{ or } 0.003b_{21} - 0.036b_{22} = 0$$

$$-0.047b_{21} + 0.413b_{22} = -0.011b_{21} + 0.038b_{22} \text{ or } -0.036b_{21} + 0.375b_{22} = 0$$

$$\text{Let } b_{21} = -1.00 \text{ then}$$

$$b_{22} = -0.083 \text{ or } b_{21} = -0.997 \text{ and } b_{22} = -0.083$$

5 FACTOR ANALYSIS.

Introduction

- In many real-life applications, the number of independent variables used in predicting a response variables will be many. The difficulties in having too many independent variables include the following:

(i) Increased computational time to get the solution ✓

(ii) Increased time in data collection ✓

(iii) Too much expenditure in data collection ✓

(iv) Presence of redundant independent variables ✓

(v) Difficulty in making inferences ✓

- In order to avoid the difficulties associated with the many input variables, factor analysis may be considered for analysis that aims at grouping the original variables into factors which underlie the input variables such that each factor will account for one or more input variables. The total number of input variables, but after performing factor analysis, the total no. of factors in the study can be reduced by dropping the insignificant factors based on certain criterion.

Exa 5.1

In the problem of studying the customer feedback about a two-wheeler produced by an industry, the marketing manager designed a questionnaire to study the customers' feedback about the two-wheel product in order to identify the factors of the study. The initial 6 variables for the study identified include the following:

(i) X_1 - Fuel efficiency ✓ (iv) X_4 - Quality of original spares ✓

(ii) X_2 - Life span of the 2-wheeler ✓ (v) X_5 - Breakdown rate ✓

(iii) X_3 - Handling convenience ✓ (vi) X_6 - Price ✓

- After administering a questionnaire among respondents, the opinion of customers can be obtained on the variables in which the range of the scores for each of the variables is assumed to be b/w 1 and 10. The score 1 means the lowest rating while 10 means the highest rating. If the application of factor analysis groups. Assume that the application of factor analysis groups these variables as follows:

(i) X_1 , X_2 , X_4 and X_5 into factor - 1

(ii) X_6 into factor - 2

(iii) X_3 into factor - 3

- If all the factors are significant, then they are retained for future analysis. A careful examination of these groupings of variables into factors reveals that the factor - 1, factor - 2 and factor - 3 can be named as technical factor, price factor and personal factor respectively. The factor analysis reveals that in future while conducting a detailed study, it is sufficient to get opinion of the customers

on the three factors which are obtained through factor analysis. If there are n variables, there will be n factors, each factor say k is represented by a linear composite. Let F_k be the linear composite of the factor k given as

$$F_k = W_{1k} X_1 + W_{2k} X_2 + W_{3k} X_3 + \dots + W_{ik} X_i + \dots + W_{nk} X_n = \sum_{i=1}^n W_{ik} X_i, k = 1, 2, \dots, n$$

Where W_{ik} is the weight of the original variable X_i in the linear composite of the factor k .

- There are several methods of factor analysis, but they do not necessarily give the same results and therefore factor analysis is not a single unique method but a set of techniques. The factor analysis methods include the centroid method, principal components method and maximum likelihood method. In this section the centroid method is discussed.

Centroid Method.

- The method maximizes the sum of absolute loadings of each factor. In this method, the coefficients of the terms in the linear composite in each factor will be either +1 or -1. This method is less cumbersome when compared to others. The commonly used terms relating to factor analysis include the following:

(i) Factor - It is an underlying dimension that accounts for several observed variables or more factors depending upon the nature of the study and the number of variables involved in it.

(ii) Factor-loadings - Are those values which explain how closely the variables are related to each one of the factors discovered.

(iii) Communality (h^2) - It indicates how much of each variable is accounted for by the underlying factor taken together. A high value of communality means that not much of the variable is left over after whatever the factors represent is taken into consideration.

(iv) Eigen value - When the sum of squared values of factor loadings relating to a factor, then such sum is referred to as Eigen value or latent root. Eigen value indicates the relative importance of each factor in accounting for a particular set of variables being analyzed.

(v) Total sum of squares - The resulting value of the sum of Eigen values of all factors.

(vi) Rotation - In the context of factor analysis is that which reveals different structures in the data and though diff. rotations give results that appear to be diff., but from a statistical point of view, all results are taken as equal. If the factors are independent orthogonal rotation is done and if the factors are correlated, an oblique rotation is made.

(vii) Factor scores - it represents the degree to which each respondent gets high scores on the group of items that load high on each factor.

Centroid Method Procedure.

- (a) Compute the matrix of correlation, R , wherein units are placed in the diagonal spaces. The product moment formula is used for working out the correlation coefficients.
- (b) If the correlation matrix obtained happens to be positive manifold (that is disregarding the diagonal elements each variable has a large sum of positive correlations than of negative correlations). The centroid method requires that the weights for all the variables be +1.0, the variables are not weighted but simply summed. In case the correlation matrix is not a positive manifold, reflections must be made before the 1st centroid factor is obtained.
- (c) The 1st centroid factor is determined as:
- The sum of coefficients (including diagonal unity) in each column of the correlation matrix is worked out.
 - The sum of these column sums (T) is obtained.
 - The sum of each column obtained in (i) above is divided by the square root of T obtained in (ii) above, resulting in what is called centroid loadings. The full set of loadings so obtained constitute the 1st centroid factor (say A).
- d) In order to obtain the 2nd centroid factor (say B), one must first obtain a matrix of residual coefficients, for this purpose, the loading for the variables on the 1st centroid factor are multiplied. This is done for all possible pairs of variables and the resulting matrix of factor cross products are named as D_1 . The D_1 is subtracted element by element from the original matrix of the correlation, R , and the result is the 1st matrix of the residual coefficients R_1 . After obtaining R_1 , one must reflect some of the variables in it, meaning thereby that some of the variables are given negative signs in the sum (this is usually done by inspection).

Exq 5.2

The random variables X_1 , X_2 and X_3 have the variance-covariance matrix $\Sigma = \begin{bmatrix} 1 & -2 & 0 \\ -2 & 5 & 0 \\ 0 & 0 & 2 \end{bmatrix}$. Determine the new factors if the factors X_1 , X_2 and X_3 are to be grouped into demographic and economical factors using the centroid factor analysis method.

Solution.

The correlation matrix R for the 3 variables is as given below:

$$r_{ik} = \frac{\sum_{j=1}^n x_{ij} x_{kj} - n \bar{x}_i \bar{x}_k}{\sqrt{\sum_{j=1}^n x_{ij}^2 - n \bar{x}_i^2} \sqrt{\sum_{j=1}^n x_{kj}^2 - n \bar{x}_k^2}}$$

$$R = \begin{bmatrix} 1.000 & -0.894 & 0.000 \\ -0.894 & 1.000 & 0.000 \\ 0.000 & 0.000 & 1.000 \end{bmatrix}$$

$$Y_{ij} = \frac{S_{ij}}{\sqrt{S_{ii} S_{jj}}}$$

Since matrix R_1 is not positive manifold, the variable X_2 is reflected to have the following correlation matrix

$$R_1 = \begin{bmatrix} 1.000 & 0.894 & 0.000 \\ 0.894 & 1.000 & 0.000 \\ 0.000 & 0.000 & 1.000 \end{bmatrix}$$

Determine the sum of columns.

$$S_i = \sum_{k=1}^3 R_{ik},$$

$$S_1 = 1.000 + 0.894 + 0.000 = 1.894$$

$$S_2 = 0.894 + 1.000 + 0.000 = 1.894$$

$$S_3 = 0.000 + 0.000 + 1.000 = 1.000$$

Determine the sum of all columns

$$T = \sum_{k=1}^3 S_i = 1.894 + 1.894 + 1.000$$

$$= 4.789$$

The 1st centroid factor A_1 is obtained by dividing column total by sq. root of sum of all columns (T), that is,

$$UL_1(j) = \frac{S_i}{\sqrt{T}} = \frac{S_i}{\sqrt{4.789}}$$

| | | | | |
|------------------------|-------|-------|-------|------------------------|
| S_i | 1.894 | 1.894 | 1.000 | $\frac{S_i}{\sqrt{T}}$ |
| $\frac{S_i}{\sqrt{T}}$ | 0.866 | 0.866 | 0.457 | ✓ |

In order to obtain the 2nd centroid factor B , develop the matrix of factor cross product Q_1

$$\begin{array}{cccc|c}
& \downarrow \rightarrow & 0.866 & -0.866 & 0.457 & \\
0.866 & & 0.749 & -0.749 & 0.396 & \left. \right\} Q_1 \\
\frac{\partial}{\partial S_i} & = 0.866 & -0.749 & 0.749 & -0.396 & \checkmark \\
& 0.457 & 0.396 & -0.396 & 0.209 &
\end{array}$$

Next obtain the matrix of residual coefficients (R_2) by subtracting Q_1 from R_1 as recorded below:

$$0.251 \quad -0.415 \quad -0.396 \quad R_2 = R_1 - Q_1$$

$$-0.145 \quad 0.251 \quad 0.396$$

$$-0.396 \quad 0.396 \quad 0.791 \quad \checkmark$$

Reflect the variables X_2 and X_3 and obtain the reflected matrix of residual coefficients R_2'

$$0.251 \quad 0.145 \quad 0.396$$

$$0.145 \quad 0.251 \quad 0.396$$

$$0.396 \quad 0.396 \quad 0.791$$

Determine the sum of columns

$$S_i = \sum_{k=1}^3 r_{ik}$$

$$S_1 = 0.251 + 0.145 + 0.396 = 0.791$$

$$S_2 = 0.145 + 0.251 + 0.396 = 0.791$$

$$S_3 = 0.396 + 0.396 + 0.791 = 1.582 \checkmark$$

Determine the sum of all columns

$$T = \sum_{i=1}^3 S_i = 0.791 + 0.791 + 1.582 = 3.165 \checkmark$$

The 2nd centroid factor (B) is obtained by dividing each of the column total by sq. root of sum of all columns (T), that is.

$$UL_1(i) = \frac{S_i}{\sqrt{T}} = \frac{S_i}{\sqrt{3.165}}$$

| S_i | 0.791 | 0.791 | 1.582 | S_i / \sqrt{T} |
|------------------|-------|-------|-------|------------------|
| S_i / \sqrt{T} | 0.445 | 0.445 | 0.889 | \checkmark |

A B
 $x_1 \checkmark$ $x_3 \checkmark$
 $x_2 \checkmark$

The loadings of all the two factors are given below.

Centroid Factor Loadings

| Variables | Factor A | Factor B | Communality (h^2) |
|-----------|------------|------------|-----------------------------------|
| x_1 | 0.866^2 | 0.445^2 | $(0.866)^2 + (0.445)^2 = 0.947$ |
| x_2 | -0.866^2 | -0.445^2 | $(-0.866)^2 + (-0.445)^2 = 0.947$ |
| x_3 | 0.457^2 | -0.889^2 | $(0.457)^2 + (-0.889)^2 = 1.000$ |
| | 1.708 | 1.187 | 2.894 |

Proportions of total and common variance:

Centroid Factor Loadings

| Variables | Factor A | Factor B | Communality (h^2) | x_3 Factor B |
|----------------------------------|----------|----------|-----------------------|----------------|
| Eigen value (Variance accounted) | 1.708 | 1.187 | 2.894 | |
| Proportion of total variance | 0.427 | 0.297 | | |
| Proportion of common variance | 0.590 | 0.410 | | |

The common variance, i.e. $\sum_{i=1}^3 h^2 = 2.894$

Total variance is equal to no. of variables = 3.000 $x_1 x_2 x_3 \checkmark$

Conclusion:

Assign each variable to the factor with which it has the maximum absolute loadings.

Factor Name of factor Variable and description

1 Demographic x_1

x_2

2 Economical x_3

The proportion of total variance of factor-1 = 0.427 (42.7%)

" " " " " factor-2 = 0.297 (29.7%)

The " " common variance of factor-1 = 0.590 (59.0%)

" " " " " factor-2 = 0.410 (41.0%)

Example 5.3

The rating of expenditure, income, education, religion, marital, age and weights for persons in a town were collected in the range 1-10 (1-least rated, while 10 is highest rated) and recorded below.

| | | | | | | | | | | |
|-------|---|---|---|---|---|---|---|---|---|----|
| X_1 | 6 | 4 | 4 | 1 | 4 | 4 | 3 | 7 | 5 | 5 |
| X_2 | 8 | 4 | 1 | 2 | 3 | 4 | 3 | 7 | 3 | 4 |
| X_3 | 9 | 6 | 6 | 6 | 5 | 6 | 0 | 6 | 1 | 2 |
| X_4 | 9 | 8 | 5 | 3 | 5 | 8 | 9 | 9 | 8 | 83 |

If the variables are grouped into demographic, economic and social factors, determine the first, second and third centroid factors using the method of factor analysis.

Solution

The correlation matrix R for 6 variables is as given below.

$$r_{ik} = \frac{\sum_{j=1}^n X_{ij} X_{kj} - n \bar{X}_i \bar{X}_k}{\sqrt{\sum_{j=1}^n X_{ij}^2 - n \bar{X}_i^2} \sqrt{\sum_{j=1}^n X_{kj}^2 - n \bar{X}_k^2}}$$

$$R_1 = \begin{bmatrix} 1.000 & 0.742 & 0.168 & 0.496 \\ 0.742 & 1.000 & 0.424 & 0.568 \\ 0.168 & 0.424 & 1.000 & 0.050 \\ 0.496 & 0.568 & 0.050 & 1.000 \end{bmatrix}$$

Since the matrix R_1 is positive manifold, none of the variables is reflected and correlation matrix given as below

$$R_1 = \begin{bmatrix} 1.000 & 0.742 & 0.168 & 0.496 \\ 0.742 & 1.000 & 0.424 & 0.568 \\ 0.168 & 0.424 & 1.000 & 0.050 \\ 0.496 & 0.568 & 0.050 & 1.000 \end{bmatrix}$$

$$\gamma = \frac{(x_1 - \bar{x}_1)(x_2 - \bar{x}_2)}{\sqrt{(x_1 - \bar{x}_1)^2 + (x_2 - \bar{x}_2)^2}} \leftarrow \text{Indirect relation}$$

Determine \bar{e} sum of columns

$$S_i = \sum_{k=1}^4 r_{ik},$$

$$S_1 = 1.000 + 0.742 + 0.168 + 0.496 = 2.406$$

$$S_2 = 0.742 + 1.000 + 0.424 + 0.568 = 2.734$$

$$S_3 = 0.168 + 0.424 + 1.000 + 0.050 = 1.642$$

$$S_4 = 0.496 + 0.568 + 0.050 + 1.000 = 2.114$$

Determine \bar{e} sum of all columns

$$T = \sum_{k=1}^4 S_i = 2.406 + 2.734 + 1.642 + 2.114 = 8.895$$

The first centroid factor A, is obtained by dividing column total by sq. root of sum of all columns (T), that is,

$$UL(i) = \frac{S_i}{\sqrt{T}} = \frac{S_i}{\sqrt{8.895}}$$

$$S_i \quad 2.406 \quad 2.734 \quad 1.642 \quad 2.114$$

$$\frac{S_i}{\sqrt{T}} \quad 0.807 \quad 0.917 \quad 0.551 \quad 0.709$$

In order to obtain \bar{e} 2nd centroid factor B, develop \bar{e} matrix of factor cross product Q_1 .

| | 0.807 | 0.917 | 0.551 | 0.709 |
|-------|----------------------|----------------------|----------------------|----------------------|
| 0.807 | 0.807×0.807 | 0.807×0.917 | 0.807×0.551 | 0.807×0.709 |
| 0.917 | 0.917×0.807 | 0.917×0.917 | 0.917×0.551 | 0.917×0.709 |
| 0.551 | 0.551×0.807 | 0.551×0.917 | 0.551×0.551 | 0.551×0.709 |
| 0.709 | 0.709×0.807 | 0.709×0.917 | 0.709×0.551 | 0.709×0.709 |

Next obtain the matrix of residual coefficients (R_2) by subtracting Q_1 from R_1 as recorded below.

$$0.349 \quad 0.003 \quad -0.276 \quad -0.076$$

$$0.003 \quad 0.160 \quad -0.081 \quad -0.082$$

$$-0.276 \quad -0.081 \quad 0.697 \quad -0.340$$

$$-0.076 \quad -0.082 \quad -0.340 \quad 0.498$$

Reflect the variables X_3 and X_4 and obtain the reflected matrix of residual coefficient R'_2

$$0.349 \quad 0.003 \quad 0.276 \quad 0.076$$

$$0.003 \quad 0.160 \quad 0.081 \quad 0.082$$

$$0.276 \quad 0.081 \quad 0.697 \quad 0.340$$

$$0.076 \quad 0.082 \quad 0.340 \quad 0.498$$

Determine the sum of columns.

$$S_1 = 0.349 + 0.003 + 0.276 + 0.076 = 0.704$$

$$S_2 = 0.003 + 0.160 + 0.081 + 0.082 = 0.325$$

$$S_3 = 0.276 + 0.081 + 0.697 + 0.340 = 1.394$$

$$S_4 = 0.076 + 0.082 + 0.340 + 0.498 = 0.996$$

↓

Determine the sum of all columns

$$T = 0.704 + 0.325 + 1.394 + 0.996 = 3.419$$

The second centroid factor B, is obtained by dividing each of the column total by sq. root of sum of all columns (T), that is,

$$UL(i) = \frac{s_i}{\sqrt{T}} = \frac{s_i}{\sqrt{3.419}}$$

$$\begin{array}{cccc} s_i & 0.704 & 0.325 & 1.394 & 0.996 \\ \frac{s_i}{\sqrt{T}} & 0.381 & 0.176 & 0.754 & 0.538 \end{array}$$

In order to obtain the third centroid factor C, develop the matrix of factor cross product Q_2 .

$$\begin{array}{ccccc} \downarrow \rightarrow & 0.381 & 0.176 & -0.754 & -0.538 \\ \text{First} & 0.381 & 0.145 & 0.067 & -0.287 & -0.205 \\ \text{Centroid} & 0.176 & 0.067 & 0.031 & -0.133 & -0.095 \\ \text{Factors} & -0.754 & -0.287 & -0.133 & 0.568 & \begin{matrix} 0.406 \\ -0.046 \end{matrix} \\ & -0.538 & -0.205 & -0.095 & 0.406 & 0.290 \end{array}$$

Next obtain the matrix of residual coefficients (R_2) by subtracting Q_2 from R_2 ;

$$\begin{array}{cccc} 0.204 & -0.064 & 0.011 & 0.129 \\ -0.064 & 0.129 & 0.052 & 0.013 \\ 0.011 & 0.052 & 0.129 & -0.746 \\ 0.129 & 0.013 & -0.746 & 0.208 \end{array}$$

Reflect the variables X_2, X_3 and X_4 and obtained reflected matrix of residual coeff. R_2'

$$\begin{array}{cccc} 0.204 & 0.064 & 0.011 & 0.129 \\ 0.064 & 0.129 & 0.052 & 0.013 \\ 0.011 & 0.052 & 0.129 & 0.746 \\ 0.129 & 0.013 & 0.746 & 0.208 \end{array}$$

Determine the sum of columns

$$S_1 = 0.204 + 0.064 + 0.011 + 0.129 = 0.409$$

$$S_2 = 0.064 + 0.129 + 0.052 + 0.013 = 0.258$$

$$S_3 = 0.011 + 0.052 + 0.129 + 0.746 = 0.937$$

$$S_4 = 0.129 + 0.013 + 0.746 + 0.208 = 1.096$$

Determine the sum of all columns

$$T = 0.409 + 0.258 + 0.937 + 1.096 = 3.699$$

The third centroid factor C_3 is obtained by dividing column total by sq. root of T

$$UL(i) = \frac{s_i}{\sqrt{T}} = \frac{s_i}{\sqrt{2.699}}$$

$$T_1 \quad s_i \quad 0.409 \quad 0.258 \quad 0.937 \quad 1.096$$

$$T_2 \quad \frac{s_i}{\sqrt{T}} \quad 0.249 \quad 0.157 \quad 0.570 \quad 0.667$$

The loadings of all $\bar{e} 3$ factors are given in table below

Centroid Factor Loadings

| Variables | Factor A | Factor B | Factor C | Communality (h^2) |
|--------------|----------|----------|----------|---|
| X_1 | 0.807 | 0.381 | 0.249 | $(0.807)^2 + (0.381)^2 + (0.249)^2 = 0.796$ |
| X_2 | 0.917 | 0.176 | -0.157 | $(0.917)^2 + (0.176)^2 + (-0.157)^2 = 0.871$ |
| X_3 | 0.551 | -0.754 | -0.570 | $(0.551)^2 + (-0.754)^2 + (-0.570)^2 = 0.871$ |
| X_4 | 0.709 | -0.538 | 0.667 | $(0.709)^2 + (-0.538)^2 + (0.667)^2 = 0.792$ |
| These totals | 2.296 | 1.034 | 0.857 | 3.330 |

Proportion of total and common variance.

Centroid Factor Loadings

| Variables | Factor A | Factor B | Factor C | Communality (h^2) |
|----------------------------------|----------|----------|----------|-----------------------|
| Eigen value (Variance accounted) | 2.296 | 1.034 | 0.857 | 3.330 |
| Proportion of total variance | 0.574 | 0.259 | 0.214 | |
| Proportion of common variance | 0.689 | 0.311 | 0.257 | |

The common variance, i.e. $\sum_{i=1}^4 h^2 = 3.330$

Total variance is equal to no. of variables = 4.000

Conclusion

Assign each variable to the factor with which it has \bar{e} maximum absolute loadings.

Factor Name of factor Variable and description

1 Demographic X_1

X_2

X_4

2 Economical X_3

3 Social

The proportion of total variance of factor - 1 = 0.574 (57.4%)

" " " " " " " - 2 = 0.259 (25.9%)

" " " " " " " - 3 = 0.214 (21.4%)

The proportion of common variance of factor - 1 = 0.689 (68.9%)

" " " " " " " - 2 = 0.311 (31.1%)

" " " " " " " - 3 = 0.257 (25.7%)

6. ANALYSIS OF VARIANCE AND COVARIANCE.

Analysis of Variance (ANOVA)

- It is an analysis technique that is used to compare more than 2 population such as in comparing the yield of 4 varieties, the duration of fluorescent tubes from 5 vendors or engine performance from 4 modes. In this section we consider one-way analysis of variance.
- In designing the experiment the experimental units are arranged in surpentine manner then treatments allocated randomly to the experimental units. Label the experimental units from 1-16 treatments T_1, T_2, T_3 and T_4 for 4 treatments. In the case of 16 experimental units and 4 treatments, the number of experiments units for each treatment is $N/t = 16/4 = 4$. The arrangement of experimental units is as shown below for 16 experimental units.

| | | | |
|----|----|----|----|
| 01 | 02 | 03 | 04 |
| 08 | 07 | 06 | 05 |
| 09 | 10 | 11 | 12 |
| 16 | 15 | 14 | 13 |

In order to assign the treatments randomly, one can use the lottery method or random number method. In the lottery method, label the 16 marbles from 1 to 16, shuffle the marbles and select one marble randomly at a time from an urn. The first 4 marbles denote the first 4 experimental units assigned treatment T_1 , second 4 marbles denote the 4 experimental units to treatment T_2 , ... last 4 marbles denote the experimental units assigned to treatment T_4 .

An alternative method to the lottery method is the random number method. In this method select 16 3-digit numbers (if the size of total experimental units is two-digit, select three-digit to avoid repetition), rank the random numbers and assign the first $n=4$ ranks that denote the experimental units to treatment T_1 , second $n=4$ ranks to treatment T_2 ..., last $n=4$ ranks to treatment T_4 resp. The three-digit random numbers selected from table below are recorded in the next table.

G2

| | | | | | | |
|-------|-------|-------|-------|-------|-------|-------|
| 3 2 8 | 5 7 1 | 4 5 8 | 3 5 2 | 8 4 2 | 1 2 9 | 2 4 4 |
| 6 1 6 | 8 3 8 | 5 4 3 | 0 9 8 | 2 5 8 | 8 8 0 | 2 7 8 |
| 3 9 4 | 4 4 1 | 8 6 5 | 0 0 0 | 1 8 7 | 3 9 1 | 5 9 1 |
| 7 6 5 | 8 7 3 | 3 0 4 | 4 8 2 | 7 6 2 | 5 9 2 | 2 5 2 |
| 8 8 3 | 9 5 0 | 9 2 4 | 9 5 4 | 3 7 2 | 1 8 6 | 6 5 6 |
| 6 6 4 | 0 9 1 | 0 0 1 | 1 3 8 | 4 0 9 | 4 5 2 | 4 6 0 |
| 0 5 5 | 5 1 9 | 4 7 1 | 2 8 2 | 8 6 6 | 5 7 1 | 2 6 7 |
| 4 0 2 | 6 0 4 | 6 3 5 | 1 0 7 | 2 3 6 | 7 9 9 | 5 6 0 |
| 9 5 1 | 4 4 6 | 0 1 5 | 5 5 7 | 6 4 7 | 1 4 8 | 5 6 9 |
| 8 6 2 | 0 2 9 | 8 1 8 | 0 0 9 | 2 5 9 | 5 7 7 | 5 2 9 |
| 5 3 3 | 8 8 1 | 8 7 8 | 5 6 3 | 7 8 1 | 6 4 0 | 4 7 2 |
| 2 6 6 | 6 0 7 | 3 7 9 | 3 3 9 | 1 4 3 | 8 7 0 | 1 6 0 |
| 4 0 0 | 4 2 3 | 1 8 2 | 7 9 2 | 2 0 9 | 7 5 3 | 2 4 5 |
| 8 0 7 | 4 3 1 | 6 5 4 | 3 8 0 | 1 0 9 | 4 7 4 | 6 5 4 |
| 9 6 6 | 9 6 1 | 6 7 9 | 2 5 5 | 5 1 0 | 2 8 2 | 3 1 0 |
| 5 9 5 | 7 9 8 | 4 8 6 | 5 1 2 | 0 2 5 | 8 8 7 | 5 4 1 |

G3

| Treatment | Random Numbers | Rank | Treatment | Random Numbers | Rank |
|----------------|----------------------------------|----------------------|----------------|----------------------------------|----------------------|
| T ₁ | 5 7 1 8 3 8 4 4 1 8 7 8 | 08 12 05 13 | T ₃ | 9 5 0 0 9 1 5 1 9 6 0 4 | 15 02 07 09 |
| T ₂ | 4 4 6 0 2 9 8 8 1 6 0 7 | 06 01 14 10 | T ₄ | 4 2 3 4 3 1 9 6 1 7 9 8 | 03 04 16 11 |

The ranks 08, 12, 05, 13 that denote the experimental units are assigned treatment T₁.

| | | | | | | |
|------------------|---|---|---|---|---|----------------|
| " 06, 01, 14, 10 | " | " | " | " | " | T ₂ |
| " 15, 02, 07, 09 | " | " | " | " | " | T ₃ |
| " 03, 04, 16, 11 | " | " | " | " | " | T ₄ |

The allocation is as shown below:

| | | | |
|--------------------|--------------------|--------------------|--------------------|
| 01, T ₃ | 02, T ₂ | 03, T ₄ | 04, T ₄ |
| 08, T ₁ | 07, T ₂ | 06, T ₃ | 05, T ₁ |
| 09, T ₂ | 10, T ₃ | 11, T ₄ | 12, T ₁ |
| 16, T ₄ | 15, T ₂ | 14, T ₃ | 13, T ₁ |

In investigating the differences in groups in one-way analysis of variance we follow the steps discussed below:

i) Arrangement of the observations

| Group | Observations | Total |
|----------------|---|----------------|
| T ₁ | y ₁₁ y ₁₂ ... y _{1n} | T ₁ |
| T ₂ | y ₂₁ y ₂₂ ... y _{2n} | T ₂ |
| ⋮ | ⋮ | ⋮ |
| T _t | y _{t1} y _{t2} ... y _{tn} | T _t |
| | | G |

(ii) Analysis of Variance Model.

$$Y_{ij} = \mu + t_i + e_{ij} \quad \begin{cases} i=1, 2, \dots, t \\ j=1, 2, \dots, n \end{cases}$$

Where: Y_{ij} - The j^{th} observation for the i^{th} group, μ - mean, t_i - i^{th} group effect.
 e_{ij} - The random error that is normally and independently distributed.

(iii) Hypothesis

$$H_0: \mu_1 = \mu_2 = \dots = \mu_t \text{ against } H_a: \mu_i \neq \mu_k \text{ for atleast 1 value (i, k) at } \alpha 1\text{-o.s}$$

(iv) Analysis of variance components

Consider corrected total sum of squares.

$$\begin{aligned} \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y})^2 &= \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i + \bar{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2 + 2 \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i)(\bar{y}_i - \bar{y}) + \sum_{i=1}^t \sum_{j=1}^n (\bar{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2 + \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2 \text{ since } 2 \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i)(\bar{y}_i - \bar{y}) = 0 \\ &= \frac{1}{n} (T_1^2 + T_2^2 + \dots + T_t^2) - \frac{G^2}{n} + \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2 \end{aligned}$$

$$SS_{\text{Total}} = SS_{\text{Treatments}} + SS_{\text{Error}}$$

$$\text{Where, } SS_{\text{Total}} = \sum_{i=1}^t \sum_{j=1}^n y_{ij}^2 - \frac{G^2}{n} \quad SS_{\text{Error}} = \sum_{i=1}^t \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2$$

$$SS_{\text{Treatments}} = \frac{1}{n} (T_1^2 + T_2^2 + \dots + T_t^2) - \frac{G^2}{n}$$

(v) Analysis of variance table.

| Source of variation | df | SS | MS | F |
|---------------------|-----|--------------------------|--------------------------------------|--------------------------|
| Treatments | t-1 | SS _{Treatments} | $\frac{SS_{\text{Treatments}}}{t-1}$ | MS _{Treatments} |
| Error | N-t | By subtraction | $\frac{SS_{\text{Error}}}{N-t}$ | MS _{Error} |
| Total | N-1 | SS _{Total} | | |

(vi) Test statistic

$$F = \frac{MS_{\text{Treatments}}}{MS_{\text{Error}}} \text{ that has approximately F-distribution with } t-1 \text{ and } N-t \text{ degrees of freedom}$$

(vii) Decision.

Reject H_0 if $F > F_\alpha(t-1, N-t)$ and conclude that the means of treatment are significantly different at $\alpha 1\text{-o.s.}$

Turkey's Test

In many practical situations, we wish to compare pairs of means and determine which means differ by testing the difference between all \bar{e} pairs of treatment means. Thus we are interested in contrasts of \bar{e} form $\Gamma = \mu_i - \mu_j$ for all $i \neq j$. The popular method for making such comparison is the Turkey's procedure proposed in 1953 for testing hypotheses for which the overall significance level is exactly α when the sample sizes are equal and at most α when the sample sizes are unequal. The procedure can also be used to construct intervals on the differences in all pairs of means. The Turkey's procedure makes use of \bar{e} distribution of the standard range statistic

$$q = \frac{\bar{y}_{\max} - \bar{y}_{\min}}{\sqrt{\text{MSE}_{\text{Error}} / n}}$$

Where \bar{y}_{\max} and \bar{y}_{\min} are the largest and smallest sample means, resp, out of q group of p sample means. Turkey's test declares two means significantly different if the absolute value of their sample differences exceeds:

$$T_q = \begin{cases} q_{\alpha}(t, f) \sqrt{\frac{\text{MSE}_{\text{Error}}}{n}} & \text{for equal sample sizes} \\ \frac{q_{\alpha}(t, f)}{\sqrt{2}} \sqrt{\text{MSE}_{\text{Error}} \left(\frac{1}{n_i} + \frac{1}{n_j} \right)} & \text{for unequal sample sizes.} \end{cases}$$

Where, t = The number of treatments

f = The number of d.f associated with the $\text{MSE}_{\text{Error}}$.

n = The number of observations for each treatment when observations are equal

(n_i, n_j) = The sample sizes for unequal sample size pair (i, j)

Ex 6.1

The yields of rice for different varieties T_1, T_2, T_3 and T_4 were collected and recorded below

| Varieties | Yield | | | | Totals. | \bar{y} |
|-----------|-------|----|----|----|---------|-----------|
| T_1 | 39 | 25 | 16 | 20 | 100 | 25 |
| T_2 | 33 | 45 | 52 | 30 | 160 | 40 |
| T_3 | 45 | 45 | 40 | 50 | 180 | 45 |
| T_4 | 41 | 45 | 60 | 54 | 200 | 50 |

Determine whether the varieties have different effects at $\alpha = 0.05$ l.o.s

Solution

(i) Analysis of variance model

$$y_{ij} = \mu + t_i + e_{ij} \quad \begin{cases} i = 1, 2, \dots, t \\ j = 1, 2, \dots, n \end{cases}$$

(ii) Hypothesis

$$H_0: \mu_1 = \mu_2 = \dots = \mu_t \text{ against } H_a: \mu_i \neq \mu_k \text{ for atleast one pair } (i, k) \text{ at } \alpha = 0.05$$

(iii) Analysis of variance components.

$$SS_{\text{Total}} = \sum_{i=1}^t \sum_{j=1}^n y_{ij}^2 - \frac{G^2}{N}$$

$$SS_{\text{Treatments}} = \frac{1}{n} (T_1^2 + T_2^2 + \dots + T_t^2) - \frac{G^2}{N} = \frac{1}{4} (100^2 + \dots + 200^2) - \frac{(640)^2}{16} = 2292.00$$

$$SS_{\text{Error}} = SS_{\text{Total}} - SS_{\text{Treatments}} = 2292 - 1400 = 892.00$$

(iv) Analysis of variance table.

| Source of variation | df | SS | MSS | F |
|---------------------|-----|----|------|--------------------------------|
| Treatments | t-1 | 3 | 1400 | $\frac{466.67}{466.67} = 6.28$ |
| Error | N-t | 12 | 892 | $\frac{74.33}{74.33} = 6.28$ |
| Total | N-1 | 15 | 2292 | |

(v) Test statistic

$$F = \frac{MS_{\text{Treatment}}}{MS_{\text{Error}}} = \frac{466.67}{74.33} = 6.28$$

(vi) Decision

Tabulated value $F_{0.05}(3, 12) = 3.49$, since $F_{\text{calc}} > F_{\text{tab}}$ reject H_0 and conclude that the means of varieties are significantly different at $\alpha = 0.05$.

- Since the null hypothesis has been rejected we wish to determine the pairs of means that are significantly different using Turkey's Test. That is determine the statistic

$$T_q = q_{\alpha}(t, f) \sqrt{\frac{MSE}{n}} = 4.20 \sqrt{\frac{74.33}{4}} = 18.11$$

| | \bar{y}_1 | \bar{y}_2 | \bar{y}_3 |
|-------------|-------------|-------------|-------------|
| | 25 | 40 | 45 |
| \bar{y}_2 | 40 | -15 | 0 |
| \bar{y}_3 | 45 | -20 | -5 |
| \bar{y}_4 | 50 | -25 | -10 |

The pairs that are significantly different include (\bar{y}_1, \bar{y}_3) and (\bar{y}_1, \bar{y}_4)

Confidence Interval

$$\bar{y}_i - \bar{y}_j - q_{\alpha}(t, f) \sqrt{\frac{MSE}{n}} \leq \mu_i - \mu_j \leq \bar{y}_i - \bar{y}_j + q_{\alpha}(t, f) \sqrt{\frac{MSE}{n}}$$

$$= -18.11 \leq \mu_i - \mu_j \leq 18.11$$

Exq 6.2

The lifespan of fluorescent tubes from 4 vendors V_1, V_2, V_3 and V_4 were collected & recorded below

| Vendors | Duration | | | | Totals |
|---------|----------|----|----|----|--------|
| | 15 | 20 | 25 | 12 | |
| V_1 | 15 | 20 | 25 | 12 | 72 |
| V_2 | 18 | 16 | 20 | 10 | 64 |
| V_3 | 45 | 45 | 40 | 50 | 180 |
| V_4 | 15 | 15 | 15 | 39 | 84 |

Determine whether there are differences of lifespan of fluorescent tubes due to the vendors at $\alpha=0.05$ l.o.s

Solution

(i) Analysis of variance model.

$$y_{ij} = \mu + t_i + e_{ij} \quad \begin{cases} i=1, 2, \dots, t \\ j=1, 2, \dots, n \end{cases}$$

(ii) Hypothesis

$$H_0: \mu_1 = \mu_2 = \dots = \mu_t \text{ against } H_a: \mu_i \neq \mu_k \text{ for atleast one pair } (i, k) \text{ at } \alpha \text{ l.o.s}$$

(iii) Analysis of variance components

$$SS_{\text{Total}} = \sum_{i=1}^t \sum_{j=1}^n y_{ij}^2 - \frac{\bar{y}^2}{n} = 15^2 + \dots + 39^2 - \frac{400^2}{16} = 2820$$

$$SS_{\text{Treatments}} = \frac{1}{n} (T_1^2 + T_2^2 + \dots + T_t^2) - \frac{\bar{y}^2}{n} = \frac{1}{4} (72^2 + \dots + 84^2) - \frac{400^2}{16} = 2184$$

$$SS_{\text{Error}} = SS_{\text{Total}} - SS_{\text{Treatments}} = 2820 - 2184 = 636$$

(iv) ANOVA table

| Source of variation | df | ss | MS | F |
|---------------------|-----|------|-----|-------|
| Treatments | t-1 | 2184 | 728 | 13.74 |
| Error | N-t | 636 | 53 | |
| Total | N-1 | 2820 | | |

(v) Test statistic; $F = 13.74$

(vi) Decision

Tabulated value $F_{0.05}(3, 12) = 3.49$, since $F_{\text{calc}} > F_{\text{tab}}$ reject H_0 and conclude that the means of vendors are significantly different at $\alpha=0.05$ l.o.s

- Since H_0 has been rejected we wish to determine the pairs of means that are significantly different using Turkey's test. Determine \bar{e} statistic.

$$T_q = q_{\alpha}(t, f) \sqrt{\frac{MSE}{n}} = 4.20 \sqrt{\frac{53}{4}} = 15.29$$

| | \bar{y}_1 | \bar{y}_2 | \bar{y}_3 |
|-------------|-------------|-------------|-------------|
| | 18 | 16 | 45 |
| \bar{y}_2 | 16 | 2 | 29 |
| \bar{y}_3 | 45 | -27 | -29 |
| \bar{y}_4 | 21 | -3 | -5 |

The pairs that are significantly different include: (\bar{y}_1, \bar{y}_3) , (\bar{y}_2, \bar{y}_3) and (\bar{y}_3, \bar{y}_4) .

Confidence Interval

$$\bar{y}_i - \bar{y}_j - q_{\alpha}(t, f) \sqrt{\frac{MSE}{n}} \leq \mu_i - \mu_j \leq \bar{y}_i - \bar{y}_j + q_{\alpha}(t, f) \sqrt{\frac{MSE}{n}}$$

$$= 15.29 \leq \mu_i - \mu_j \leq 15.29$$

Analysis of Covariance.

- In " " " one may be interested in studying one independent variable, X and may wish to control the influence of uncontrollable variable (covariate), Y , which is known to be correlated with the dependent variable Y . Then the technique of analysis of covariance is considered for this type of analysis for a valid evaluation of the outcome of analysis. The arrangement of data is given below.

| Treatments | T_1 | T_2 | \dots | T_t | |
|--------------|-----------------------|-----------------------|----------|-----------------------|-------------------|
| | $X \quad Y$ | $X \quad Y$ | \dots | $X \quad Y$ | |
| Observations | $X_{11} \quad Y_{11}$ | $X_{21} \quad Y_{21}$ | | $X_{t1} \quad Y_{t1}$ | |
| | $X_{12} \quad Y_{12}$ | $X_{22} \quad Y_{22}$ | | $X_{t2} \quad Y_{t2}$ | |
| | $\vdots \quad \vdots$ | $\vdots \quad \vdots$ | | $\vdots \quad \vdots$ | |
| | $X_{1n} \quad Y_{1n}$ | $X_{2n} \quad Y_{2n}$ | | $X_{tn} \quad Y_{tn}$ | |
| Total | T_{x1} | T_{y1} | T_{x2} | T_{y2} | T_{xt} T_{yt} |

Analysis of covariance model

$$Y_{ij} = \mu + t_i + b(X_{ij} - \bar{X}) + e_{ij}$$

Analysis of covariance components

$$G_{xx} = \sum_{i=1}^t \sum_{j=1}^n X_{ij}^2 - \frac{G_x^2}{N}, \quad G_{xy} = \sum_{i=1}^t \sum_{j=1}^n X_{ij} Y_{ij} - \frac{G_x G_y}{N}, \quad G_{yy} = \sum_{i=1}^t \sum_{j=1}^n Y_{ij}^2 - \frac{G_y^2}{N}$$

$$T_{xx} = \frac{1}{n} \sum_{i=1}^t T_{xi}^2 - \frac{G_x^2}{N}, \quad T_{xy} = \frac{1}{n} \sum_{i=1}^t T_{xi} T_{yi} - \frac{G_x G_y}{N}, \quad T_{yy} = \frac{1}{n} \sum_{i=1}^t T_{yi}^2 - \frac{G_y^2}{N}$$

$$E_{xx} = G_{xx} - T_{xx}, \quad E_{xy} = G_{xy} - T_{xy}, \quad E_{yy} = G_{yy} - T_{yy}$$

Adjusted values of X are given as

$$G'_{xx} = G_{xx} - \frac{G_{xy}^2}{G_{xx}}, \quad E'_{xx} = E_{xx} - \frac{E_{xy}^2}{E_{xx}}$$

$$G'_{xy} = G_{xy} - \frac{G_{xy}^2}{G_{xx}}, \quad E'_{xy} = E_{xy} - \frac{E_{xy}^2}{E_{xx}}, \quad T'_{xx} = G'_{xx} - E'_{xx}$$

Adjusted values of Y are given as

$$G'_{yy} = G_{yy} - \frac{G_{xy}^2}{G_{yy}}, \quad E'_{yy} = E_{yy} - \frac{E_{xy}^2}{E_{yy}}, \quad T'_{yy} = G'_{yy} - E'_{yy}$$

Adjusted degrees of freedom.

The degrees of freedom for the error is adjusted for adjusted values of X and Y by reducing it by one and therefore d.f for adjusted values of X & Y for total are also reduced by one.

- If X is not affected by the treatments then there should not be significant difference b/w treatments with reference to X . Then the regression coefficient with treatments is computed as $b = \frac{E_{xy}}{E_{yy}}$ and therefore adjust treatment means for Y as:

$$\bar{X}_i = \bar{x}_i - b(\bar{y}_i - \bar{Y})$$

Conventionally the results are combined and given in a single table as shown below.

| Source of Variation | Sum of products | | | | Adjusted values of Y | | | |
|---|-----------------|----------|----------------|----------------|------------------------|----------|-------|-------------------|
| | df | XX | $X\bar{Y}$ | $\bar{Y}Y$ | df | SS | MS | F |
| Treatments | $t-1$ | T_{xx} | $T_{x\bar{y}}$ | $T_{\bar{y}y}$ | | | | |
| Error | $n-t$ | E_{xx} | $E_{x\bar{y}}$ | $E_{\bar{y}y}$ | $n-t-1$ | E_{xx} | EMS | |
| Total | $n-1$ | G_{xx} | $G_{x\bar{y}}$ | $G_{\bar{y}y}$ | $n-2$ | G_{xx} | | |
| Treatments adjusted for the average regression within treatments. | | | | | $t-1$ | T_{xx} | TMS | $\frac{TMS}{EMS}$ |

The regression coefficient within treatments is computed as $b = \frac{E_{xy}}{E_{xx}}$ if X is affected by treatments and therefore adjust the treatment means of Y as

$$\bar{Y}'_i = \bar{Y}_i - b(\bar{x}_i - \bar{x})$$

Exq. 6.3

The data of the variables X, Y and treatments $T_1, T_2 \& T_3$ were collected & recorded as shown

| Treatments | x | y | $x+y$ | $x-y$ | \bar{x} | \bar{y} | \bar{X} | \bar{Y} | | |
|------------|-----|-----|-------|-------|-----------|-----------|-----------|-----------|-------|------|
| T_1 | 10 | 2 | 16 | 5 | 14 | 7 | 40 | 14 | 13.33 | 4.67 |
| T_2 | 15 | 8 | 14 | 2 | 15 | 5 | 44 | 15 | 14.67 | 5.0 |
| T_3 | 20 | 5 | 25 | 6 | 21 | 9 | 66 | 20 | 22 | 6.67 |

Determine whether the treatments are significant at $\alpha = 0.05$ to 0.01 .

Solution.

(i) Analysis of covariance model.

$$Y_{ij} = \mu + t_i + b(x_{ij} - \bar{x}) + e_{ij}$$

(ii) Hypothesis

$$H_0: \mu_1 = \mu_2 = \mu_3 \text{ against } H_a: \mu_i \neq \mu_k \text{ for atleast one pair } (i, k) \text{ at } \alpha = 0.05 \text{ to } 0.01$$

(iii) Analysis of variance components for X .

$$G_{xx} = \sum_{i=1}^t \sum_{j=1}^n X_{ij}^2 - \frac{G_x^2}{N} = 10^2 + \dots + 21^2 - \frac{150^2}{9} = 164$$

$$G_{x\bar{y}} = \sum_{i=1}^t \sum_{j=1}^n X_{ij}\bar{y}_{ij} - \frac{G_x G_y}{N} = 10(2) + \dots + 21(9) - \frac{150(49)}{9} = 43.33$$

$$G_{\bar{y}y} = \sum_{i=1}^t \sum_{j=1}^n \bar{y}_{ij}^2 - \frac{G_y^2}{N} = 2^2 + \dots + 9^2 - \frac{49^2}{9} = 46.22$$

$$T_{xx} = \frac{1}{n} \sum_{i=1}^t T_{xi}^2 - \frac{G_x^2}{N} = \frac{1}{3} (40^2 + \dots + 66^2) - \frac{150^2}{9} = 130.67$$

$$T_{x\bar{y}} = \frac{1}{n} \sum_{i=1}^t T_{xi} \bar{y}_{i\bar{y}} - \frac{G_x G_y}{N} = \frac{1}{3} (40(14) + \dots + 66(20)) - \frac{150(49)}{9} = 30$$

$$T_{\bar{y}y} = \frac{1}{n} \sum_{i=1}^t T_{\bar{y}y}^2 - \frac{G_y^2}{N} = \frac{1}{3} (14^2 + \dots + 20^2) - \frac{49^2}{9} = 6.89$$

$$E_{xx} = G_{xx} - T_{xx} = 164 - 130.67 = 33.33$$

$$E_{xy} = G_{xy} - T_{xy} = 43.33 - 30 = 13.33$$

$$E_{yy} = G_{yy} - T_{yy} = 46.22 - 6.89 = 39.33$$

(iv) Adjusted values of X are given as

$$G'_{xx} = G_{xx} - \frac{G_{xy}^2}{G_{yy}} = 164 - \frac{43.33^2}{46.22} = 123.38$$

$$E'_{xx} = E_{xx} - \frac{E_{xy}^2}{E_{yy}} = 33.33 - \frac{13.33^2}{39.33} = 28.81$$

$$T'_{xx} = G'_{xx} - E'_{xx} = 123.38 - 28.81 = 94.56$$

Adjusted

| Source of variation | df | SS_x | SS_y | SS_{xy} |
|---------------------|-------|--------|--------|-----------|
| Treatments | $t-1$ | 2 | 130.67 | 6.89 |
| Error | $n-t$ | 6 | 33.33 | 39.33 |
| Total | $n-1$ | 8 | 164.00 | 46.22 |

| | df | SS_x | MS_x | F_x |
|--|---------|--------|--------|-------|
| | 2 | 94.56 | 47.28 | 8.20 |
| | $n-t-1$ | 5 | 28.81 | 5.76 |
| | $n-2$ | 7 | 173.38 | |

v) Test statistic

$$F = 8.20$$

vi) Decision (reject H_0)

Tabulated value $F_{0.05}(2, 5) = 5.79$, since $F_{\text{calc}} > F_{\text{tab}}$ then differences in treatments means are significant at $\alpha = 0.05$ l.o.s and therefore there are variations due to factor X .

vii) Adjusted means of X groups

$$\bar{X}'_i = \bar{X}_i - b(\bar{Y}_i - \bar{Y})$$

$$b = \frac{E_{xy}}{E_{yy}} = \frac{13.33}{39.33} = 0.3389$$

$$\bar{X}'_1 = \bar{X}_1 - b(\bar{Y}_1 - \bar{Y}) = 13.33 - 0.3389(4.67 - 5.44) = 13.59$$

$$\bar{X}'_2 = \bar{X}_2 - b(\bar{Y}_2 - \bar{Y}) = 14.67 - 0.3389(5.00 - 5.44) = 14.82$$

$$\bar{X}'_3 = \bar{X}_3 - b(\bar{Y}_3 - \bar{Y}) = 22.00 - 0.3389(6.67 - 5.44) = 21.58$$

viii) Adjusted components of Y are given as

$$G'_{yy} = G_{yy} - \frac{G_{xy}^2}{G_{xx}} = 46.22 - \frac{43.33^2}{164} = 34.77$$

$$E'_{yy} = E_{yy} - \frac{E_{xy}^2}{E_{xx}} = 39.33 - \frac{13.33^2}{33.33} = 34.00$$

$$T'_{yy} = G'_{yy} - E'_{yy} = 34.77 - 34.00 = 0.77$$

Source of variation df SS_x SS_y SS_{xy}

Treatments 2 130.7 6.89 30.00

Error 6 33.3 39.33 13.33

Total 8 164.0 46.22 43.33

Treatments adjusted for the avg. regression within treatments.

df SS_y MS_y F_y

2

E'_{yy} 34.00 6.80

5 34.00 6.80

7 34.77

T'_{yy} 0.77

2 0.77 0.39 0.057 = $\frac{0.39}{6.80}$

ix) Decision

Tabulated value $F_{0.05}(2,5) = 5.79$, since $F_{\text{calc}} > F_{\text{tab}}$ then the differences in treatments means are significant at $\alpha = 0.05$.

* Exercise.

Unit 7. CLUSTER AND CONJOINT ANALYSIS.

Cluster Analysis.

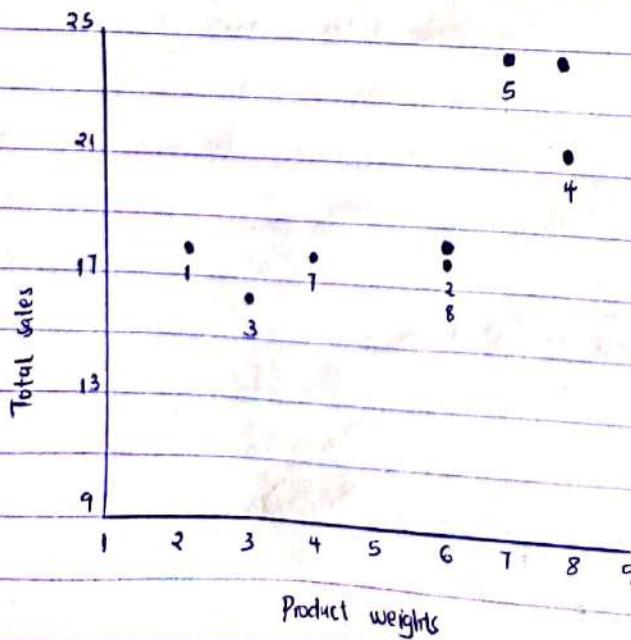
The " " technique groups persons, items or objects into unknown number of groups such that members of each group have similar attributes or simply groups that are homogeneous. The different cluster analysis techniques form groups such that the similarity among the members of the groups is maximised.

The cluster analysis techniques are similar to discriminant analysis except that whereas in discriminant analysis the number of groups are known, in cluster analysis the groups are unknown. In order to explain the concept of cluster analysis in illustration is conducted using the data of weights of products (X_1) and total sales (X_2) recorded in table below for vendors in a duration of one week.

Sales of product for vendors

| Id | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|------------------|----|----|----|----|----|----|----|----|
| Weight (X_1) | 2 | 6 | 3 | 8 | 7 | 8 | 4 | 6 |
| Sales (X_2) | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 |

Graph of product sales against product weights



The observations in the table are represented in a scatter diagram and the graph reveals the possible 3 groups that may be given as group one with individuals 1, 3, 7, group two with individual 2, 8 and group 3 with individuals 4, 5, 6 resp. The measures of similarity include distance measure, correlation coefficients and association coefficients.

Distance Measure.

- The similarity measure commonly used is the Euclidean distance measure that is given by the formula for instance for k points is given as $d_{ij} = \sqrt{\sum_{k=1}^n (x_{ik} - x_{jk})^2}$
Where d_{ij} is the distance btwn \bar{e} points i & j
 x_{ik} is the coordinate of \bar{e} point i along the axis k
 x_{ik} is the " " " " " " " " " " " " " " "
 n is the total number of variables (axes)

- Considering 2 points $(1, 2)$, the distance btwn the 2 points is given as:

$$d_{12} = \sqrt{(x_{11} - x_{21})^2 + (x_{12} - x_{22})^2}$$

Exq. 7.1.

The data of product weights and sales were collected and recorded in table below.

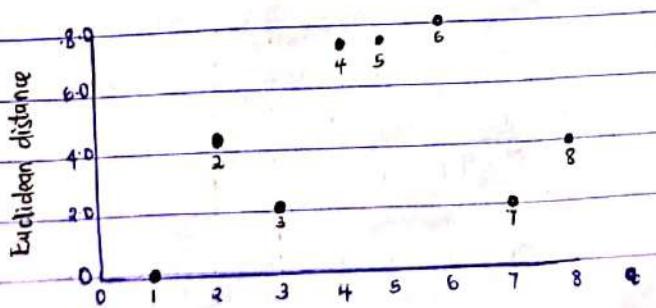
| | | | | | | | | |
|-------------------|----|----|----|----|----|----|----|----|
| Id | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| Weights (X_1) | 2 | 6 | 3 | 8 | 7 | 8 | 4 | 6 |
| Sales (X_2) | 18 | 18 | 18 | 18 | 18 | 18 | 18 | 18 |

Determine the possible cluster groups for the data using the Euclidean distance measure.

Solution.

| Id | X_1 | X_2 | Euclidean, d_{ij} |
|----|-------|-------|---------------------|
| 1 | 2 | 18 | 0.0 |
| 2 | 6 | 20 | 4.5 |
| 3 | 3 | 16 | 2.2 |
| 4 | 8 | 22 | 7.2 |
| 5 | 7 | 24 | 1.8 |
| 6 | 8 | 24 | 8.5 |
| 7 | 4 | 18 | 2.0 |
| 8 | 6 | 19 | 4.1 |

Graph of Euclidean distance against Id



Groups when using Euclidean distance are given as

Group 1: 1, 3, 7 Group 2: 2, 8 Group 3: 4, 5, 6

The representative point for all the objects within a cluster may be the centroid of that cluster.

The formula for the coordinates of the centroid of a cluster is given as;

$$Y_k = \frac{1}{n} \sum_{i=1}^n X_{ik}$$

Where Y_k is the coordinate of the centroid of the cluster along the axis k , X_{ik} the coordinate of the point i along the axis k , and n is the total no. of points in the cluster.

- If the distance btwn a point | partial cluster and another point | partial cluster is the least among all such values, then they can be grouped together.

| Id | X_1 | X_2 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----|-------|-------|-----|-----|-----|-----|-----|-----|-----|-----|
| 1 | 2 | 18 | 0.0 | 4.5 | 2.2 | 7.2 | 7.8 | 8.5 | 2.0 | 4.1 |
| 2 | 6 | 20 | 4.5 | 0.0 | 5.0 | 2.8 | 4.1 | 4.5 | 2.8 | 1.0 |
| 3 | 3 | 16 | 2.2 | 5.0 | 0.0 | 7.8 | 8.9 | 9.4 | 2.2 | 4.2 |
| 4 | 8 | 22 | 7.2 | 2.8 | 7.8 | 0.0 | 2.2 | 2.0 | 5.7 | 3.6 |
| 5 | 7 | 24 | 7.8 | 4.1 | 8.9 | 2.2 | 0.0 | 1.0 | 6.7 | 5.1 |
| 6 | 8 | 24 | 8.5 | 4.5 | 9.4 | 2.0 | 1.0 | 0.0 | 7.2 | 5.4 |
| 7 | 4 | 18 | 2.0 | 2.8 | 2.2 | 5.7 | 6.7 | 7.2 | 0.0 | 2.2 |
| 8 | 6 | 19 | 4.1 | 1.0 | 4.2 | 3.6 | 5.1 | 5.4 | 2.2 | 0.0 |

Correlation Coefficient

The correlation coefficient btwn variables can be used as a measure in cluster analysis, however, the measure is not commonly used in many fields. The measure when used in grouping variables considers variables with maximum coefficient as the variables that belong to same group.

Association Coefficient.

The measure that considers the presence or absence of attributes in objects such that if a particular attribute is present in an object, then the respective association coefficient is assumed as 1 otherwise it is assumed as 0. The association coefficient is a kind of similarity coefficient and grouping when done using this measure groups objects on the maximum value of such coefficient.

| Id | X_1 | X_2 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|----|-------|-------|-----|-----|-----|-----|-----|-----|-----|-----|
| 1 | 2 | 18 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 2 | 6 | 20 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 3 | 3 | 16 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| 4 | 8 | 22 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 5 | 7 | 24 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |
| 6 | 8 | 24 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |
| 7 | 4 | 18 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 8 | 6 | 19 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 |

Conjoint Analysis

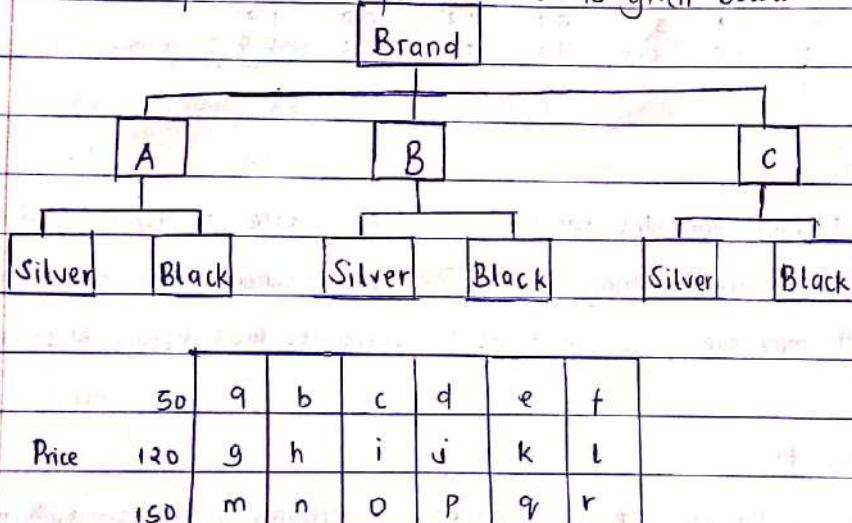
- The " " is a technique normally used to measure the customer preferences on products. Let the number of factors for the product be denoted by m and number of levels for each factor be denoted by n_1, n_2, \dots, n_m , then the total no. of profiles for the product is $n_1 \times n_2 \times \dots \times n_m$ (Kothari, 2004). A research project manager wish to buy computers for a project but before awarding the contract to supplier he/she wish to carry out a study on preference of project participants. The preference is affected by the brand, color, and price and the levels of each of the factors of the product are recorded as below:

Brand : Acer (A), HP (B) and Toshiba (C)

Color : Silver and Black

Price : 50, 120 and 150

The complete factorial combinations of the product profiles in for the factors (brand, color, and price) and respective levels is given below-



The description of the profiles are given as:

Profile a: Brand - A, Color - Silver, Price - 50 , Profile b: Brand - A, Color - Black, Price - 50
!

Profile g: Brand - C, Color - Silver , Price - 150 , Profile r: Brand - C, Color - Black , Price - 150

The total no. of product profile is $3 \times 2 \times 3 = 18$, i.e., all the 18 product profiles are considered for the study. A card defining the product profile, for instance one shown in Figure 5.4 is represented to the customer and he/she is required to give preference value using an interval or ordinal scale. In the interval scale, the respondent response is in the scale 0 to 10 where 0 means that the customer has absolutely no interest in the product with the profile shown on the respective card and 10 means the customer has extremely high interest to buy that product.

If one choose to use an ordinal scale, then each respondent is expected to rank all the product profiles shown to him/her from the most desirable product profile (rank 1) to the least desirable product profile (rank 18) such that rank 1 means that the customer has extremely high interest to buy the product with the profile shown on the respective card and the rank 18 means that the customer has absolutely no interest to buy that product.

7.4 Product profile card for the combination a (Brand - A, Color - silver, Price - \$ 50)

| Card 1 | |
|--------------------|----------------------|
| Brand: | Acer |
| Color: | Silver |
| Price: | \$ 50 |
| Respondent Rating: | <input type="text"/> |

The traditional conjoint analysis techniques is thought of as a multiple regression problem whose respondents' ratings for the product concepts are observations on the dependent variable. The xtics of the product or attribute levels are observations on the independent or predictor variables. The estimated regression coefficients associated with the independent variables are the part-worth utilities or preference scores for the levels and coefficient of determination R^2 for the regression xtics is internal consistency of the respondent.

- Consider the following coding for the levels of factors given as:

A-1, B-2, C-3, Silver-1, Black-2

- The factorial experiment design is then given in table below with respective codes.

| Id | Raw | | | Preference | Coding | | |
|----|-------|--------|-------|------------|--------|-------|-------|
| | Brand | Color | Price | | Brand | Color | Price |
| 1 | A | Silver | 50 | 5.0 | 1 | 1 | 50 |
| 2 | A | Silver | 120 | 6.0 | 1 | 1 | 120 |
| 3 | A | Silver | 150 | 7.0 | 1 | 1 | 150 |
| 4 | A | Black | 50 | 9.5 | 1 | 2 | 50 |
| 5 | A | Black | 120 | 8.5 | 1 | 2 | 120 |
| 6 | A | Black | 150 | 9.0 | 1 | 2 | 150 |
| 7 | B | Silver | 50 | 4.0 | 2 | 1 | 50 |
| 8 | B | Silver | 120 | 5.5 | 2 | 1 | 120 |
| 9 | B | Silver | 150 | 6.5 | 2 | 1 | 150 |
| 10 | B | Black | 50 | 8.0 | 2 | 2 | 50 |
| 11 | B | Black | 120 | 7.5 | 2 | 2 | 120 |
| 12 | B | Black | 150 | 9.5 | 2 | 2 | 150 |
| 13 | C | Silver | 50 | 3.0 | 3 | 1 | 50 |
| 14 | C | Silver | 120 | 4.5 | 3 | 1 | 120 |
| 15 | C | Silver | 150 | 5.0 | 3 | 1 | 150 |
| 16 | C | Black | 50 | 6.5 | 3 | 2 | 50 |
| 17 | C | Black | 120 | 7.0 | 3 | 2 | 120 |
| 18 | C | Black | 150 | 9.0 | 3 | 2 | 150 |

Now record the respondent's preference using an interval scale in which the respondent rates the profile in a scale 0 to 10. Assume the first card is made up of 1st levels of brand, color and price, that is Acer, silver and \$50 and the respondent gives preference of 4, then 2nd card has profile brand, color and price of Acer, silver and \$120, and respondent has preference rating of 5 and so on, the preferences are recorded resp. for each of the profiles.

- Consider using the dummy coding procedure for the independent variables or product xtics. The procedure considers the coding in which the dummy coding uses 0 to represent absence of xtic and 1 to represent the presence of the xtic.
- The brand factor will be coded in 3 separate columns (Acer, HP and Toshiba), color would be recorded in 2 column (silver & black) and price in 3 columns (\$50, \$120 & \$150).

Dummy Coding for the factor levels.

| Acer | HP | Toshiba | Silver | Black | \$50 | \$120 | \$150 | Preference |
|------|----|---------|--------|-------|------|-------|-------|------------|
| 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 4 |
| 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 8 |
| 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 2 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 6 |
| 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 5 |
| 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 3 |
| 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 9 |
| 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 6 |
| 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 5 |
| 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 8 |
| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 7 |
| 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 5 |
| 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 8 |
| 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 7 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 6 |
| 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 99 |

Consider fitting multiple linear regression

$$Y_i = \beta_0 + \beta_1 X_{1i} + \dots + \beta_{p-1} X_{(p-1)i} + \epsilon_i$$

- In multiple regression analysis for the conjoint analysis there is no independent variable may be perfectly predicted based on the state of any other independent variable or combination of independent variables and therefore regression procedure cannot separate the effects of the independent variables, for instance, we can perfectly predict the state of Brand A based on the states of Brand B and C, this is called linear dependency.

- In order to solve the problem of linear dependency, omit one column from each attribute (any of the attributes, it does not matter which attribute is to be omitted) and therefore omit the first attribute for each of the variables (factors) and produce a table as shown in table 7.4

In order to estimate for the regression parameters by least squares or Maximum likelihood method may be this section of least squares method is used to solve for the parameters using the following steps:

Rewrite the model as

$$Y = X\beta + e$$

Table 7.4 Recorded data for attributes after omitting one attribute on each factor.

| HP | Toshiba | Black | 120 | 150 | Preference |
|----|---------|-------|-----|-----|------------|
| 0 | 0 | 0 | 0 | 0 | 4 |
| 0 | 0 | 0 | 1 | 1 | 5 |
| 0 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 8 |
| 0 | 0 | 1 | 1 | 0 | 5 |
| 0 | 0 | 1 | 0 | 1 | 2 |
| 1 | 0 | 0 | 0 | 0 | 6 |
| 1 | 0 | 0 | 1 | 0 | 5 |
| 1 | 0 | 0 | 0 | 1 | 3 |
| 1 | 0 | 1 | 0 | 0 | 9 |
| 1 | 0 | 1 | 1 | 0 | 6 |
| 1 | 0 | 1 | 0 | 1 | 5 |
| 0 | 1 | 0 | 0 | 0 | 8 |
| 0 | 1 | 0 | 1 | 0 | 7 |
| 0 | 1 | 0 | 0 | 1 | 5 |
| 0 | 1 | 1 | 0 | 0 | 8 |
| 0 | 1 | 1 | 1 | 0 | 7 |
| 0 | 1 | 1 | 0 | 1 | 6 |

Multiply both sides by X' such that

$$X'Y = X'X\beta$$

Multiply both sides by inverse of $X'X$ such that

$$\beta = (X'X)^{-1}X'Y$$

| | | | | | | | | | | | | | |
|---------|----|---|---|----|-----------------|--------|---------|--------|--------|--------|--------|--------|--------|
| | 18 | 6 | 6 | -9 | 6 | 6 | | 0.333 | -0.167 | -0.167 | -0.111 | -0.167 | -0.167 |
| | 6 | 6 | 0 | 3 | 2 | 2 | | -0.167 | 0.333 | 0.167 | 0.000 | 0.000 | 0.000 |
| $X'X =$ | 6 | 0 | 6 | 3 | 2 | 2 | $X'X =$ | -0.167 | 0.167 | 0.333 | 0.000 | 0.000 | 0.000 |
| | 9 | 3 | 3 | 9 | 3 | 3 | | -0.111 | 0.000 | 0.000 | 0.222 | 0.000 | 0.000 |
| | 6 | 2 | 3 | 3 | 6 | 0 | | -0.167 | 0.000 | 0.000 | 0.000 | 0.333 | 0.167 |
| | 6 | 2 | 2 | 3 | 0 | 6 | | -0.167 | 0.000 | 0.000 | 0.000 | 0.167 | 0.333 |
| | 99 | | | | 4.944 | | | | | | | | |
| | 34 | | | | 1.667 | | | | | | | | |
| $X'Y =$ | 41 | | | | $\hat{\beta} =$ | 2.833 | | | | | | | |
| | 56 | | | | | 1.444 | | | | | | | |
| | 35 | | | | | -1.333 | | | | | | | |
| | 21 | | | | | 3.667 | | | | | | | |

$$\hat{Y} = 4.944 + 1.667X_1 + 2.833X_2 + 1.444X_3 - 1.333X_4 + 3.667X_5$$

Hypothesis

$H_0: \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$ against $H_a: \beta_i \neq \beta_k$ for at least one pair (i, k) at $\alpha = 0.05$ & 0.05
From the table 7.4 above, include
 $(\bar{y} - y)^2 : 0.892 | 1.929 | 1.633 | 2.596 | 0.003 | 0.522 | 0.313 | 0.071 | 0.003 | 0.892 | 0.522 | 0.313 | 0.049 | 0.309 | 0.790$
 $1.494 | 0.790 | 0.198 | = 13.44$

Analysis of Variance table.

| Source of variation | df | SS | MS | F |
|---------------------|----|-------|--------|--------|
| Regression | 5 | 75.06 | 15.011 | 13.398 |
| Error | 12 | 13.44 | 1.120 | |
| Total | 17 | 88.50 | | |

Test statistic

$$F = 13.398$$

Decision

Tabulated value of $F_{0.05}(5, 12) = 3.106$. Since $F_{\text{calc}} > F_{\text{tab}}$, reject H_0 and conclude $\beta_i \neq \beta_k$ for atleast one pair (i, k) at $\alpha = 0.05$ & 0.05

Performance of model.

Coefficient of determination, R^2

$$R^2 = \frac{SS_{\text{Regression}}}{SS_{\text{Total}}} = \frac{75.06}{88.50} = 0.8481 \text{ (the regression accounts for } 84.81\% \text{ of variation)}$$

Adjusted coefficient of determination, R_{adj}^2

$$R_{\text{adj}}^2 = 1 - \frac{\frac{SS_{\text{Error}}}{df_{\text{Error}}}}{\frac{SS_{\text{Total}}}{df_{\text{Total}}}} = 1 - \frac{\frac{13.44}{12}}{\frac{88.50}{17}} = 0.7848 \text{ (the regression accounts for } 74.48\% \text{ of variation)}$$

In this problem we have considered solving separate regression eqtns for each respondent, therefore to estimate utilities, the respondent must have to evaluate at least as many profiles as parameters to be evaluated.

5/7/2021

Distributions of Linear functions of a Random Vector or Quadratic Forms

Distribution of Linear functions of q Random Vector

Properties of Linear Function of a random vector.

a) If a single random variable such as X_1 is multiplied by a constant c , then,

$$E(cx_1) = cE(x_1) = c\mu_1$$

$$\begin{aligned} V(cx_1) &= E(cx_1 - c\mu_1)^2 = c^2 E(x_1 - \mu_1)^2 \\ &= c^2 V(x_1) = c^2 \delta_{11} \end{aligned}$$

$$E(cx_1) = cE(x_1) = c\mu_1$$

$$V(cx_1) = E(cx_1 - c\mu_1)^2 = c^2 E(x_1 - \mu_1)^2$$

$$= c^2 V(x_1) = c^2 \delta_{11}$$

$$c^2 V(x_1) = c^2 \delta_{11}$$

b) If X_2 is a second random variable where a & b are constants then using additional properties of expectation we get

$$\text{Cov}(ax_1, bx_2) = E(ax_1 - a\mu_1)(bx_2 - b\mu_2) = aE(x_1 - \mu_1)(x_2 - \mu_2)b$$

$$= a \text{Cov}(x_1, x_2) b = a \delta_{12} b$$

c) If $Z = ax_1 + bx_2$ is a linear combination,

$$\begin{aligned} E(Z) &= E(ax_1 + bx_2) = aE(x_1) + bE(x_2) = a\mu_1 + b\mu_2 \end{aligned}$$

$$\begin{aligned} V(Z) &= V(ax_1 + bx_2) = E(ax_1 + bx_2 - a\mu_1 - b\mu_2)^2 \\ &= E(ax_1 - a\mu_1 + bx_2 - b\mu_2)^2 \\ &= a^2 E(x_1 - \mu_1)^2 + ab E(x_1 - \mu_1)(x_2 - \mu_2) + b^2 E(x_2 - \mu_2)^2 \\ &= a^2 \delta_{11} + 2ab \delta_{12} + b^2 \delta_{22} \end{aligned}$$

d) If $C' = [a \ b]$ and variance covariance matrix of $\mathbf{X} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ is $\Sigma = \begin{bmatrix} \delta_{11} & \delta_{12} \\ \delta_{21} & \delta_{22} \end{bmatrix}$ then we have

$$E(C'X) = C'E \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = [a \ b] \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = a\mu_1 + b\mu_2$$

$$\begin{aligned} V(C'X) &= C'V \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} C = [a \ b] \begin{bmatrix} \delta_{11} & \delta_{12} \\ \delta_{21} & \delta_{22} \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \\ &= C^2 \delta_{11} + 2ab \delta_{12} + b^2 \delta_{22} \end{aligned}$$

e) Consider a general linear combination of p random variables x_1, x_2, \dots, x_p such that

$$Z_1 = c_{11}x_1 + c_{12}x_2 + \dots + c_{1p}x_p$$

$$Z_2 = c_{21}x_1 + c_{22}x_2 + \dots + c_{2p}x_p$$

$$Z_q = c_{q1}x_1 + c_{q2}x_2 + \dots + c_{qp}x_p$$

OR

$$Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_q \end{bmatrix} = \begin{bmatrix} c_{11} & c_{12} & \dots & c_{1p} \\ c_{21} & c_{22} & \dots & c_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ c_{q1} & c_{q2} & \dots & c_{qp} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} = CX$$

$$E(z) = E[CX] = CE(x) = C E \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} = CMx$$

$$V(z) = V(CX) = CV(x)C' = CV \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix} C' = C \begin{bmatrix} \delta_{11} & \delta_{12} & \dots & \delta_{1p} \\ \delta_{21} & \delta_{22} & \dots & \delta_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \delta_{p1} & \delta_{p2} & \dots & \delta_{pp} \end{bmatrix} C'$$

where M_x and Σ_x are the mean vector and variance covariance matrix of x .

Examples:

1. Let $x' = [x_1 \ x_2 \ x_3 \ x_4]$ be a random vector with mean vector $\mu_x' = [\mu_1 \ \mu_2 \ \mu_3 \ \mu_4]$ and variance covariance matrix Σ_x .

$$\Sigma_x = \begin{bmatrix} \delta_{11} & \delta_{12} & \delta_{13} & \delta_{14} \\ \delta_{21} & \delta_{22} & \delta_{23} & \delta_{24} \\ \delta_{31} & \delta_{32} & \delta_{33} & \delta_{34} \\ \delta_{41} & \delta_{42} & \delta_{43} & \delta_{44} \end{bmatrix}$$

- i) Find the mean vector and variance covariance matrix for the linear combination given as

$$z = x_1 - x_3$$

Solution:

$$E(z) = E(cx) = CE(x) = [1 \ 0 \ -1 \ 0] E \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} = [1 \ 0 \ -1 \ 0] \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{bmatrix} = \mu_1 - \mu_3$$

$$= \mu_1 - \mu_3$$

$$V(z) = V(cx) = CV(x)C'$$

$$= [1 \ 0 \ -1 \ 0] \begin{bmatrix} \delta_{11} & \delta_{12} & \delta_{13} & \delta_{14} \\ \delta_{21} & \delta_{22} & \delta_{23} & \delta_{24} \\ \delta_{31} & \delta_{32} & \delta_{33} & \delta_{34} \\ \delta_{41} & \delta_{42} & \delta_{43} & \delta_{44} \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \end{bmatrix}$$

$$= [1 \ 0 \ -1 \ 0] \begin{bmatrix} \delta_{11} - \delta_{13} \\ \delta_{21} - \delta_{23} \\ \delta_{31} - \delta_{33} \\ \delta_{41} - \delta_{43} \end{bmatrix}$$

$$= \delta_{11} - \delta_{13} - \delta_{31} + \delta_{33}$$

ii) Find the mean vector and variance covariance matrix for the linear combination given as:

$$Z_1 = X_1 - X_3, \quad Z_2 = X_2 + X_4$$

$$\underline{Z} = \begin{bmatrix} Z_1 \\ Z_2 \end{bmatrix}$$

$$\begin{aligned} E(\underline{Z}) &= C E(\underline{X}) = \begin{bmatrix} 1 & 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{bmatrix} \\ &= \begin{bmatrix} \mu_1 - \mu_3 \\ \mu_2 + \mu_4 \end{bmatrix} \end{aligned}$$

$$V(\underline{Z}) = V(C\underline{X}) = CV(\underline{X})C'$$

$$= \begin{bmatrix} 1 & 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} \delta_{11} & \delta_{12} & \delta_{13} & \delta_{14} \\ \delta_{21} & \delta_{22} & \delta_{23} & \delta_{24} \\ \delta_{31} & \delta_{32} & \delta_{33} & \delta_{34} \\ \delta_{41} & \delta_{42} & \delta_{43} & \delta_{44} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ -1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & -1 & 0 \end{bmatrix} \begin{bmatrix} \delta_{11} - \delta_{13} \\ \delta_{21} - \delta_{23} \\ \delta_{31} - \delta_{33} \\ \delta_{41} - \delta_{43} \end{bmatrix} \begin{bmatrix} \delta_{12} + \delta_{14} \\ \delta_{22} + \delta_{24} \\ \delta_{32} + \delta_{34} \\ \delta_{42} + \delta_{44} \end{bmatrix}$$

$$= \begin{bmatrix} \delta_{11} - \delta_{13} - \delta_{31} + \delta_{33} & \delta_{12} + \delta_{14} - \delta_{32} - \delta_{34} \\ \delta_{21} - \delta_{23} - \delta_{41} + \delta_{43} & \delta_{22} + \delta_{24} + \delta_{42} + \delta_{44} \end{bmatrix}$$

2. Derive expressions of the mean and variance covariance for the following linear combinations in terms of the mean and variance of the random variables X_1, X_2, X_3, X_4

$$a) X_1 - X_2 + X_3$$

$$b) X_1 + X_2 - X_3 + X_4$$

If X_1, X_3 and X_4 are independent.

Solution

$$a) E(\underline{Z}) = E(C\underline{X}) = \begin{bmatrix} 1 & -1 & 1 & 0 \end{bmatrix} E \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \begin{bmatrix} 1 & -1 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{bmatrix}$$

$$= \mu_1 - \mu_2 + \mu_3$$

$$V(\underline{Z}) = V(C\underline{X}) = CV(\underline{X})C' = \begin{bmatrix} 1 & -1 & 1 & 0 \end{bmatrix} V(\underline{X}) \begin{bmatrix} 1 \\ -1 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -1 & 1 & 0 \end{bmatrix} \begin{bmatrix} \delta_{11} & \delta_{12} & \delta_{13} & \delta_{14} \\ \delta_{21} & \delta_{22} & \delta_{23} & \delta_{24} \\ \delta_{31} & \delta_{32} & \delta_{33} & \delta_{34} \\ \delta_{41} & \delta_{42} & \delta_{43} & \delta_{44} \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ 1 \\ 0 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & -1 & 1 & 0 \end{bmatrix} \begin{bmatrix} \delta_{11} - \delta_{12} + \delta_{13} \\ \delta_{21} - \delta_{22} + \delta_{23} \\ \delta_{31} - \delta_{32} + \delta_{33} \\ \delta_{41} - \delta_{42} + \delta_{43} \end{bmatrix}$$

$$= \delta_{11} - \delta_{12} + \delta_{13} - \delta_{21} + \delta_{22} - \delta_{23} + \delta_{31} - \delta_{32} + \delta_{33}$$

$$= \delta_{11} - 2\delta_{12} + 2\delta_{13} + \delta_{23} - 2\delta_{23} + \delta_{33}$$

$$\begin{aligned}
 b) E(z) &= E(cx) = cE(x) = cE\begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = c\begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{bmatrix} \\
 &= [1 \ 1 \ -1 \ 1] \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \end{bmatrix} \\
 &= \mu_1 + \mu_2 - \mu_3 + \mu_4
 \end{aligned}$$

Q

$$V(z) = V(cx) = cv(x)c^T = [1 \ 1 \ -1 \ 1] \begin{bmatrix} \delta_{11} & \delta_{12} & \delta_{13} & \delta_{14} \\ \delta_{21} & \delta_{22} & \delta_{23} & \delta_{24} \\ \delta_{31} & \delta_{32} & \delta_{33} & \delta_{34} \\ \delta_{41} & \delta_{42} & \delta_{43} & \delta_{44} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ -1 \\ 1 \end{bmatrix}$$

$$= [1 \ 1 \ -1 \ 1] \begin{bmatrix} \delta_{11} + \delta_{12} - \delta_{13} + \delta_{14} \\ \vdots \\ \delta_{21} + \delta_{22} - \delta_{23} + \delta_{24} \\ \delta_{31} + \delta_{32} - \delta_{33} + \delta_{34} \\ \delta_{41} + \delta_{42} - \delta_{43} + \delta_{44} \end{bmatrix}$$

$$\begin{aligned}
 &= \delta_{11} + \delta_{12} - \delta_{13} + \delta_{14} + \delta_{21} + \delta_{22} - \delta_{23} + \delta_{24} + \delta_{31} + \delta_{32} - \delta_{33} + \delta_{34} + \delta_{41} + \delta_{42} - \delta_{43} + \delta_{44} \\
 &= \delta_{11} + 2\delta_{12} - 2\delta_{13} + 2\delta_{14} + \delta_{22} - 2\delta_{23} + \delta_{33} - 2\delta_{34} + 2\delta_{24} + \delta_{44} \\
 &= \delta_{11} + 2\delta_{12} + \delta_{22} - 2\delta_{23} + \delta_{33} + 2\delta_{24} + \delta_{44}
 \end{aligned}$$

Exercise: Derive the expression for the mean & variance-covariance for the following linear combinations in terms of the mean and variance of the random variables X_1, X_2, X_3 and X_4 .

$X_1 - X_2 + X_3 + X_4$. If X_2, X_3 and X_4 are independent.

The equation of the best line of fit is an approximation of the simple linear regression where β_0 is the y-intercept and A is the slope gradient of the best line representation.

Example

The age of the household head (x) and expenditure (y) for households were collected and recorded below.

| | | | | | | | | |
|---|----|----|----|----|----|----|----|----|
| x | 25 | 28 | 30 | 25 | 30 | 50 | 60 | 72 |
| y | 10 | 12 | 15 | 10 | 20 | 15 | 20 | 18 |

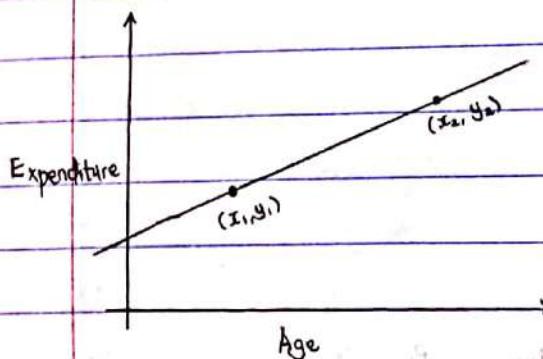
Determine the simple linear regression for estimating the expenditure of the household.

Solution

The simple linear regression model takes the general form

$$Y_i = \beta_0 + \beta_1 X_i + \theta_i$$

Construct a scatter diagram and draw the best line of fit. Estimate the unknown parameters.



$$\hat{\beta}_0 = y\text{-intercept} = 9.27$$

$$\hat{\beta}_1 = \frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{20.70 - 9.27}{30 - 0} = 0.143$$

Multivariate Normal Linear Regression and Correlation Analysis.

Simple Linear Regression:

- Simple linear regression model is one that has single response & single predictor variable that takes the general form

$$Y_i = \beta_0 + \beta_1 X_i + \theta_i, i = 1, 2, \dots, n$$

(Response) = [mean depending on X_i] + [error]

- In order to predict the response variable, one has to estimate first the unknown parameters β_0 & β_1 . The methods commonly used to estimate the unknown parameters include ; graphical, least squares and maximum likelihood methods as described below.

a) Graphical Method.

Construct scatter diagram for the response & predictor variables : The best line of fit is drawn such that no. of points above the best line and no. of points below the best line are nearly equal. The equation of the best line of fit is an approximation of the simple linear regression where the β_0 is the y-intercept & β_1 is the slope / gradient of the best line representation.

Example:

The example & solution is above.

b) Least Square Method.

Consider the simple linear regression

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

Make e_i the subject, square both sides.

$$e_i^2 = (y_i - \beta_0 - \beta_1 x_i)^2$$

Sum over all possible sample observations-

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\text{Let } \Phi = \sum_{i=1}^n e_i^2$$

$$\Phi = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

Differentiate Φ w.r.t β_0 and evaluate to zero to solve for β_0 .

$$\frac{\partial \Phi}{\partial \beta_0} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\frac{\partial \beta_0}{\partial \beta_0}$$

$$\hat{\beta}_0 = \sum_{i=1}^n (y_i - \beta_1 x_i) = \bar{y} - \beta_1 \bar{x}$$

Differentiate Φ w.r.t β_1 and evaluate to 0 to solve for β_1 .

$$\frac{\partial \Phi}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\frac{\partial \beta_1}{\partial \beta_1}$$

But, $\beta_0 = \bar{y} - \beta_1 \bar{x}$ so that,

$$(y_i - \bar{y} + \beta_1 \bar{x} - \beta_1 x_i) x_i = 0$$

$$\beta_1 (x_i - \bar{x}) x_i = \sum_{i=1}^n (y_i - \bar{y}) x_i$$

$$= \sum_{i=1}^n (y_i - \bar{y}) x_i$$

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

$$= \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

Example:

The age and weight of persons were collected and recorded below.

| | X | 25 | 28 | 30 | 25 | 30 | 50 | 60 | 72 |
|--|---|----|----|----|----|----|----|----|----|
| | Y | 10 | 12 | 15 | 10 | 20 | 15 | 20 | 18 |

Use the simple linear regression to estimate weight.

Solution:

| x | y | xy | x^2 |
|-----|--------------|----|-----|
| 25 | | | |
| 28 | | | |
| 30 | | | |
| 25 | | | |
| 30 | | | |
| 50 | | | |
| 60 | | | |
| 72 | | | |
| 320 | 120 532 1518 | | |

a) Maximum Likelihood Method.

Consider the probability function.

$$f(x_i, \beta, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp -\frac{1}{2\sigma^2} (y_i - \beta_0 - \beta_1 x_i)^2$$

Consider the joint probability function for x_1, x_2, \dots, x_n

$$f(x, \beta, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp -\frac{1}{2\sigma^2} (y_1 - \beta_0 - \beta_1 x_1)^2$$

$$\frac{1}{\sqrt{2\pi}\sigma} \exp -\frac{1}{2\sigma^2} (y_2 - \beta_0 - \beta_1 x_2)^2$$

⋮

$$\frac{1}{\sqrt{2\pi}\sigma} \exp -\frac{1}{2\sigma^2} (y_n - \beta_0 - \beta_1 x_n)^2$$

$$L = (2\pi)^{-n/2} \sigma^{-n/2} \exp -\frac{1}{2\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

Find natural way of both sides

$$\ln L = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

Differentiate w.r.t. β_0 and equate to zero.

$$\frac{\partial}{\partial \beta_0} \ln L = -2 \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

Differentiate w.r.t. β_1 and equate to zero.

$$\sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i) x_i = 0$$

$$\hat{\beta}_1 = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

Differentiate w.r.t. σ^2

$$\frac{\partial}{\partial \sigma^2} \ln L = \frac{1}{2\sigma} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

$$n \sigma^2 = \sum (y_i - \beta_0 - \beta_1 x_i)^2 \quad \hat{\sigma}^2 = \frac{1}{n} \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

Under Least Squares.

Example:

| | |
|---|--|
| x | |
| y | |

We refer the simple linear regression to estimate weight.

Solution:

| x | y | xy | x^2 |
|----|----|------|-------|
| 25 | 10 | 250 | 625 |
| 28 | 12 | 336 | 784 |
| 30 | 15 | 450 | 900 |
| 25 | 10 | 250 | 625 |
| 30 | 20 | 600 | 900 |
| 50 | 15 | 750 | 2500 |
| 60 | 20 | 1200 | 3600 |
| 72 | 18 | 1296 | 5184 |
| | | 5132 | 15118 |

$$\hat{\beta}_1 = \frac{\sum xy - n \bar{x} \bar{y}}{\sum x^2 - n \bar{x}^2} = \frac{5132 - 8 \left(\frac{820}{8} \right) \left(\frac{140}{8} \right)}{15118 - 8 \left(\frac{820}{8} \right)^2} = 0.1432$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = 9.271$$

Test Significance of simple Linear Regression Models

Hypothesis: $H_0: \beta_1 = 0$ against $H_a: \beta_1 \neq 0$ at α level of significance.

Analysis of Variance Components

Consider the corrected total sum of squares.

$$\begin{aligned}\sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y})^2 + 2 \sum_{i=1}^n (y_i - \hat{y})(\hat{y}_i - \bar{y}) + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (y_i - \hat{y})^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ \text{since } 2 \sum_{i=1}^n (y_i - \hat{y})(\hat{y}_i - \bar{y}) &= 0\end{aligned}$$

$$SS_{\text{Total}} = SS_{\text{Error}} + SS_{\text{Reg}}$$

$$= SS_{\text{Reg}} + SS_{\text{Error}}$$

Degrees of freedom:

The total sum of squares has $n-1$, regression has $p-1$ and error has $n-p$. Where p is the value of regression parameters, n is the sample size.

| Source of variation | df | SS | MS | F |
|---------------------|-------|----------------------------|-------------------------------|-------------------------------|
| Regression | $p-1$ | By subtraction | $\frac{SS_{\text{Reg}}}{p-1}$ | $\frac{MS_{\text{Reg}}}{MSE}$ |
| Error | $n-p$ | $\sum (y_i - \hat{y}_i)^2$ | $\frac{SSE}{n-p}$ | |
| Total | $n-1$ | $\sum (y_i - \bar{y})^2$ | | |

Test statistic

$F = \frac{MS_{\text{Reg}}}{MS_{\text{Error}}}$ that has approximately F distribution with $p-1, n-p$ d.f.

Decision:

Reject H_0 if $F > F_{\alpha}(p-1, n-p)$ and conclude that $\beta_1 \neq 0$ at α level of significance.

Coefficient of determination (R^2)

$$R^2 = 1 - \frac{SS_{\text{Error}}}{SS_{\text{Total}}} = \frac{SS_{\text{Reg}}}{SS_{\text{Total}}}$$

Adjusted Coefficient of determination:

$$R_{\text{adj}}^2 = 1 - \frac{SS_{\text{Error}}}{n-p}$$

$$\frac{SS_{\text{Total}}}{n-1}$$

Test significance of rejection parameter using t -test.

Suppose we wish to test the hypothesis. Hypothesis: $H_0: \beta_1 = \beta_{10}$ against $H_a: \beta_1 \neq \beta_{10}$ at α level of significance.

Test statistic:

$$t_1 = \frac{\hat{\beta}_1 - \beta_{10}}{\sqrt{V(\hat{\beta}_1)}} \text{ but } V(\hat{\beta}_1) = \sigma^2 \frac{1}{S_{xx}} \text{ where } S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$$

t is approximately t distribution with $n-2$ degrees of freedom.

Decision:

Reject H_0 if $|t_1| > t_{\alpha/2}, n-p$ at α level of significance.

Similarly, we wish to test the hypothesis,

Hypothesis: $H_0: \beta_0 = \beta_{00}$ against $H_a: \beta_0 \neq \beta_{00}$ at $\alpha = 0.05$

Test statistic

$$t_0 = \frac{\hat{\beta}_0 - \beta_{00}}{\sqrt{V(\hat{\beta}_0)}} \text{ but } V(\hat{\beta}_0) = \sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right) \text{ where } S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$$

t_0 is approximately t distribution with $n-2$ degrees of freedom.

Decision:

Reject H_0 if $|t_0| > t_{\alpha/2}, n-p$ at α -level of significance.

Confidence Intervals

The $100(1-\alpha)\%$ C.I. for $\bar{\beta}$ parameters are given as

$$\hat{\beta}_0 \pm t_{\alpha/2}, n-p \sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} \text{ such that,}$$

$$\hat{\beta}_0 - t_{\alpha/2}, n-p \sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)} < \beta_0 < \hat{\beta}_0 + t_{\alpha/2}, n-p \sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}$$

$$\hat{\beta}_1 \pm t_{\alpha/2}, n-p \sqrt{\frac{\sigma^2}{S_{xx}}} \text{ such that}$$

$$\hat{\beta}_1 - t_{\alpha/2}, n-p \sqrt{\frac{\sigma^2}{S_{xx}}} < \beta_1 < \hat{\beta}_1 + t_{\alpha/2}, n-p \sqrt{\frac{\sigma^2}{S_{xx}}}$$

Example

The amount of sand and floor area in urban area construction were collected and recorded below.

| | | | | | | | | |
|----------|----|----|----|----|----|----|----|----|
| Sand (x) | 25 | 28 | 30 | 25 | 30 | 50 | 60 | 72 |
| Area (y) | 10 | 12 | 15 | 10 | 20 | 15 | 20 | 18 |

Determine the following for determining the floor area:

- The simple linear regression model using least squares

- b) Whether the CLR model is significant.
 c) Whether the SLR parameters are significant.
 d) 95% C.I. of the regression parameters.

Solution.

Consider the simple linear regression model.

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

| x | y | xy | x^2 | y^2 | \bar{y} |
|----|-----|------|-------|-------|-----------|
| 25 | 10 | 250 | 625 | 100 | |
| 28 | 12 | 336 | 784 | 144 | |
| 30 | 15 | 450 | 900 | 225 | |
| 25 | 10 | 250 | 625 | 100 | |
| 30 | 20 | 600 | 900 | 400 | |
| 50 | 15 | 750 | 2500 | 225 | |
| 60 | 20 | 1200 | 3600 | 400 | |
| 72 | 18 | 1296 | 5184 | 324 | |
| 80 | 120 | 5132 | 15118 | 1918 | |

$$\hat{y}_i = \beta_0 + \beta_1 x_i \text{ where,}$$

$$\hat{\beta}_1 = \frac{\sum xy - n\bar{x}\bar{y}}{\sum x^2 - n\bar{x}^2} = \frac{5132 - 8(120/8)(120/8)}{15118 - 8(120/8)^2} = 0.1432.$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$= \frac{120}{8} - 0.1432 \left(\frac{120}{8}\right)$$

$$= 9.271$$

Simple linear regression model is given as:

$$\hat{y}_i = 9.271 + 0.1432x$$

b) Hypothesis: $H_0: \beta_1 = 0$ against $H_a: \beta_1 \neq 0$.

Analysis of variance component.

Corrected total sum of squares.

$$SS_{Total} = \sum (y_i - \bar{y})^2 = \sum y_i - n\bar{y}^2$$

$$= 10^2 + 12^2 + \dots + 18^2 - 8 \left(\frac{120}{8}\right)^2 = 1918 - 8 \left(\frac{120}{8}\right)^2$$

$$\begin{aligned}
 SSe_{\text{error}} &= \sum (y_i - \hat{y}_i)^2 \\
 &= 2.85^2 + 1.28^2 + \dots + 0.68^2 = 70.45 \\
 SSe_{\text{reg}} &= SSt_{\text{total}} - SSe_{\text{error}} \\
 &= 118 - 70.45 = 47.55
 \end{aligned}$$

Anova table.

| Source of variation | D.F | S.S | M.S | F |
|------------------------------|-----|-------|-------|--------|
| Regression | 1 | 47.55 | 47.55 | 4.05 |
| Error | 6 | 70.45 | 11.74 | |
| Total | 7 | 118 | | |
| T.S = F = $\frac{MS_R}{MSE}$ | | | | = 4.05 |

Decision

Tabulated value $f_{0.05}(1, 6) = 5.99$, since $f < f_{0.05}(1, 6)$ fail to reject H_0 & conclude that $\beta_1 = 0$ at $\alpha = 0.05$ l.o.s

Significance of the Simple Regression Parameters:

Hypothesis : $H_0: \beta_0 = \beta_{00}$ vs $H_a: \beta_0 \neq \beta_{00}$ at $\alpha = 0.05$ l.o.s

$$t_0 = \frac{\beta_0 - \beta_{00}}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} = \frac{9.271}{\sqrt{11.74 \left(\frac{1}{8} + \frac{320}{15118 - 8 \cdot \frac{320}{18}} \right)}} = 2.997$$

Decision

Tabulated $t_{0.025}(6) = t_{0.025}(6) = 2.45$, since $t_0 > t_{0.025}(6)$ reject H_0 and conclude that $\beta_0 \neq 0$ at $\alpha = 0.05$ l.o.s.

Hypothesis : $H_0: \beta_1 = \beta_{10}$ vs $H_a: \beta_1 \neq \beta_{10}$ at $\alpha = 0.05$ l.o.s

$$t_1 = \frac{\beta_1 - \beta_{10}}{\sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}} = \frac{0.1432}{\sqrt{11.74 \left(\frac{1}{8} + \frac{320}{15118 - 8 \cdot \frac{320}{18}} \right)}} = 2.012$$

Decision: $t_{0.025}, n-p = t_{0.025}(6) = 2.45$

Since $t_1 < t_{0.025}(6)$ fail to reject H_0 and conclude that $\beta_1 = 0$ at $\alpha = 0.05$ l.o.s.

d) Confidence Intervals:

The $100(1-\alpha)\%$ confidence intervals are given as:

$$\hat{\beta}_0 = 9.271, 9.271 \pm t_{0.025}(6) \sqrt{\sigma^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}}\right)}$$

$$= 9.271 \pm 2.45 \sqrt{11.74 \left(\frac{1}{8} + \left(\frac{320}{8}\right)^2\right) \left(\frac{1}{15118} - 8 \left(\frac{320}{8}\right)^2\right)}$$

$$= 9.271 \pm 7.580$$

$$= (1.7000, 10.841)$$

$$\hat{\beta}_1 = 0.1432, 0.1432 \pm t_{0.025}(6) \sqrt{\frac{\sigma^2}{S_{xx}}}$$

$$= 0.1432 \pm 2.45 \frac{11.74}{15118 - 8 \left(\frac{320}{8}\right)^2}$$

$$= 0.1432 \pm 0.174$$

$$= (0.031, 0.13172)$$

Multivariate Normal Linear Regression and Correlation Analysis.

Multivariate Normal Linear Regression Model.

The "Multivariate Linear Regression" analysis is statistical methodology for predicting values of one or more response variables from the collection of predictor (independent) variables.

The multivariate linear regression model with a single response takes the form:

$$y_i = \underbrace{\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_{p-1} x_{pi}}_{\text{Response}} + \underbrace{\epsilon_i}_{\text{error}}$$

The term linear refers to the fact that the mean is a linear function of the unknown parameters $\beta_0, \beta_1, \dots, \beta_{p-1}$.

Consider, $p=2$

$$y_i = \beta_0 + \beta_1 x_{1i} + \epsilon_i$$

Consider $p=3$

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \beta_3 x_{3i} + \epsilon_i$$

When n independent observations on y and the associated values of x are collected then complete model is given as:

$$y_1 = \beta_0 + \beta_1 x_{11} + \dots + \beta_{p-1} x_{p-11} + \epsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_{12} + \dots + \beta_{p-1} x_{p-12} + \epsilon_2$$

:

$$y_n = \beta_0 + \beta_1 x_{1n} + \dots + \beta_{p-1} x_{p-1n} + \epsilon_n$$

Where $\beta_0, \beta_1, \dots, \beta_{p-1}$ are the unknown parameters and the error term e_i has the following properties: $E(e_i) = 0$

$$\text{Cov}(e_i, e_j) = 0$$

$$V(e_i) = \sigma^2$$

a) Least Squares Method.

- Consider multivariate linear regression.

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{(p-1)i} + e_i$$

- Make e_i the subject, square both sides over the observations

$$\sum e_i^2 = \sum (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_{p-1} x_{(p-1)i})^2$$

$$\text{Let } \phi = \sum_{i=1}^n e_i^2$$

- Differentiate ϕ w.r.t. β_0 and equate to zero to solve for β_0 .

$$\sum (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_{p-1} x_{(p-1)i}) = 0$$

$$\sum (\beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{(p-1)i}) = \sum y_i$$

Differentiate ϕ w.r.t. β_1 and equate to zero to solve for β_1 .

$$\sum (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_{p-1} x_{(p-1)i}) x_{1i} = 0$$

$$\sum (\beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{(p-1)i}) x_{1i} = \sum y_i \cdot x_{1i}$$

Differentiate ϕ w.r.t. β_{p-1} and equate to zero to solve for β_{p-1} .

$$\sum (y_i - \beta_0 - \beta_1 x_{1i} - \dots - \beta_{p-1} x_{(p-1)i}) x_{(p-1)i} = 0$$

$$\sum (\beta_0 + \beta_1 x_{1i} + \dots + \beta_{p-1} x_{(p-1)i}) x_{(p-1)i} = \sum y_i x_{(p-1)i}$$

$$\begin{bmatrix} n & \sum_{i=1}^n x_{1i} & \dots & \sum_{i=1}^n x_{(p-1)i} \\ \sum_{i=1}^n x_{1i} & \sum_{i=1}^n x_{1i}^2 & \dots & \sum_{i=1}^n x_{1i} x_{(p-1)i} \\ \vdots & & & \vdots \\ \sum_{i=1}^n x_{(p-1)i} & \sum_{i=1}^n x_{1i} x_{(p-1)i} & \dots & \sum_{i=1}^n x_{(p-1)i}^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum y_i x_{1i} \\ \vdots \\ \sum y_i x_{(p-1)i} \end{bmatrix}$$

Such that:

$$X\beta = Y$$

Multiply both sides by X'

$$X'X\beta = X'Y$$

Multiply both sides by $(X'X)^{-1}$

$$\hat{\beta} = (X'X)^{-1} X'Y$$

$$\text{Where } \hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_{p-1} \end{bmatrix}, \quad X'X = \begin{bmatrix} n & \sum \hat{x}_{1i} & \dots & \sum \hat{x}_{p-1i} \\ \sum \hat{x}_{1i} & \sum \hat{x}_{1i}^2 & \dots & \sum \hat{x}_{1i} \hat{x}_{p-1i} \\ \vdots & & & \\ \sum \hat{x}_{p-1i} & \sum \hat{x}_{p-1i} \hat{x}_{1i} & \dots & \sum \hat{x}_{p-1i}^2 \end{bmatrix}$$

$$X'y = \begin{bmatrix} \sum y_i \\ \sum x_{1i} y_i \\ \vdots \\ \sum x_{p-1i} y_i \end{bmatrix}$$

Analysis of Variance.

Hypothesis: $H_0: \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$ against $H_a: \beta_i \neq \beta_j$ for atleast one pair.

Analysis of variance component.

Corrected sum of squares.

$$\begin{aligned} \sum_{i=1}^n (y_i - \bar{y})^2 &= \sum_{i=1}^n (y_i - \hat{y}_i + \hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + 2 \sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \\ &= \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 \end{aligned}$$

$$SS_{\text{Total}} = SS_{\text{Reg}} + SS_{\text{Error}}$$

ANOVA table.

| Source of variation | d.f | SS | MS | F |
|---------------------|-----|------------------------------------|---------------------------|-------------------------------------|
| Regression | p-1 | by subtraction | $SS_{\text{Reg}}/(p-1)$ | $MS_{\text{Reg}}/MS_{\text{Error}}$ |
| Error | n-p | $\sum_{i=1}^n (y_i - \hat{y}_i)^2$ | $SS_{\text{Error}}/(n-p)$ | |
| Total | n-1 | $\sum_{i=1}^n (y_i - \bar{y})^2$ | | |

Test Significance of the Regression Model.

Test statistic: $F = \frac{MS_{\text{Reg}}}{MS_{\text{Error}}}$ approximately with p-1, n-p d.f.

Decision:

Reject H_0 if $F > F_{\alpha}(p-1, n-p)$ at α level of significance.

Confidence Interval of Parameters.

The $100(1-\alpha)\%$ confidence interval of the parameters are given as:

$$\hat{\beta}_i \pm t_{\alpha/2} \frac{S}{n-p} \sqrt{C_{ii}}$$

$$\text{where : } S^2 = \frac{1}{n-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$C_{ii} = i^{\text{th}} \text{ value of } (X'X)^{-1}$$

Coefficient of Determination.

$$R^2 = 1 - \frac{SS_{\text{Error}}}{SS_{\text{Total}}} = \frac{SS_{\text{Reg}}}{SS_{\text{Total}}}$$

Adjusted Coefficient of Determination.

$$R^2_{\text{adj}} = 1 - \frac{SS_{\text{Error}} / df}{SS_{\text{Total}} / df} = 1 - \frac{SS_{\text{Error}} / n-p}{SS_{\text{Total}} / n-1}$$

Example :

The data of the dependent variable (y) and independent variables x_1 & x_2 were collected and recorded below.

| | | | | | | | | | | | |
|---|-------|----|----|----|----|----|----|----|----|----|----|
| ✓ | x_1 | 15 | 15 | 16 | 12 | 15 | 17 | 14 | 15 | 15 | 14 |
| | x_2 | 5 | 6 | 7 | 6 | 6 | 8 | 6 | 7 | 6 | 6 |
| | y | 25 | 24 | 23 | 20 | 25 | 28 | 22 | 24 | 25 | 24 |

Determine the following :

- The multiple linear regression model for estimating y using the sample data.
- Test the significance of the multiple linear regression, $\alpha = 0.05$.
- The 95% C.I of the regressor parameters.

Solution.

Consider multiple linear regression model:

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$$

Rewrite the model in matrix form as:

$$Y = X\beta$$

| x_1 | x_2 | y | x_1^2 | $x_1 x_2$ | x_2^2 | $x_1 y$ | $x_2 y$ | y^2 | \hat{y} | $(y - \hat{y})^2$ |
|-------|-------|-----|---------|-----------|---------|---------|---------|-------|-----------|-------------------|
| 15 | 5 | 25 | 225 | 75 | 25 | 375 | 125 | 625 | | |
| 15 | 6 | 24 | 225 | 90 | 36 | 360 | 144 | 516 | | |
| 16 | 7 | 23 | 256 | 112 | 49 | 363 | 161 | 529 | | |
| 12 | 6 | 20 | 144 | 72 | 36 | 240 | 120 | 400 | | |
| 15 | 6 | 25 | 225 | 90 | 36 | 375 | 150 | 625 | | |
| 17 | 8 | 26 | 289 | 136 | 64 | 476 | 224 | 784 | | |
| 14 | 6 | 22 | 196 | 84 | 36 | 308 | 132 | 484 | | |
| 15 | 7 | 24 | 225 | 105 | 49 | 360 | 168 | 576 | | |
| 15 | 6 | 25 | 225 | 90 | 36 | 375 | 150 | 625 | | |
| 14 | 6 | 24 | 196 | 84 | 36 | 336 | 144 | 576 | | |
| | | | 2206 | 938 | 403 | 3573 | 1518 | 5800 | | |