



Development of e-Courses for B.Sc.(Agriculture) Degree Program



STAM 101
STATISTICS

STATISTICS

AUTHOR

TNAU, Tamil Nadu



AGRIMOON.COM

All About Agriculture...

Index

SN	Lecture	Page No
1	Data - definition - Collection of data - Primary and secondary data - Classification of data - Qualitative and quantitative data	5-14
2	Diagrammatic representation of data - uses and limitations - simple, Multiple, Component and percentage bar diagrams - pie chart	15-22
3	Graphical representation - Histogram - Frequency polygon and Frequency curve	23-27
4	Measures of averages - Mean - median - mode - geometric mean - harmonic mean - computation of the above statistics for raw and grouped data - merits and demerits - measures of location - percentiles - quartiles - computation of the above statistics for raw and grouped data	28-50
5	Measures of dispersion - Range, Variance -Standard deviation - co-efficient of variation - computation of the above statistics for raw and grouped data	51-58
6	Probability - Basic concepts-trial- event-equally likely- mutually exclusive - independent event, additive and multiplicative laws. Theoretical distributions discrete and continuous distributions, Binomial distributions-properties	59-71
7	Poisson Distributions - properties, Normal Distributions- properties	72-84
8	Sampling-basic concepts- sampling vs complete enumeration parameter and statistic-sampling methods-simple random sampling and stratified random sampling	85-92
9	Test of significance - Basic concepts - null hypothesis - alternative hypothesis - level of significance - Standard error and its importance - steps in testing	93-106
10	T-test - definition - assumptions - test for equality of two means- independent and paired t test	107-119
11	Attributes- Contingency table - 2x2 contingency table - Test for independence of attributes - test for goodness of fit of mendalian ratio	120-129
12	Correlation - definition - Scatter diagram -Pearson's correlation co-efficient - properties of correlation coefficient	130-135
13	Regression - definition - fitting of simple linear regression equation - testing the significance of the regression coefficient	136-141
14	Design of experiments - basic concepts - treatment - experimental unit - experimental error - basic principle - replication, randomization and local control.	142-144
15	Completely randomized design - description - layout - analysis - advantages and disadvantages	145-148
16	Randomized blocks design - description - layout - analysis - advantages and disadvantages	149-151

17	Latin square design – description – layout – analysis – advantages and disadvantages	152-155
18	Factorial experiments – factor and levels – types – symmetrical and asymmetrical – simple, main and interaction effects – advantages and disadvantages	156-158
19	2 ² Factorial Experiments in RBD – lay out – analysis	159-161
20	2 ³ factorial experiments in RBD – lay out – analysis	162-166
21	Split plot design – layout – ANOVA Table	167-170
22	Strip plot design – layout – ANOVA Table	171-175
23	Long term experiments – ANOVA table – guard rows – optimum plot size – determination methods.	176-179
P1	Diagrammatic and graphic representation – simple, multiple, component and percentage bar diagram – pie chart – histogram. Frequency polygon, frequency curve	180-189
P2	Measures of central tendency – mean median, mode, geometric mean, harmonic mean for raw data	190-195
P3	Measures of central tendency – mean, median, mode, geometric mean and harmonic mean for grouped data	196-205
P4	Measures of dispersion – variance, standard deviation and coefficient of variation for raw data	206-208
P5	Measures of dispersion – variance, standard deviation and coefficient of variation for grouped data	209-212
P6	Selection of simple random sampling using lottery method and random numbers	213-215
P7	Students’s t test – paired and independent t test	216-229
P8	Chi square test – test for association and goodness of fit	230-238
P9	Calculation of Karl Pearson’s correlation coefficient	239-241
P10	Fitting of simple linear regression of y on x	242-245
P11	Formation of ANOVA table for completely Randomised design (CRD) with equal replication and comparison of means using critical difference values	246-252
P12	Formation of ANOVA table for Randomised blocks design (RBD) and comparison of means using critical difference values	253-259
P13	Formation of ANOVA table for Latin square design (LSD) and comparison of means using critical difference values	260-266

Lecture.1

Data – definition – Collection of data – Primary and secondary data – Classification of data – Qualitative and quantitative data.

Basic Concepts

Statistics (Definition)

Quantitative figures are known as data.

Statistics is the science which deals with the

- (i) Collection of data
- (ii) Organization of data or Classification of data
- (iii) Presentation of data
- (iv) Analysis of data
- (v) Interpretation of data

Data and statistics are not same as used commonly.

Example for data

1. No. of farmers in a block.
2. The rainfall over a period of time.
3. Area under paddy crop in a state.

Functions of statistics

Statistics simplifies complexity, presents facts in a definite form, helps in formulation of suitable policies, facilitates comparison and helps in forecasting.

Uses of statistics

Statistics has pervaded almost all spheres of human activities. Statistics is useful in the administration of various states, Industry, business, economics, research workers, banking, insurance companies etc.

Limitations of Statistics

1. Statistical theories can be applied only when there is variability in the experimental material.
2. Statistics deals with only aggregates or groups and not with individual objects.
3. Statistical results are not exact.
4. Statistics can be misused.

Collection of data

Data can be collected by using sampling methods or experiments.

Data

The information collected through censuses and surveys or in a routine manner or other sources is called a raw data. When the raw data are grouped into groups or classes, they are known as grouped data.

There are two types of data

1. Primary data
2. Secondary data.

Primary data

The data which is collected by actual observation or measurement or count is called primary data.

Methods of collection of primary data

Primary data is collected in any one of the following methods

1. Direct personal interviews.
2. Indirect oral interviews
3. Information from correspondents.
4. Mailed questionnaire method.
5. Schedules sent through enumerators.

1. Direct personal interviews

The persons from whom information are collected are known as informants or respondents. The investigator personally meets them and asks questions to gather the necessary information.

Merits

1. The collected informations are likely to be uniform and accurate. The investigator is there to clear the doubts of the informants.
2. People willingly supply information because they are approached personally. Hence more response is noticed in this method than in any other method.

Limitations

It is likely to be very costly and time consuming if the number of persons to be interviewed is large and the persons are spread over a wide area.

2. Indirect oral interviews

Under this method, the investigator contacts witnesses or neighbors or friends or some other third parties who are capable of supplying the necessary information.

Merits

For almost all the surveys of this kind, the informants live within a closed area. Hence, the time and the cost are less. For certain surveys, this is the only method available.

Limitations

The information obtained by this method is not very reliable. The informants and the person who conducts a survey easily distort the truth.

3. Information from correspondents

The investigator appoints local agents or correspondents in different places and compiles the information sent by them.

Merits

- For certain kinds of primary data collection, this is the only method available.
- This method is very cheap and expeditious.
- The quality of data collected is also good due to long experience of local representatives.

Limitations

Local agents and correspondents are not likely to be serious and careful.

4. Mailed Questionnaire method

Under this method a list of questions is prepared and is sent to all the informants by post. The list of questions is technically called questionnaire.

Merits

1. It is relatively cheap.
2. It is preferable when the informants are spread over a wide area.
3. It is fast if the informants respond duly.

Limitations

1. Were the informants are illiterate people, this method cannot be adopted.
2. It is possible that some of the persons who receive the questionnaires do not return them. Their action is known as non – response.

5. Schedules sent through enumerators

Under this method, enumerators or interviewers take the schedules, meet the informants and fill in their replies. A schedule is filled by the interviewer in a face to face situation with the informant.

Merits

1. It can be adopted even if the informants are illiterate.

2. Non-response is almost nil as the enumerators go personally and contact the informants.
3. The informations collected are reliable. The enumerators can be properly trained for the same.

Limitations

1. It is costliest method.
2. Extensive training is to be given to the enumerators for collecting correct and uniform informations.

Secondary data

The data which are compiled from the records of others is called secondary data.

The data collected by an individual or his agents is primary data for him and secondary data for all others. The secondary data are less expensive but it may not give all the necessary information.

Secondary data can be compiled either from published sources or from unpublished sources.

Sources of published data

1. Official publications of the central, state and local governments.
2. Reports of committees and commissions.
3. Publications brought about by research workers and educational associations.
4. Trade and technical journals.
5. Report and publications of trade associations, chambers of commerce, bank etc.
6. Official publications of foreign governments or international bodies like U.N.O, UNESCO etc.

Sources of unpublished data

All statistical data are not published. For example, village level officials maintain records regarding area under crop, crop production etc. They collect details for

administrative purposes. Similarly details collected by private organizations regarding persons, profit, sales etc become secondary data and are used in certain surveys.

Characteristics of secondary data

The secondary data should possess the following characteristics. They should be reliable, adequate, suitable, accurate, complete and consistent.

Variables

Variability is a common characteristic in biological Sciences. A quantitative or qualitative characteristic that varies from observation to observation in the same group is called a variable.

Quantitative data

The basis of classification is according to differences in quantity. In case of quantitative variables the observations are made in terms of kgs, Lt, cm etc. Example weight of seeds, height of plants.

Qualitative data

When the observations are made with respect to quality is called qualitative data.
Eg: Crop varieties, Shape of seeds, soil type.
The qualitative variables are termed as attributes.

Classification of data

Classification is the process of arranging data into groups or classes according to the common characteristics possessed by the individual items.

Data can be classified on the basis of one or more of the following kinds namely

1. Geography
2. Chronology
3. Quality
4. Quantity.

1. Geographical classification (or) Spatial Classification

Some data can be classified area-wise, such as states, towns etc.

Data on area under crop in India can be classified as shown below

Region	Area (in hectares)
Central India	-
West	-
North	-
East	-
South	-

2. Chronological or Temporal or Historical Classification

Some data can be classified on the basis of time and arranged chronologically or historically.

Data on Production of food grains in India can be classified as shown below

Year	Tonnes
1990-91	-
1991-92	-
1992-93	-
1993-94	-
1994-95	-

3. Qualitative Classification

Some data can be classified on the basis of attributes or characteristics. The number of farmers based on their land holdings can be given as follows

Type of farmers	Number of farmers
Marginal	907
Medium	1041
Large	1948
Total	3896

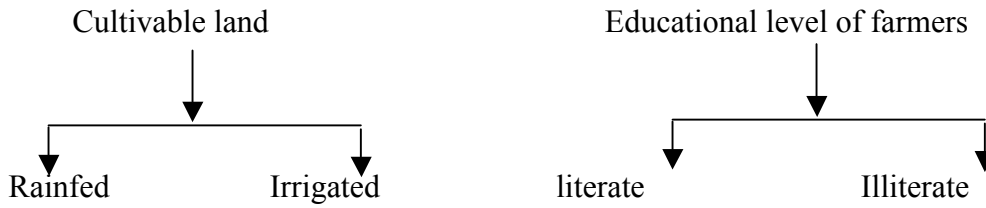
Qualitative classification can be of two types as follows

- (i) Simple classification
- (ii) Manifold classification

(i) Simple Classification

This is based on only one quality.

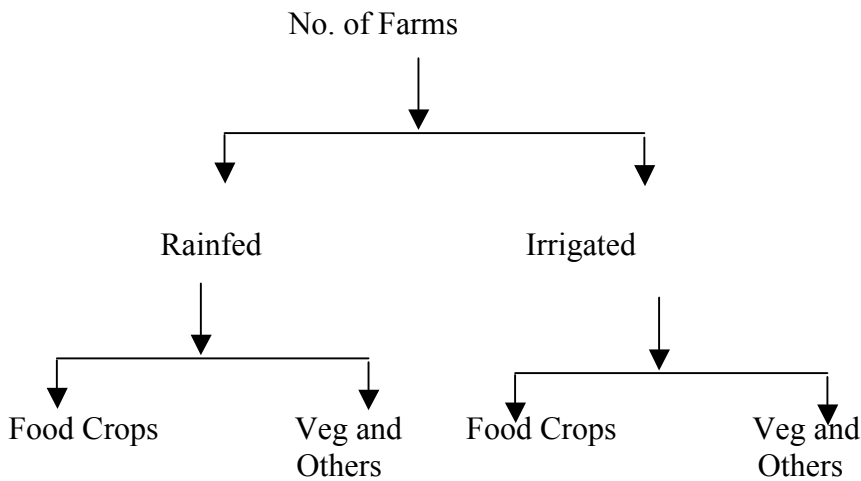
Eg:



(ii) Manifold Classification

This is based on more than one quality.

Eg:



4. Quantitative classification

Some data can be classified in terms of magnitude. The data on land holdings by farmers in a block. Quantitative classification is based the land holding which is the variable in this example.

Land holding (hectare)	Number of Farmers
< 1	442
1-2	908
2-5	471
>5	124
Total	1945

Difference between Primary and secondary data

	Primary Data	Secondary Data
1. Original data	Primary data are original because investigation himself collects them.	Secondary data are not original since investigator makes use of the other agencies.
2. Suitability	If these data are collected accurately and systematically their suitability will be very positive.	These might or might not suit the objectives of enquiry.
3. Time and labour	These data involve large expenses in terms of money, time and manpower	These data are relatively less costly.
4. Precaution	don't need any great precaution while using these data.	These should be used with great care and caution.

Questions

1. A simple table contains data on

- a) Two characteristics
- b) Several characteristics
- c) One characteristic
- d) Three characteristics

Ans: One characteristic

2. When the collected data is grouped with reference to time, we have

- a) Quantitative classification
- b) Qualitative classification
- c) Geographical Classification
- d) Chorological Classification

Ans: Chorological Classification

3. Geographical classification means, classification of data according to Region.

Ans: True

4. An arrangement of data into rows and columns is known as Tabulation.

Ans: True

5. Data on yield is a quantitative variable

Ans: True

6. Qualitative variables are also called as attributes.

Ans: True

7. Define primary and secondary data

8. Give the advantages of tabulation.

9. Write a detail note on the types of classification

10. What are the essential characteristics of a good table?

Lecture.2

Diagrammatic representation of data – uses and limitations – simple, Multiple, Component and percentage bar diagrams – pie chart

Diagrams

Diagrams are various geometrical shape such as bars, circles etc. Diagrams are based on scale but are not confined to points or lines. They are more attractive and easier to understand than graphs.

Merits

1. Most of the people are attracted by diagrams.
2. Technical Knowledge or education is not necessary.
3. Time and effort required are less.
4. Diagrams show the data in proper perspective.
5. Diagrams leave a lasting impression.
6. Language is not a barrier.
7. Widely used tool.

Demerits (or) limitations

1. Diagrams are approximations.
2. Minute differences in values cannot be represented properly in diagrams.
3. Large differences in values spoil the look of the diagram.
4. Some of the diagrams can be drawn by experts only. eg. Pie chart.
5. Different scales portray different pictures to laymen.

Types of Diagrams

The important diagrams are

1. Simple Bar diagram.
2. Multiple Bar diagram.
3. Component Bar diagram.
4. Percentage Bar diagram.

5. Pie chart
6. Pictogram
7. Statistical maps or cartograms.

In all the diagrams and graphs, the groups or classes are represented on the x-axis and the volumes or frequencies are represented in the y-axis.

Simple Bar diagram

If the classification is based on attributes and if the attributes are to be compared with respect to a single character we use simple bar diagram.

Example

1. The area under different crops in a state.
2. The food grain production of different years.
3. The yield performance of different varieties of a crop.
4. The effect of different treatments etc.

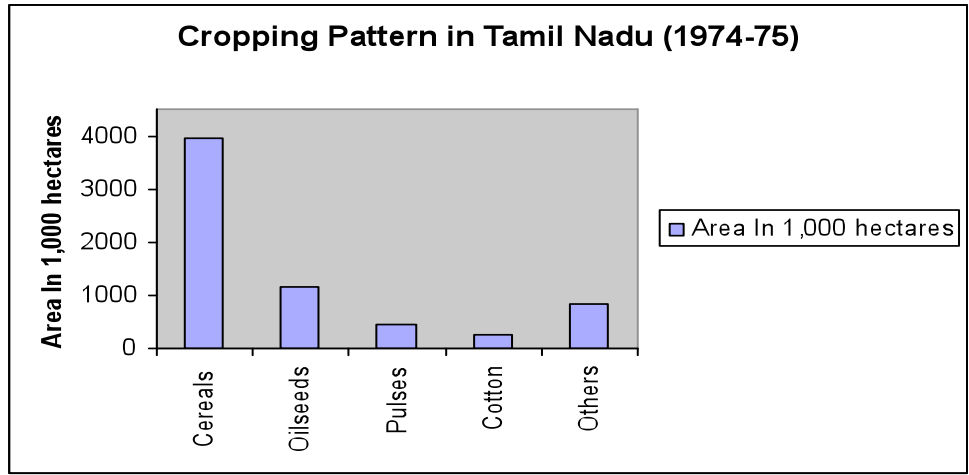
Simple bar diagrams Consists of vertical bars of equal width. The heights of these bars are proportional to the volume or magnitude of the attribute. All bars stand on the same baseline. The bars are separated from each others by equal intervals. The bars may be coloured or marked.

Example

The cropping pattern in Tamil Nadu in the year 1974-75 was as follows.

Crops	Area In 1,000 hectares
Cereals	3940
Oilseeds	1165
Pulses	464
Cotton	249
Others	822

The simple bar diagram for this data is given below.



Multiple bar diagram

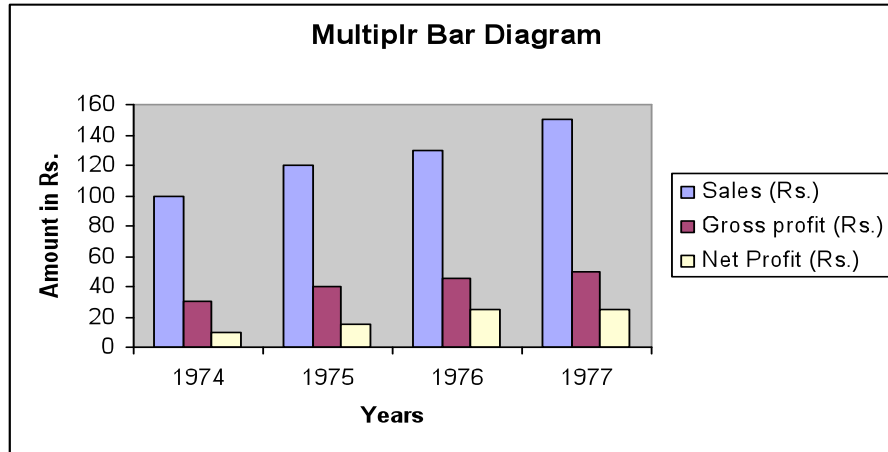
If the data is classified by attributes and if two or more characters or groups are to be compared within each attribute we use multiple bar diagrams. If only two characters are to be compared within each attribute, then the resultant bar diagram used is known as double bar diagram.

The multiple bar diagram is simply the extension of simple bar diagram. For each attribute two or more bars representing separate characters or groups are to be placed side by side. Each bar within an attribute will be marked or coloured differently in order to distinguish them. Same type of marking or colouring should be done under each attribute. A footnote has to be given explaining the markings or colourings.

Example

Draw a multiple bar diagram for the following data which represented agricultural production for the period from 2000-2004

Year	Food grains (tones)	Vegetables (tones)	Others (tones)
2000	100	30	10
2001	120	40	15
2002	130	45	25
2003	150	50	25
2004			



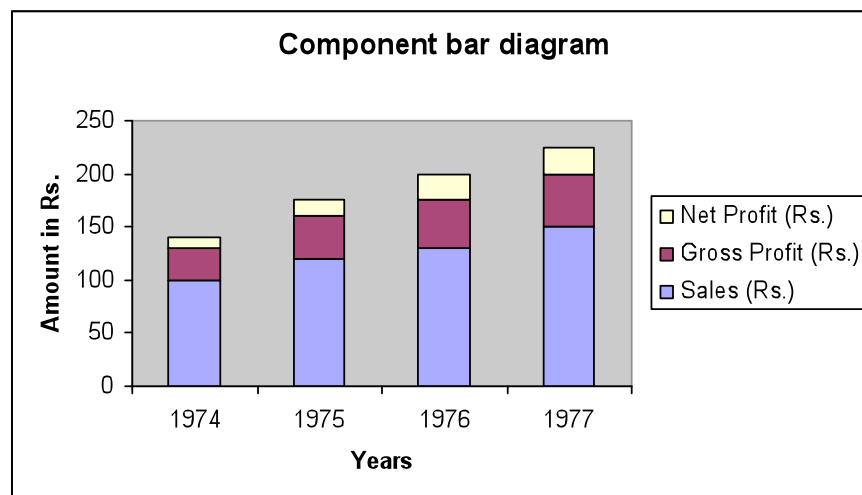
Component bar diagram

This is also called sub – divided bar diagram. Instead of placing the bars for each component side by side we may place these one on top of the other. This will result in a component bar diagram.

Example:

Draw a component bar diagram for the following data

Year	Sales (Rs.)	Gross Profit (Rs.)	Net Profit (Rs.)
1974	100	30	10
1975	120	40	15
1976	130	45	25
1977	150	50	25



Percentage bar diagram

Sometimes when the volumes of different attributes may be greatly different for making meaningful comparisons, the attributes are reduced to percentages. In that case each attribute will have 100 as its maximum volume. This sort of component bar chart is known as percentage bar diagram.

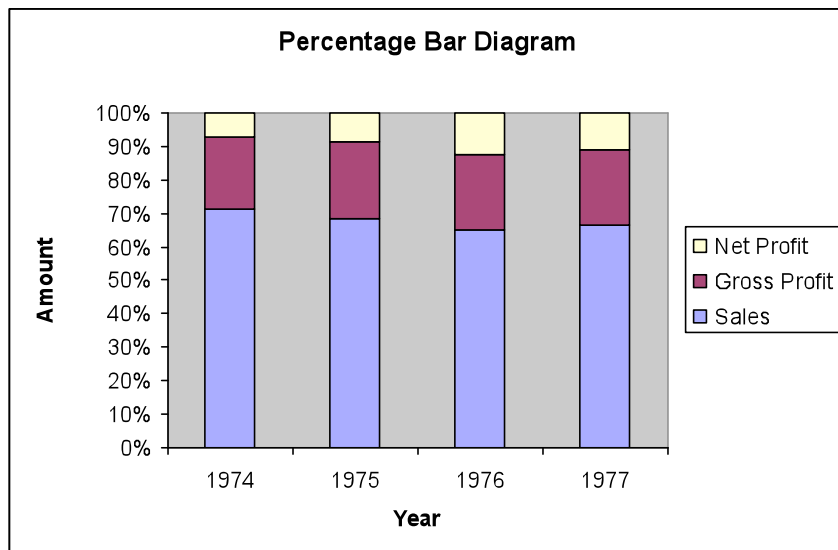
$$\text{Percentage} = \frac{\text{Actual value}}{\text{Total of the actual value}} \times 100,$$

Example:

Draw a Percentage bar diagram for the following data

Using the formula $\text{Percentage} = \frac{\text{Actual value}}{\text{Total of the actual value}} \times 100$, the above table is converted.

Year	Sales (Rs.)	Gross Profit (Rs.)	Net Profit (Rs.)
1974	71.43	21.43	7.14
1975	68.57	22.86	8.57
1976	65	22.5	12.5
1977	66.67	22.22	11.11



Pie chart / Pie Diagram

Pie diagram is a circular diagram. It may be used in place of bar diagrams. It consists of one or more circles which are divided into a number of sectors. In the construction of pie diagram the following steps are involved.

Step 1:

Whenever one set of actual value or percentage are given, find the corresponding angles in degrees using the following formula

$$\text{Angle} = \frac{\text{Actual value}}{\text{Total of the actual value}} \times 360^\circ$$

$$\text{(or) Angle} = \frac{\text{Percentage}}{100} \times 360^\circ$$

Step 2:

Find the radius using the area of the circle πr^2 where value of π is $22/7$ or 3.14

Example

Given the cultivable land area in four southern states of India. Construct a pie diagram for the following data.

State	Cultivable area(in hectares)
Andhra Pradesh	663
Karnataka	448
Kerala	290
Tamil Nadu	556
Total	1957

Using the formula

$$\text{Angle} = \frac{\text{Actual value}}{\text{Total of the actual value}} \times 360^\circ$$

(or)

$$\text{Angle} = \frac{\text{Percentage}}{100} \times 360^\circ$$

The table value becomes

State	Cultivable area
Andhra Pradesh	121.96
Karnataka	82.41
Kerala	53.35
Tamil Nadu	102.28

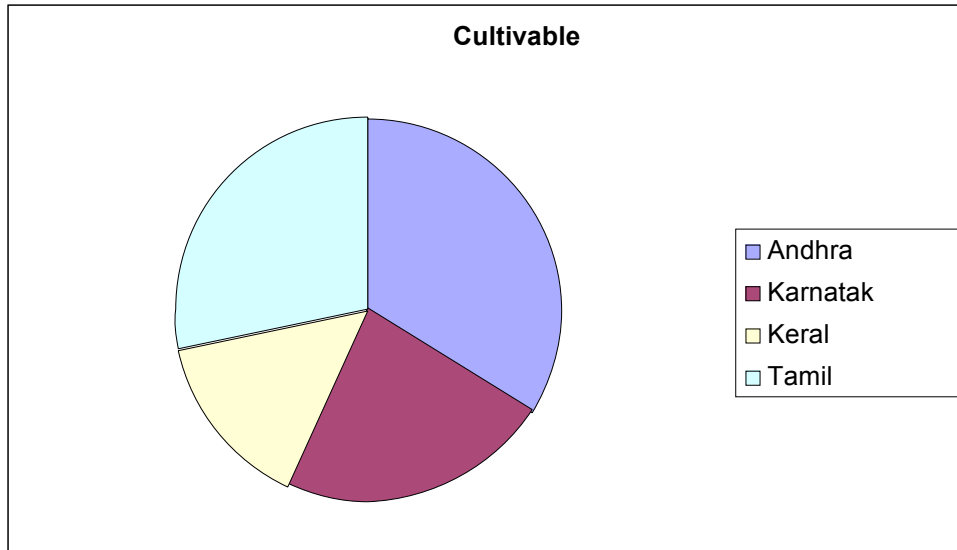
$$\text{Radius} = \pi r^2$$

Here $\pi r^2 = 1957$

$$r^2 = \frac{1957}{\pi} = \frac{1957}{3.14} = 623.24$$

$$r = 24.96$$

$$r = 25 \text{ (approx)}$$



Questions

- In a component bar diagram the length of the bar
 - Will be same for all
 - Depends on the total
 - will not be same
 - none of these

Ans: Depends on the total

2. The length of the bar will be same for all categories in
- a) Multiple bar diagram
 - b) component bar diagram
 - c) Percentage bar diagram
 - d) none of these

Ans: Percentage bar diagram

3. Sub-divided bar diagram are also called Component bar diagram.

Ans: True

4. The multiple bar diagram is the extension of simple bar diagram.

Ans: True

5. In a bar the width of the bars should be equal.

Ans: True

6. In a percentage bar diagram the length of the bars will not be equal.

Ans: False

7. How diagrams are useful in representing statistical data?

8. How to draw a pie chart?

9. Explain how to draw simple and multiple bar diagrams.

10. Explain how to draw Component and percentage bar diagrams.

Lecture.3

Graphical representation – Histogram – Frequency polygon and Frequency curve

Graphs

Graphs are charts consisting of points, lines and curves. Charts are drawn on graph sheets. Suitable scales are to be chosen for both x and y axes, so that the entire data can be presented in the graph sheet. Graphical representations are used for grouped quantitative data.

Histogram

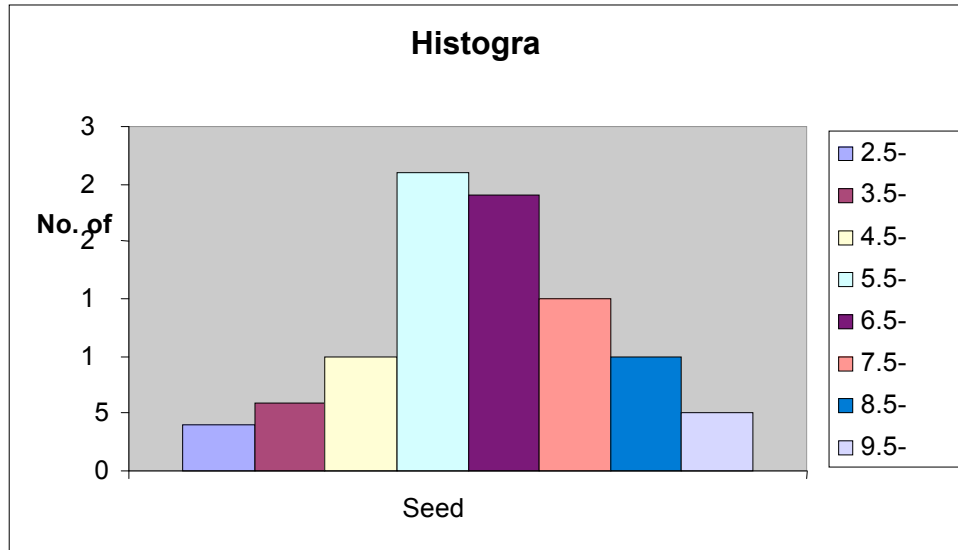
When the data are classified based on the class intervals it can be represented by a histogram. Histogram is just like a simple bar diagram with minor differences. There is no gap between the bars, since the classes are continuous. The bars are drawn only in outline without colouring or marking as in the case of simple bar diagrams. It is the suitable form to represent a frequency distribution.

Class intervals are to be presented in x axis and the bases of the bars are the respective class intervals. Frequencies are to be represented in y axis. The heights of the bars are equal to the corresponding frequencies.

Example

Draw a histogram for the following data

Seed Yield (gms)	No. of Plants
2.5-3.5	4
3.5-4.5	6
4.5-5.5	10
5.5-6.5	26
6.5-7.5	24
7.5-8.5	15
8.5-9.5	10
9.5-10.5	5



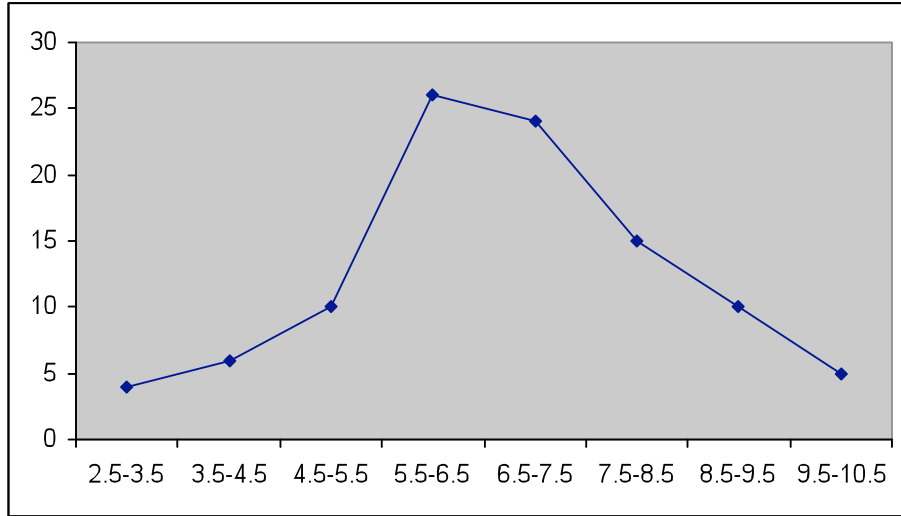
Frequency Polygon

The frequencies of the classes are plotted by dots against the mid-points of each class. The adjacent dots are then joined by straight lines. The resulting graph is known as frequency polygon.

Example

Draw frequency polygon for the following data

Seed Yield (gms)	No. of Plants
2.5-3.5	4
3.5-4.5	6
4.5-5.5	10
5.5-6.5	26
6.5-7.5	24
7.5-8.5	15
8.5-9.5	10
9.5-10.5	5



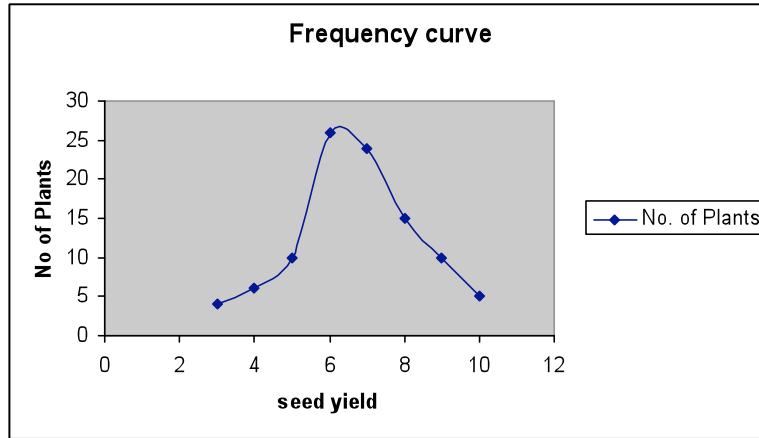
Frequency curve

The procedure for drawing a frequency curve is same as for frequency polygon. But the points are joined by smooth or free hand curve.

Example

Draw frequency curve for the following data

Seed Yield (gms)	No. of Plants
2.5-3.5	4
3.5-4.5	6
4.5-5.5	10
5.5-6.5	26
6.5-7.5	24
7.5-8.5	15
8.5-9.5	10
9.5-10.5	5



Ogives

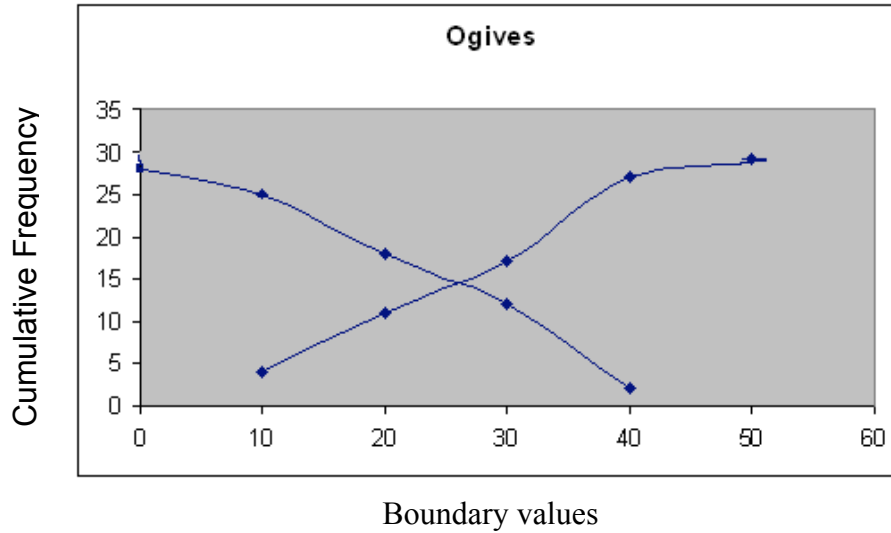
Ogives are known also as cumulative frequency curves and there are two kinds of ogives. One is less than ogive and the other is more than ogive.

Less than ogive: Here the cumulative frequencies are plotted against the upper boundary of respective class interval.

Greater than ogive: Here the cumulative frequencies are plotted against the lower boundaries of respective class intervals.

Example

Continuous Interval	Mid Point	Frequency	< cumulative Frequency	> cumulative frequency
0-10	5	4	4	29
10-20	15	7	11	25
20-30	25	6	17	18
30-40	35	10	27	12
40-50	45	2	29	2



Questions

- With the help of histogram we can draw
 - Frequency polygon
 - frequency curve
 - Frequency distribution
 - all the above
- Ogives for more than type and less than type distribution intersect at
 - Mean
 - median
 - Mode
 - origin

Ans: all the above

Ans: median

- To draw the frequency polygon we take the mid values in the X axis.
- To draw the frequency polygon we take the mid values in the X axis.
- In a frequency curve the points are joined by bits of straight lines
Ans: False
- He stogram can be drawn for equal and unequal classes
Ans: True
- Explain how to draw frequency curve
- Explain how to draw histogram.
- Explain the diagrams that can be drawn for a frequency distribution table
- Explain how to draw less than and more than Ogives.

Lecture.4

Measures of averages - Mean – median – mode – geometric mean – harmonic mean – computation of the above statistics for raw and grouped data - merits and demerits - measures of location – percentiles – quartiles - computation of the above statistics for raw and grouped data

In the study of a population with respect to one in which we are interested we may get a large number of observations. It is not possible to grasp any idea about the characteristic when we look at all the observations. So it is better to get one number for one group. That number must be a good representative one for all the observations to give a clear picture of that characteristic. Such representative number can be a central value for all these observations. This central value is called a measure of central tendency or an average or a measure of locations. There are five averages. Among them mean, median and mode are called simple averages and the other two averages geometric mean and harmonic mean are called special averages.

Arithmetic mean or mean

Arithmetic mean or simply the mean of a variable is defined as the sum of the observations divided by the number of observations. It is denoted by the symbol \bar{x} . If the variable x assumes n values $x_1, x_2 \dots x_n$ then the mean is given by

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n x_i$$

This formula is for the ungrouped or raw data.

Example 1

Calculate the mean for pH levels of soil 6.8, 6.6, 5.2, 5.6, 5.8

Solution

$$\bar{x} = \frac{6.8 + 6.6 + 5.2 + 5.6 + 5.8}{5} = \frac{30}{5} = 6$$

Grouped Data

The mean for grouped data is obtained from the following formula:

$$\bar{x} = \frac{\sum fx}{n}$$

Where x = the mid-point of individual class

f = the frequency of individual class

n = the sum of the frequencies or total frequencies in a sample.

Short-cut method

$$\bar{x} = A + \frac{\sum fd}{n} \times c$$

Where $d = \frac{x - A}{c}$

A = any value in x

n = total frequency

c = width of the class interval

Example 2

Given the following frequency distribution, calculate the arithmetic mean

Marks	:	64	63	62	61	60	59
Number of Students	} :	8	18	12	9	7	6

Solution

X	F	F _x	D=x-A	Fd
64	8	512	2	16
63	18	1134	1	18
62	12	744	0	0
61	9	549	-1	-9
60	7	420	-2	-14
59	6	354	-3	-18
	60	3713		-7

Direct method

$$\bar{x} = \frac{\sum fx}{n}$$

$$\bar{x} = \frac{3713}{60} = 61.88$$

Short-cut method

$$\bar{x} = A + \frac{\sum fd}{n} \times c$$

Here A = 62

$$\bar{x} = 62 - \frac{7}{60} \times 1 = 61.88$$

Example 3

For the frequency distribution of seed yield of sesamum given in table, calculate the mean yield per plot.

Yield per plot in(in g)	64.5-84.5	84.5-104.5	104.5-124.5	124.5-144.5
No of plots	3	5	7	20

Solution

Yield (in g)	No of Plots (f)	Mid X	$d = \frac{x - A}{c}$	Fd
64.5-84.5	3	74.5	-1	-3
84.5-104.5	5	94.5	0	0
104.5-124.5	7	114.5	1	7
124.5-144.5	20	134.5	2	40
Total	35			44

A=94.5

The mean yield per plot is

Direct method:

$$\begin{aligned} \bar{x} &= \frac{\sum fx}{n} = \frac{(74.5 \times 3) + (94.5 \times 5) + (114.5 \times 7) + (134.5 \times 20)}{35} \\ &= \frac{4187.5}{35} = 119.64 \text{ gms} \end{aligned}$$

Shortcut method

$$\bar{x} = A + \frac{\sum fd}{n} \times c$$

$$\bar{x} = 94.5 + \frac{44}{35} \times 20 = 119.64 \text{ g}$$

Merits and demerits of Arithmetic mean

Merits

1. It is rigidly defined.
2. It is easy to understand and easy to calculate.
3. If the number of items is sufficiently large, it is more accurate and more reliable.
4. It is a calculated value and is not based on its position in the series.
5. It is possible to calculate even if some of the details of the data are lacking.
6. Of all averages, it is affected least by fluctuations of sampling.
7. It provides a good basis for comparison.

Demerits

1. It cannot be obtained by inspection nor located through a frequency graph.
2. It cannot be in the study of qualitative phenomena not capable of numerical measurement i.e. Intelligence, beauty, honesty etc.,
3. It can ignore any single item only at the risk of losing its accuracy.
4. It is affected very much by extreme values.
5. It cannot be calculated for open-end classes.
6. It may lead to fallacious conclusions, if the details of the data from which it is computed are not given.

Median

The median is the middle most item that divides the group into two equal parts, one part comprising all values greater, and the other, all values less than that item.

Ungrouped or Raw data

Arrange the given values in the ascending order. If the number of values are odd, median is the middle value

If the number of values are even, median is the mean of middle two values.

By formula

$$\text{When } n \text{ is odd, Median} = Md = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ value}$$

When n is even, Average of $\left(\frac{n}{2}\right)$ and $\left(\frac{n}{2} + 1\right)^{th}$ value

Example 4

If the weights of sorghum ear heads are 45, 60, 48, 100, 65 gms, calculate the median

Solution

Here n = 5

First arrange it in ascending order

45, 48, 60, 65, 100

$$\begin{aligned} \text{Median} &= \left(\frac{n+1}{2}\right)^{th} \text{ value} \\ &= \left(\frac{5+1}{2}\right) = 3^{rd} \text{ value} = 60 \end{aligned}$$

Example 5

If the sorghum ear- heads are 5, 48, 60, 65, 65, 100 gms, calculate the median.

Solution

Here n = 6

$$\begin{aligned} \text{Median} &= \text{Average of } \left(\frac{n}{2}\right) \text{ and } \left(\frac{n}{2} + 1\right)^{th} \text{ value} \\ \left(\frac{n}{2}\right) &= \frac{6}{2} = 3^{rd} \text{ value} = 60 \quad \text{and} \quad \left(\frac{n}{2} + 1\right) = \frac{6}{2} + 1 = 4^{th} \text{ value} = 65 \\ \text{Median} &= \frac{60 + 65}{2} = 62.5 \text{ g} \end{aligned}$$

Grouped data

In a grouped distribution, values are associated with frequencies. Grouping can be in the form of a discrete frequency distribution or a continuous frequency distribution. Whatever may be the type of distribution, cumulative frequencies have to be calculated to know the total number of items.

Cumulative frequency (cf)

Cumulative frequency of each class is the sum of the frequency of the class and the frequencies of the pervious classes, ie adding the frequencies successively, so that the last cumulative frequency gives the total number of items.

Discrete Series

Step1: Find cumulative frequencies.

Step2 : Find $\left(\frac{n}{2} + 1\right)$

Step3: See in the cumulative frequencies the value just greater than $\left(\frac{n}{2} + 1\right)$

Step4: Then the corresponding value of x is median.

Example 6

The following data pertaining to the number of insects per plant. Find median number of insects per plant.

Number of insects per plant (x)	1	2	3	4	5	6	7	8	9	10	11	12
No. of plants(f)	1	3	5	6	10	13	9	5	3	2	2	1

Solution

Form the cumulative frequency table

x	f	cf
1	1	1
2	3	4
3	5	9
4	6	15
5	10	25
6	13	38
7	9	47
8	5	52
9	3	55
10	2	57
11	2	59
12	1	60
	60	

Median = size of $\left(\frac{n+1}{2}\right)^{th}$ item

Here the number of observations is even. Therefore median = average of (n/2)th item and (n/2+1)th item.

$$= (30^{\text{th}} \text{ item} + 31^{\text{st}} \text{ item}) / 2 = (6+6)/2 = 6$$

Hence the median size is 6 insects per plant.

Continuous Series

The steps given below are followed for the calculation of median in continuous series.

Step1: Find cumulative frequencies.

Step2: Find $\left(\frac{n}{2}\right)$

Step3: See in the cumulative frequency the value first greater than $\left(\frac{n}{2}\right)$, Then the corresponding

class interval is called the Median class. Then apply the formula

$$\text{Median} = l + \frac{\frac{n}{2} - m}{f} \times c$$

where l = Lower limit of the medianal class

m = cumulative frequency preceding the medianal class

c = width of the class

f = frequency in the median class.

n = Total frequency.

Example 7

For the frequency distribution of weights of sorghum ear-heads given in table below. Calculate the median.

Weights of ear heads (in g)	No of ear heads (f)	Less than class	Cumulative frequency (m)
60-80	22	<80	22
80-100	38	<100	60
100-120	45	<120	105
120-140	35	<140	140
140-160	24	<160	164
Total	164		

Solution

$$\text{Median} = l + \frac{\frac{n}{2} - m}{f} \times c$$

$$\left(\frac{n}{2}\right) = \left(\frac{164}{2}\right) = 82$$

It lies between 60 and 105. Corresponding to 60 the less than class is 100 and corresponding to 105 the less than class is 120. Therefore the medianal class is 100-120. Its lower limit is 100.

Here $l = 100, n=164, f = 45, c = 20, m = 60$

$$\text{Median} = 100 + \frac{82 - 60}{45} \times 20 = 109.78 \text{ gms}$$

Merits of Median

1. Median is not influenced by extreme values because it is a positional average.
2. Median can be calculated in case of distribution with open-end intervals.
3. Median can be located even if the data are incomplete.

Demerits of Median

1. A slight change in the series may bring drastic change in median value.
2. In case of even number of items or continuous series, median is an estimated value other than any value in the series.
3. It is not suitable for further mathematical treatment except its use in calculating mean deviation.
4. It does not take into account all the observations.

Mode

The mode refers to that value in a distribution, which occur most frequently. It is an actual value, which has the highest concentration of items in and around it. It shows the centre of concentration of the frequency in around a given value. Therefore, where the purpose is to know the point of the highest concentration it is preferred. It is, thus, a positional measure.

Its importance is very great in agriculture like to find typical height of a crop variety, maximum source of irrigation in a region, maximum disease prone paddy variety. Thus the mode is an important measure in case of qualitative data.

Computation of the mode

Ungrouped or Raw Data

For ungrouped data or a series of individual observations, mode is often found by mere inspection.

Example 8

Find the mode for the following seed weight

2 , 7, 10, 15, 10, 17, 8, 10, 2 gms

∴ Mode = 10

In some cases the mode may be absent while in some cases there may be more than one mode.

Example 9

(1) 12, 10, 15, 24, 30 (no mode)

(2) 7, 10, 15, 12, 7, 14, 24, 10, 7, 20, 10

the modal values are 7 and 10 as both occur 3 times each.

Grouped Data

For Discrete distribution, see the highest frequency and corresponding value of x is mode.

Example:

Find the mode for the following

Weight of sorghum in gms (x)	No. of ear head(f)
50	4
65	6
75	16
80	8
95	7
100	4

Solution

The maximum frequency is 16. The corresponding x value is 75.

∴ mode = 75 gms.

Continuous distribution

Locate the highest frequency the class corresponding to that frequency is called the modal class.

Then apply the formula.

$$\text{Mode} = l + \frac{f_s}{f_p + f_s} \times c$$

Where l = lower limit of the modal class

f_p = the frequency of the class preceding the modal class

f_s = the frequency of the class succeeding the modal class

and c = class interval

Example 10

For the frequency distribution of weights of sorghum ear-heads given in table below. Calculate the mode

Weights of ear heads (g)	No of ear heads (f)	
60-80	22	
80-100	38	f_p
100-120	45	f
120-140	35	f_s
140-160	20	
Total	160	

Solution

$$\text{Mode} = l + \frac{f_s}{f_p + f_s} \times c$$

Here $l = 100$, $f = 45$, $c = 20$, $m = 60$, $f_p = 38$, $f_s = 35$

$$\begin{aligned} \text{Mode} &= 100 + \frac{35_s}{38 + 35} \times 20 \\ &= 100 + \frac{35_s}{73} \times 20 \\ &= 109.589 \end{aligned}$$

Geometric mean

The geometric mean of a series containing n observations is the n th root of the product of the values.

If x_1, x_2, \dots, x_n are observations then

$$\begin{aligned} \text{G.M} &= \sqrt[n]{x_1, x_2, \dots, x_n} \\ &= (x_1, x_2, \dots, x_n)^{1/n} \end{aligned}$$

$$\begin{aligned} \text{Log GM} &= \frac{1}{n} \log(x_1, x_2, \dots, x_n) \\ &= \frac{1}{n} (\log x_1 + \log x_2 + \dots + \log x_n) \\ &= \frac{\sum \log x_i}{n} \end{aligned}$$

$$\text{GM} = \text{Antilog } \frac{\sum \log x_i}{n}$$

For grouped data

$$\text{GM} = \text{Antilog } \left[\frac{\sum f \log x_i}{n} \right]$$

GM is used in studies like bacterial growth, cell division, etc.

Example 11

If the weights of sorghum ear heads are 45, 60, 48, 100, 65 gms. Find the Geometric mean for the following data

Weight of ear head x (g)	Log x
45	1.653
60	1.778
48	1.681
100	2.000
65	1.813
Total	8.925

Solution

Here n = 5

$$\begin{aligned} \text{GM} &= \text{Antilog } \frac{\sum \log x_i}{n} \\ &= \text{Antilog } \frac{8.925}{5} \\ &= \text{Antilog } 1.785 \end{aligned}$$

= 60.95

Grouped Data

Example 12

Find the Geometric mean for the following

Weight of sorghum (x)	No. of ear head(f)
50	4
65	6
75	16
80	8
95	7
100	4

Solution

Weight of sorghum (x)	No. of ear head(f)	Log x	f x log x
50	5	1.699	8.495
63	10	10.799	17.99
65	5	1.813	9.065
130	15	2.114	31.71
135	15	2.130	31.95
Total	50	9.555	99.21

Here n= 50

$$GM = \text{Antilog} \left[\frac{\sum f \log x_i}{n} \right]$$

$$= \text{Antilog} \left[\frac{99.21}{50} \right]$$

$$= \text{Antilog } 1.9842 = 96.43$$

Continuous distribution

Example 13

For the frequency distribution of weights of sorghum ear-heads given in table below.

Calculate the Geometric mean

Weights of ear heads (in g)	No of ear heads (f)
60-80	22
80-100	38
100-120	45

120-140	35
140-160	20
Total	160

Solution

Weights of ear heads (in g)	No of ear heads (f)	Mid x	Log x	f log x
60-80	22	70	1.845	40 59
80-100	38	90	1.954	74.25
100-120	45	110	2.041	91.85
120-140	35	130	2.114	73.99
140-160	20	150	2.176	43.52
Total	160			324.2

Here $n = 160$

$$\begin{aligned}
 \text{GM} &= \text{Antilog} \left[\frac{\sum f \log x_i}{n} \right] \\
 &= \text{Antilog} \left[\frac{324.2}{160} \right] \\
 &= \text{Antilog} [2.02625] \\
 &= 106.23
 \end{aligned}$$

Harmonic mean (H.M)

Harmonic mean of a set of observations is defined as the reciprocal of the arithmetic average of the reciprocal of the given values. If x_1, x_2, \dots, x_n are n observations,

$$\text{H.M} = \frac{n}{\sum_{i=1}^n \left(\frac{1}{x_i} \right)}$$

For a frequency distribution

$$\text{H.M} = \frac{n}{\sum_{i=1}^n f \left(\frac{1}{x_i} \right)}$$

H.M is used when we are dealing with speed, rates, etc.

Example 13

From the given data 5, 10,17,24,30 calculate H.M.

X	$\frac{1}{x}$
5	0.2000
10	0.1000
17	0.0588
24	0.0417
30	0.4338

$$H.M = \frac{5}{0.4338} = 11.526$$

Example 14

Number of tomatoes per plant are given below. Calculate the harmonic mean.

Number of tomatoes per plant	20	21	22	23	24	25
Number of plants	4	2	7	1	3	1

Solution

Number of tomatoes per plant (x)	No of plants(f)	$\frac{1}{x}$	$f\left(\frac{1}{x}\right)$
20	4	0.0500	0.2000
21	2	0.0476	0.0952
22	7	0.0454	0.3178
23	1	0.0435	0.0435
24	3	0.0417	0.1251
25	1	0.0400	0.0400
	18		0.8216

$$H.M = \frac{n}{\sum f\left(\frac{1}{x_i}\right)} = \frac{18}{0.1968} = 21.91$$

Merits of H.M

1. It is rigidly defined.
2. It is defined on all observations.
3. It is amenable to further algebraic treatment.
4. It is the most suitable average when it is desired to give greater weight to smaller observations and less weight to the larger ones.

Demerits of H.M

1. It is not easily understood.
2. It is difficult to compute.
3. It is only a summary figure and may not be the actual item in the series
4. It gives greater importance to small items and is therefore, useful only when small items have to be given greater weightage.
5. It is rarely used in grouped data.

Percentiles

The percentile values divide the distribution into 100 parts each containing 1 percent of the cases. The x^{th} percentile is that value below which x percent of values in the distribution fall. It may be noted that the median is the 50^{th} percentile.

For raw data, first arrange the n observations in increasing order. Then the x^{th} percentile is given by

$$P_x = \left(\frac{x(n+1)}{100} \right)^{\text{th}} \text{ item}$$

For a frequency distribution the x^{th} percentile is given by

$$P_x = l + \left(\frac{(x.n/100) - cf}{f} \times C \right)$$

Where

l = lower limit of the percentile class which contains the x^{th} percentile value ($x.n/100$)

cf = cumulative frequency upto l

f = frequency of the percentile class

C = class interval

N = total number of observations

Percentile for Raw Data or Ungrouped Data

Example 15

The following are the paddy yields (kg/plot) from 14 plots:

30,32,35,38,40,42,48,49,52,55,58,60,62, and 65 (after arranging in ascending order). The

computation of 25^{th} percentile (Q_1) and 75^{th} percentile (Q_3) are given below:

$$\begin{aligned}
 P_{25} \text{ (or } Q_1) &= \left(\frac{25(14+1)}{100} \right)^{\text{th}} \text{ item} \\
 &= \left(3\frac{3}{4} \right)^{\text{th}} \text{ item} \\
 &= 3^{\text{rd}} \text{ item} + (4^{\text{th}} \text{ item} - 3^{\text{rd}} \text{ item}) \left(\frac{3}{4} \right) \\
 &= 35 + (38-35) \left(\frac{3}{4} \right) \\
 &= 35 + 3 \left(\frac{3}{4} \right) = 37.25 \text{ kg}
 \end{aligned}$$

$$\begin{aligned}
 P_{75} \text{ (or } Q_3) &= \left(\frac{75(14+1)}{100} \right)^{\text{th}} \text{ item} \\
 &= \left(11\frac{1}{4} \right)^{\text{th}} \text{ item} \\
 &= 11^{\text{th}} \text{ item} + (12^{\text{th}} \text{ item} - 11^{\text{th}} \text{ item}) \left(\frac{1}{4} \right) \\
 &= 55 + (58-55) \left(\frac{1}{4} \right) \\
 &= 55 + 3 \left(\frac{1}{4} \right) = 55.75 \text{ kg}
 \end{aligned}$$

Example 16

The frequency distribution of weights of 190 sorghum ear-heads are given below. Compute 25th percentile and 75th percentile.

Weight of ear-heads (in g)	No of ear heads
40-60	6
60-80	28
80-100	35
100-120	55
120-140	30
140-160	15
160-180	12
180-200	9
Total	190

Solution

Weight of ear-heads (in g)	No of ear heads	Less than class	Cumulative frequency
40-60	6	< 60	6
60-80	28	< 80	34
80-100	35	<100	69
100-120	55	<120	124
120-140	30	<140	154
140-160	15	<160	169
160-180	12	<180	181
180-200	9	<200	190
Total	190		

or P_{25} , first find out $\left(\frac{25(190)}{100}\right)$, and for P_{75} , $\left(\frac{75(190)}{100}\right)$, and proceed as in the case of median.

For P_{25} , we have $\left(\frac{25(190)}{100}\right) = 47.5$.

The value 47.5 lies between 34 and 69. Therefore, the percentile class is 80-100. Hence,

$$\begin{aligned}
 P_{25} = Q_1 &= l + \left(\frac{(25.n/100) - cf}{f} \times C \right) \\
 &= 80 + \left(\frac{(47.5) - 34}{35} \times 20 \right) \\
 &= 80 + \left(\frac{(13.5)}{35} \times 20 \right) \\
 &= 80 + 7.71 \text{ or } 87.71 \text{ g.}
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 P_{75} &= l + \left(\frac{(75.n/100) - cf}{f} \times C \right) \text{ Class} \\
 &= 120 + \left(\frac{(142.5) - 121}{30} \times 20 \right) \\
 &= 120 + \left(\frac{(21.5)}{30} \times 20 \right) \\
 &= 120 + 14.33 = 134.33 \text{ g.}
 \end{aligned}$$

Quartiles

The quartiles divide the distribution in four parts. There are three quartiles. The second quartile divides the distribution into two halves and therefore is the same as the median. The first (lower) quartile (Q1) marks off the first one-fourth, the third (upper) quartile (Q3) marks off the three-fourth. It may be noted that the second quartile is the value of the median and 50th percentile.

Raw or ungrouped data

First arrange the given data in the increasing order and use the formula for Q1 and Q3 then quartile deviation, Q.D is given by

$$Q.D = \frac{Q_3 - Q_1}{2}$$

Where $Q_1 = \left(\frac{n+1}{4}\right)^{th}$ item and $Q_3 = 3\left(\frac{n+1}{4}\right)^{th}$ item

Example 18

Compute quartiles for the data given below (grains/panicles) 25, 18, 30, 8, 15, 5, 10, 35, 40, 45

Solution

5, 8, 10, 15, 18, 25, 30, 35, 40, 45

$$\begin{aligned} Q_1 &= \left(\frac{n+1}{4}\right)^{th} \\ &= \left(\frac{10+1}{4}\right)^{th} \\ &= (2.75)^{th} \text{ item} \\ &= 2^{nd} \text{ item} + \left(\frac{3}{4}\right)(3^{rd} \text{ item} - 2^{nd} \text{ item}) \\ &= 8 + \frac{3}{4}(10-8) \\ &= 8 + \frac{3}{4} \times 2 \end{aligned}$$

$$\begin{aligned}
 &= 8 + 1.5 \\
 &= 9.5 \\
 Q_3 &= 3\left(\frac{n+1}{4}\right)^{\text{th}} \\
 &= 3 \times (2.75)^{\text{th}} \text{ item} \\
 &= (8.75)^{\text{th}} \text{ item} \\
 &= 8^{\text{th}} \text{ item} + \left(\frac{1}{4}\right)(9^{\text{th}} \text{ item} - 8^{\text{th}} \text{ item}) \\
 &= 35 + \frac{1}{4}(40 - 35) \\
 &= 35 + 1.25 \\
 &= 36.25
 \end{aligned}$$

Discrete Series

Step1: Find cumulative frequencies.

Step2: Find $\left(\frac{n+1}{4}\right)$

Step3: See in the cumulative frequencies, the value just greater than $\left(\frac{n+1}{4}\right)$, then the corresponding value of x is Q_1

Step4: Find $3\left(\frac{n+1}{4}\right)$

Step5: See in the cumulative frequencies, the value just greater than $3\left(\frac{n+1}{4}\right)$, then the corresponding value of x is Q_3

Example 19

Compute quartiles for the data given bellow (insects/plant).

X	5	8	12	15	19	24	30
f	4	3	2	4	5	2	4

Solution

x	f	cf
5	4	4
8	3	7
12	2	9
15	4	13
19	5	18
24	2	20

$$Q_1 = \left(\frac{n+1}{4}\right)^{th} \text{ item} = \left(\frac{24+1}{4}\right) = \left(\frac{25}{4}\right) = 6.25^{th} \text{ item}$$

$$Q_3 = 3\left(\frac{n+1}{4}\right)^{th} \text{ item} = 3\left(\frac{24+1}{4}\right) = 18.75^{th} \text{ item} \therefore Q_1 = 8; Q_3 = 24$$

Continuous series

Step1: Find cumulative frequencies

Step2: Find $\left(\frac{n}{4}\right)$

Step3: See in the cumulative frequencies, the value just greater than $\left(\frac{n}{4}\right)$, then the corresponding class interval is called first quartile class.

Step4: Find $3\left(\frac{n}{4}\right)$ See in the cumulative frequencies the value just greater than $3\left(\frac{n}{4}\right)$ then the corresponding class interval is called 3rd quartile class. Then apply the respective formulae

$$Q_1 = l_1 + \frac{\frac{n}{4} - m_1}{f_1} \times c_1$$

$$Q_3 = l_3 + \frac{3\left(\frac{n}{4}\right) - m_3}{f_3} \times c_3$$

Where l_1 = lower limit of the first quartile class

f_1 = frequency of the first quartile class

c_1 = width of the first quartile class

m_1 = c.f. preceding the first quartile class

l_3 = lower limit of the 3rd quartile class

f_3 = frequency of the 3rd quartile class

c_3 = width of the 3rd quartile class

m_3 = c.f. preceding the 3rd quartile class

Example 20: The following series relates to the marks secured by students in an examination.

Marks	No. of Students
0-10	11
10-20	18
20-30	25
30-40	28
40-50	30
50-60	33
60-70	22
70-80	15
80-90	12
90-100	10

Find the quartiles

Solution

C.I	f	cf
0-10	11	11
10-20	18	29
20-30	25	54
30-40	28	82
40-50	30	112
50-60	33	145
60-70	22	167
70-80	15	182
80-90	12	194
90-100	10	204
	204	

$$\left(\frac{n}{4}\right) = \left(\frac{204}{4}\right) = 51$$

$$3\left(\frac{n}{4}\right) = 153$$

$$Q_1 = l_1 + \frac{\frac{n}{4} - m_1}{f_1} \times c_1$$
$$= 20 + \frac{51 - 29}{25} \times 10 = 20 + 8.8 = 28.8$$

$$Q_3 = l_3 + \frac{3\left(\frac{n}{4}\right) - m_3}{f_3} \times c_3$$
$$= 60 + \frac{153 - 145}{22} \times 12 = 60 + 4.36 = 64.36$$

Questions

1. The middle value of an ordered series is called
- a) 2nd quartile b) 5th decile
- c) 50th percentile d) all the above

Ans: all the above

2. For a set of values the modal value can be
- a) Unimodal b) bimodal
- c) Trimodal d) All of these

d) Ans: all the above

3. Mode is suitable for qualitative data.

Ans: True

4. Decile divides the group into ten equal parts.

Ans: True

5. Mean is affected by extreme values.

Ans: True

6. Geometric mean can be calculated for negative values.

Ans: False

7. Define mean and median

Statistics

8. For what type of data mode can be calculated.
9. Explain how to calculate the arithmetic mean for raw and grouped data.
10. Explain how to calculate median and mode for grouped data.

Lecture.5

Measures of dispersion - Range, Variance -Standard deviation – co-efficient of variation - computation of the above statistics for raw and grouped data

Measures of Dispersion

The averages are representatives of a frequency distribution. But they fail to give a complete picture of the distribution. They do not tell anything about the scatterness of observations within the distribution.

Suppose that we have the distribution of the yields (kg per plot) of two paddy varieties from 5 plots each. The distribution may be as follows

Variety I	45	42	42	41	40
Variety II	54	48	42	33	30

It can be seen that the mean yield for both varieties is 42 kg but cannot say that the performances of the two varieties are same. There is greater uniformity of yields in the first variety whereas there is more variability in the yields of the second variety. The first variety may be preferred since it is more consistent in yield performance.

Form the above example it is obvious that a measure of central tendency alone is not sufficient to describe a frequency distribution. In addition to it we should have a measure of scatterness of observations. The scatterness or variation of observations from their average are called the dispersion. There are different measures of dispersion like the range, the quartile deviation, the mean deviation and the standard deviation.

Characteristics of a good measure of dispersion

An ideal measure of dispersion is expected to possess the following properties

1. It should be rigidly defined
2. It should be based on all the items.
3. It should not be unduly affected by extreme items.
4. It should lend itself for algebraic manipulation.
5. It should be simple to understand and easy to calculate

Range

This is the simplest possible measure of dispersion and is defined as the difference between the largest and smallest values of the variable.

- In symbols, Range = $L - S$.
- Where L = Largest value.
- S = Smallest value.

In individual observations and discrete series, L and S are easily identified.

In continuous series, the following two methods are followed.

Method 1

L = Upper boundary of the highest class

S = Lower boundary of the lowest class.

Method 2

L = Mid value of the highest class.

S = Mid value of the lowest class.

Example 1

The yields (kg per plot) of a cotton variety from five plots are 8, 9, 8, 10 and 11. Find the range

Solution

L=11, S = 8.

Range = $L - S = 11 - 8 = 3$

Example 2

Calculate range from the following distribution.

Size:	60-63	63-66	66-69	69-72	72-75
Number:	5	18	42	27	8

Solution

L = Upper boundary of the highest class = 75

S = Lower boundary of the lowest class = 60

Range = L – S = 75 – 60 = 15

Merits and Demerits of Range

Merits

1. It is simple to understand.
2. It is easy to calculate.
3. In certain types of problems like quality control, weather forecasts, share price analysis, etc.,
range is most widely used.

Demerits

1. It is very much affected by the extreme items.
2. It is based on only two extreme observations.
3. It cannot be calculated from open-end class intervals.
4. It is not suitable for mathematical treatment.
5. It is a very rarely used measure.

Standard Deviation

It is defined as the positive square-root of the arithmetic mean of the Square of the deviations of the given observation from their arithmetic mean.

The standard deviation is denoted by s in case of sample and Greek letter σ (sigma) in case of population.

The formula for calculating standard deviation is as follows

$$s = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n - 1}} \quad \text{for raw data}$$

And for grouped data the formulas are

$$s = \sqrt{\frac{\sum fx^2}{N} - \left(\frac{\sum fx}{N}\right)^2} \text{ for discrete data}$$

$$s = C \times \sqrt{\frac{\sum fd^2}{N} - \left(\frac{\sum fd}{N}\right)^2} \text{ for continuous data}$$

Where $d = \frac{x - A}{C}$

C = class interval

Example 3

Raw Data

The weights of 5 ear-heads of sorghum are 100, 102, 118, 124, 126 gms. Find the standard deviation.

Solution

x	x ²
100	10000
102	10404
118	13924
124	15376
126	15876
570	65580

$$\begin{aligned} \text{Standard deviation } s &= \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}} \\ &= \sqrt{\frac{65580 - \frac{(570)^2}{5}}{5-1}} = \sqrt{150} = 12.25 \text{ gms} \end{aligned}$$

Example 4

Discrete distribution

The frequency distributions of seed yield of 50 sesamum plants are given below. Find the standard deviation.

Seed yield in gms (x)	3	4	5	6	7
Frequency (f)	4	6	15	15	10

Solution

Seed yield in gms (x)	f	fx	fx ²
3	4	12	36
4	6	24	96
5	15	75	375
6	15	90	540
7	10	70	490
Total	50	271	1537

Here $n = 50$

$$\begin{aligned}
 \text{Standard deviation } s &= \sqrt{\frac{\sum fx^2}{n} - \left(\frac{\sum fx}{n}\right)^2} \\
 &= \sqrt{\frac{1537}{50} - \left(\frac{271}{50}\right)^2} \\
 &= \sqrt{30.74 - 29.3764} \\
 &= 1.1677 \text{ gms}
 \end{aligned}$$

Example 5
Continuous distribution

The Frequency distributions of seed yield of 50 sesamum plants are given below. Find the standard deviation.

Seed yield in gms (x)	2.5-3.5	3.5-4.5	4.5-5.5	5.5-6.5	6.5-7.5
No. of plants (f)	4	6	15	15	10

Solution

Seed yield in gms (x)	No. of Plants f	Mid x	$d = \frac{x - A}{C}$	df	d ² f
2.5-3.5	4	3	-2	-8	16
3.5-4.5	6	4	-1	-6	6
4.5-5.5	15	5	0	0	0
5.5-6.5	15	6	1	15	15

6.5-7.5	10	7	2	20	40
Total	50	25	0	21	77

A=Assumed mean = 5

n=50, C=1

$$\begin{aligned}
 s &= C \times \sqrt{\frac{\sum fd^2}{n} - \left(\frac{\sum fd}{n}\right)^2} \\
 &= 1 \times \sqrt{\frac{77}{50} - \left(\frac{21}{50}\right)^2} \\
 &= \sqrt{1.54 - 0.1764} \\
 &= \sqrt{1.3636} = 1.1677
 \end{aligned}$$

Merits and Demerits of Standard Deviation

Merits

1. It is rigidly defined and its value is always definite and based on all the observations and the actual signs of deviations are used.
2. As it is based on arithmetic mean, it has all the merits of arithmetic mean.
3. It is the most important and widely used measure of dispersion.
4. It is possible for further algebraic treatment.
5. It is less affected by the fluctuations of sampling and hence stable.
6. It is the basis for measuring the coefficient of correlation and sampling.

Demerits

1. It is not easy to understand and it is difficult to calculate.
2. It gives more weight to extreme values because the values are squared up.
3. As it is an absolute measure of variability, it cannot be used for the purpose of comparison.

Variance

The square of the standard deviation is called variance

(i.e.) variance = (SD)².

Coefficient of Variation

The Standard deviation is an absolute measure of dispersion. It is expressed in terms of units in which the original figures are collected and stated. The standard deviation of heights of plants cannot be compared with the standard deviation of weights of the grains, as both are expressed in different units, i.e heights in centimeter and weights in kilograms. Therefore the standard deviation must be converted into a relative measure of dispersion for the purpose of comparison. The relative measure is known as the coefficient of variation. The coefficient of variation is obtained by dividing the standard deviation by the mean and expressed in percentage. Symbolically, Coefficient of

$$\text{variation (C.V)} = \frac{SD}{\text{mean}} \times 100$$

If we want to compare the variability of two or more series, we can use C.V. The series or groups of data for which the C.V. is greater indicate that the group is more variable, less stable, less uniform, less consistent or less homogeneous. If the C.V. is less, it indicates that the group is less variable or more stable or more uniform or more consistent or more homogeneous.

Example 6

Consider the measurement on yield and plant height of a paddy variety. The mean and standard deviation for yield are 50 kg and 10 kg respectively. The mean and standard deviation for plant height are 55 am and 5 cm respectively.

Here the measurements for yield and plant height are in different units. Hence the variabilities can be compared only by using coefficient of variation.

$$\text{For yield, CV} = \frac{10}{50} \times 100 = 20\%$$

$$\text{For plant height, CV} = \frac{5}{55} \times 100 = 9.1\%$$

The yield is subject to more variation than the plant height.

Questions

1. Which measure is affected most by the presence of extreme values.

- a) Range
- b) Standard Deviation
- b) Quartile Deviation
- d) Mean deviation

Ans: Standard Deviation

2. Variance is square of _____

- a) Range
- b) Standard Deviation
- c) Quartile Deviation
- d) Mean deviation

Ans: Standard Deviation

3. If the CV of variety I is 30% and variety II is 25% then Variety II is more consistent.

Ans: True

4. For the set of data 5, 5, 5,5,5,5 the Standard deviation value is zero.

Ans: True

5. The absolute measures of dispersion will have the original units.

Ans: True

6. The mean deviation value for a set of data can take even negative value.

Ans: False

7. Define dispersion.

8. Define C.V. What are its uses?

9. What are the differences between absolute measure and relative measure of dispersion?

10. How to calculate the standard deviation for raw and grouped data?

Lecture.6

Probability – Basic concepts-trial- event-equally likely- mutually exclusive – independent event, additive and multiplicative laws. Theoretical distributions- discrete and continuous distributions, Binomial distributions-properties

Probability

The concept of probability is difficult to define in precise terms. In ordinary language, the word probable means likely (or) chance. Generally the word, probability, is used to denote the happening of a certain event, and the likelihood of the occurrence of that event, based on past experiences. By looking at the clear sky, one will say that there will not be any rain today. On the other hand, by looking at the cloudy sky or overcast sky, one will say that there will be rain today. In the earlier sentence, we aim that there will not be rain and in the latter we expect rain. On the other hand a mathematician says that the probability of rain is '0' in the first case and that the probability of rain is '1' in the second case. In between 0 and 1, there are fractions denoting the chance of the event occurring. In ordinary language, the word probability means uncertainty about happenings. In Mathematics and Statistics, a numerical measure of uncertainty is provided by the important branch of statistics – called theory of probability. Thus we can say, that the theory of probability describes certainty by 1 (one), impossibility by 0 (zero) and uncertainties by the co-efficient which lies between 0 and 1.

Trial and Event An experiment which, though repeated under essentially identical (or) same conditions does not give unique results but may result in any one of the several possible outcomes. Performing an experiment is known as a trial and the outcomes of the experiment are known as events.

Example 1: Seed germination – either germinates or does not germinates are events.

2. In a lot of 5 seeds none may germinate (0), 1 or 2 or 3 or 4 or all 5 may germinate.

Sample space (S)

A set of all possible outcomes from an experiment is called sample space. For example, a set of five seeds are sown in a plot, none may germinate, 1, 2, 3, 4 or all five may germinate. i.e the possible outcomes are {0, 1, 2, 3, 4, 5}. The set of numbers is called a sample space. Each possible outcome (or) element in a sample space is called sample point.

Exhaustive Events

The total number of possible outcomes in any trial is known as exhaustive events (or) exhaustive cases.

Example

1. When pesticide is applied a pest may survive or die. There are two exhaustive cases namely (survival, death)
2. In throwing of a die, there are six exhaustive cases, since anyone of the 6 faces
1, 2, 3, 4, 5, 6 may come uppermost.
3. In drawing 2 cards from a pack of cards the exhaustive number of cases is $52C_2$, since 2 cards can be drawn out of 52 cards in $52C_2$ ways

Trial	Random Experiment	Total number of trials	Sample Space
(1)	One pest is exposed to pesticide	$2^1=2$	{S,D}
(2)	Two pests are exposed to pesticide	$2^2=4$	{SS, SD, DS, DD}
(3)	Three pests are exposed to pesticide	$2^3=8$	{SSS, SSD, SDS, DSS, SDD, DSD, DDS, DDD}
(4)	One set of three seeds	$4^1= 4$	{0,1,2,3}
(5)	Two sets of three seeds	$4^2=16$	{0,1},{0,2},{0,3} etc

Favourable Events

The number of cases favourable to an event in a trial is the number of outcomes which entail the happening of the event.

Example

1. When a seed is sown if we observe non germination of a seed, it is a favourable event. If we are interested in germination of the seed then germination is the favourable event.

Mutually Exclusive Events

Events are said to be mutually exclusive (or) incompatible if the happening of any one of the events excludes (or) precludes the happening of all the others i.e.) if no two or more of the events can happen simultaneously in the same trial. (i.e.) The joint occurrence is not possible.

Example

1. In observation of seed germination the seed may either germinate or it will not germinate. Germination and non germination are mutually exclusive events.

Equally Likely Events

Outcomes of a trial are said to be equally likely if taking in to consideration all the relevant evidences, there is no reason to expect one in preference to the others. (i.e.) Two or more events are said to be equally likely if each one of them has an **equal chance of occurring**.

Independent Events

Several events are said to be independent if the happening of an event is not affected by the happening of one or more events.

Example

1. When two seeds are sown in a pot, one seed germinates. It would not affect the germination or non germination of the second seed. One event does not affect the other event.

Dependent Events

If the happening of one event is affected by the happening of one or more events, then the events are called dependent events.

Example

If we draw a card from a pack of well shuffled cards, if the first card drawn is not replaced then the second draw is dependent on the first draw.

Note: In the case of independent (or) dependent events, the joint occurrence is possible.

Definition of Probability

Mathematical (or) Classical (or) a-priori Probability

If an experiment results in 'n' exhaustive cases which are mutually exclusive and equally likely cases out of which 'm' events are favourable to the happening of an event 'A', then the probability 'p' of happening of 'A' is given by

$$p = P(A) = \frac{\text{Favourable number of cases}}{\text{Exhaustive number of cases}} = \frac{m}{n}$$

Note

1. If $m = 0 \Rightarrow P(A) = 0$, then 'A' is called an impossible event. (i.e.) also by $P(\phi) = 0$.
2. If $m = n \Rightarrow P(A) = 1$, then 'A' is called assure (or) certain event.
3. The probability is a non-negative real number and cannot exceed unity (i.e.) lies between 0 to 1.
4. The probability of non-happening of the event 'A' (i.e.) $P(\bar{A})$. It is denoted by 'q'.

$$P(\bar{A}) = \frac{n - m}{n} = 1 - \frac{m}{n} = 1 - P(A)$$

$$\Rightarrow q = 1 - p$$

$$\Rightarrow p + q = 1$$

$$\text{(or) } P(A) + P(\bar{A}) = 1.$$

Statistical (or) Empirical Probability (or) a-posteriori Probability

If an experiment is repeated a number (n) of times, an event ‘A’ happens ‘m’ times then the statistical probability of ‘A’ is given by

$$p = P(A) = \lim_{n \rightarrow \infty} \frac{m}{n}$$

Axioms for Probability

1. The probability of an event ranges from 0 to 1. If the event cannot take place its probability shall be ‘0’ if it certain, its probability shall be ‘1’.

Let E_1, E_2, \dots, E_n be any events, then $P(E_i) \geq 0$.

2. The probability of the entire sample space is ‘1’. (i.e.) $P(S) = 1$.

$$\text{Total Probability, } \sum_{i=1}^n P(E_i) = 1$$

3. If A and B are mutually exclusive (or) disjoint events then the probability of occurrence of either A (or) B denoted by $P(A \cup B)$ shall be given by

$$P(A \cup B) = P(A) + P(B)$$

$$P(E_1 \cup E_2 \cup \dots \cup E_n) = P(E_1) + P(E_2) + \dots + P(E_n)$$

If E_1, E_2, \dots, E_n are mutually exclusive events.

Example 1: Two dice are tossed. What is the probability of getting (i) Sum 6 (ii) Sum 9?

Solution

When 2 dice are tossed. The exhaustive number of cases is 36 ways.

(i) Sum 6 = {(1, 5), (2, 4), (3, 3), (4, 2), (5, 1)}

∴ Favourable number of cases = 5

$$P(\text{Sum } 6) = \frac{5}{36}$$

(ii) Sum 9 = {(3, 6), (4, 5), (5, 4), (6, 3)}

∴ Favourable number of cases = 4

$$P(\text{Sum } 9) = \frac{4}{36} = \frac{1}{9}$$

Example 2: A card is drawn from a pack of cards. What is a probability of getting (i) a king (ii) a spade (iii) a red card (iv) a numbered card?

Solution

There are 52 cards in a pack.

One can be selected in $52C_1$ ways.

\therefore Exhaustive number of cases is $= 52C_1 = 52$.

(i) A king

There are 4 kings in a pack.

One king can be selected in $4C_1$ ways.

\therefore Favourable number of cases is $= 4C_1 = 4$

Hence the probability of getting a king $= \frac{4}{52}$

(ii) A spade

There are 13 spades in a pack.

One spade can be selected in $13C_1$ ways.

\therefore Favourable number of cases is $= 13C_1 = 13$

Hence the probability of getting a spade $= \frac{13}{52}$

(iii) A red card

There are 26 red cards in a pack.

One red card can be selected in $26C_1$ ways.

\therefore Favourable number of cases is $= 26C_1 = 26$

Hence the probability of getting a red card $= \frac{26}{52}$

(iv) A numbered card

There are 36 numbered cards in a pack.

One numbered card can be selected in $36C_1$ ways.

\therefore Favourable number of cases is $= 36C_1 = 36$

Hence the probability of getting a numbered card $= \frac{36}{52}$

Example 3: What is the probability of getting 53 Sundays when a leap year selected at random?

Solution

A leap year consists of 366 days.

This has 52 full weeks and 2 days remained.

The remaining 2 days have the following possibilities.

- (i) Sun, Mon (ii) Mon, Tues (iii) Tues, Wed (iv) Wed, Thurs (v) Thurs, Fri (vi) Fri, Sat (vii) Sat, Sun.

In order that a lap year selected at random should contain 53 Sundays, one of the 2 over days must be Sunday.

∴ Exhaustive number of cases is = 7

∴ Favourable number of cases is = 2

∴ Required Probability is = $\frac{2}{7}$

Conditional Probability

Two events A and B are said to be dependent, when B can occur only when A is known to have occurred (or vice versa). The probability attached to such an event is called the conditional probability and is denoted by P (A/B) (read it as: A given B) or, in other words, probability of A given that B has occurred.

$$P(A/B) = \frac{P(A \cap B)}{P(B)} = \frac{P(AB)}{P(B)}$$

If two events A and B are **dependent**, then the conditional probability of B given A is,

$$P(B/A) = \frac{P(A \cap B)}{P(A)} = \frac{P(AB)}{P(A)}$$

Theorems of Probability

There are two important theorems of probability namely,

1. The addition theorem on probability
2. The multiplication theorem on probability.

I. Addition Theorem on Probability

(i) Let A and B be any two events which are **not mutually exclusive**

$$P(A \text{ or } B) = P(A \cup B) = P(A + B) = P(A) + P(B) - P(A \cap B) \quad (\text{or})$$

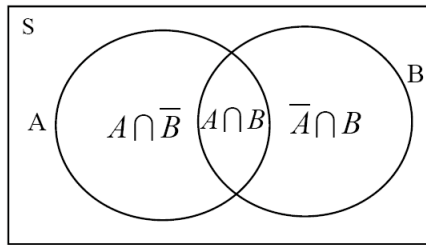
$$= P(A) + P(B) - P(AB)$$

Proof

Let us take a random experiment with a sample space S of N sample points.

Then by the definition of probability ,

$$P(A \cup B) = \frac{n(A \cup B)}{n(S)} = \frac{n(A \cup B)}{N}$$



From the diagram, using the axiom for the mutually exclusive events, we write

$$P(A \cup B) = \frac{n(A) + n(\bar{A} \cap B)}{N}$$

Adding and subtracting $n(A \cap B)$ in the numerator,

$$= \frac{n(A) + n(\bar{A} \cap B) + n(A \cap B) - n(A \cap B)}{N}$$

$$= \frac{n(A) + n(B) - n(A \cap B)}{N}$$

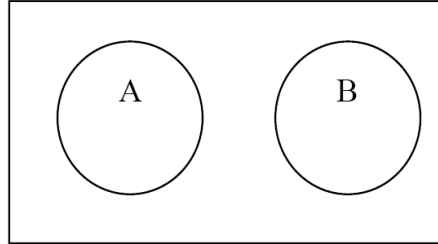
$$= \frac{n(A)}{N} + \frac{n(B)}{N} - \frac{n(A \cap B)}{N}$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

(ii) Let A and B be any two events which are **mutually exclusive**

$$P(A \text{ or } B) = P(A \cup B) = P(A + B) = P(A) + P(B)$$

Proof



We know that, $n(A \cup B) = n(A) + n(B)$

$$\begin{aligned} P(A \cup B) &= \frac{n(A \cup B)}{n} \\ &= \frac{n(A) + n(B)}{n} \\ &= \frac{n(A)}{n} + \frac{n(B)}{n} \end{aligned}$$

$$P(A \cup B) = P(A) + P(B)$$

Note

(i) In the case of 3 events, **(not mutually exclusive events)**

$$\begin{aligned} P(A \text{ or } B \text{ or } C) &= P(A \cup B \cup C) = P(A + B + C) \\ &= P(A) + P(B) + P(C) - P(A \cap B) - P(B \cap C) - P(A \cap C) + P(A \cap B \cap C) \end{aligned}$$

(ii) In the case of 3 events, **(mutually exclusive events)**

$$P(A \text{ or } B \text{ or } C) = P(A \cup B \cup C) = P(A + B + C) = P(A) + P(B) + P(C)$$

Example

Using the additive law of probability we can find the probability that in one roll of a die, we will obtain either a one-spot or a six-spot. The probability of obtaining a one-spot is $1/6$. The probability of obtaining a six-spot is also $1/6$. The probability of rolling a die and getting a side that has both a one-spot with a six-spot is 0. There is no side on a

die that has both these events. So substituting these values into the equation gives the following result:

$$\frac{1}{6} + \frac{1}{6} - 0 = \frac{2}{6} = \frac{1}{3} = 0.3333$$

Finding the probability of drawing a 4 of hearts or a 6 of any suit using the additive law of probability would give the following:

$$\frac{1}{52} + \frac{4}{52} - 0 = \frac{5}{52} = 0.0962$$

There is only a single 4 of hearts, there are 4 sixes in the deck and there isn't a single card that is both the 4 of hearts and a six of any suit.

Now using the additive law of probability, you can find the probability of drawing either a king or any club from a deck of shuffled cards. The equation would be completed like this:

$$\frac{4}{52} + \frac{13}{52} - \frac{1}{52} = \frac{16}{52} = 0.3077$$

There are 4 kings, 13 clubs, and obviously one card is both a king and a club. We don't want to count that card twice, so you must subtract one of its occurrences away to obtain the result.

II. Multiplication Theorem on Probability

(i) If A and B be any two events which are **not independent**, then (i.e.) **dependent**.

$$P(A \text{ and } B) = P(A \cap B) = P(AB) = P(A) \cdot P(B/A) \longrightarrow \text{(I)}$$

$$= P(B) \cdot P(A/B) \longrightarrow \text{(II)}$$

Where P (B/A) and P (A/B) are the conditional probability of B given A and A given B respectively.

Proof

Let n is the total number of events

$n(A)$ is the number of events in A

$n(B)$ is the number of events in B

$n(A \cup B)$ is the number of events in $(A \cup B)$

$n(A \cap B)$ is the number of events in $(A \cap B)$

$$\begin{aligned} P(A \cap B) &= \frac{n(A \cap B)}{n} \\ &= \frac{n(A \cap B)}{n} \times \frac{n(A)}{n(A)} \\ &= \frac{n(A)}{n} \times \frac{n(A \cap B)}{n(A)} \end{aligned}$$

$$P(A \cap B) = P(A) \cdot P(B/A) \longrightarrow (I)$$

$$\begin{aligned} P(A \cap B) &= \frac{n(A \cap B)}{n} \\ &= \frac{n(A \cap B)}{n} \times \frac{n(B)}{n(B)} \\ &= \frac{n(B)}{n} \times \frac{n(A \cap B)}{n(B)} \end{aligned}$$

$$P(A \cap B) = P(B) \cdot P(A/B) \longrightarrow (II)$$

(ii) If A and B be any two events which are **independent**, then,

$$P(B/A) = P(B) \text{ and } P(A/B) = P(A)$$

$$P(A \text{ and } B) = P(A \cap B) = P(AB) = P(A) \cdot P(B)$$

Note

(i) In the case of 3 events, (**dependent**)

$$P(A \cap B \cap C) = P(A) \cdot P(B/A) \cdot P(C/AB)$$

(ii) In the case of 3 events, **(independent)**

$$P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C)$$

Example

So in finding the probability of drawing a 4 and then a 7 from a well shuffled deck of cards, this law would state that we need to multiply those separate probabilities together. Completing the equation above gives:

$$p(4 \text{ and } 7) = \frac{4}{52} \times \frac{4}{52} = \frac{16}{2704} = 0.0059$$

Given a well shuffled deck of cards, what is the probability of drawing a Jack of Hearts, Queen of Hearts, King of Hearts, Ace of Hearts, and 10 of Hearts?

$$p(10, J, Q, K, A \text{ of hearts}) = \frac{1}{52} \times \frac{1}{52} \times \frac{1}{52} \times \frac{1}{52} \times \frac{1}{52} = 0.0000000026$$

In any case, given a well shuffled deck of cards, obtaining this assortment of cards, drawing one at a time and returning it to the deck would be highly unlikely (it has an exceedingly low probability).

Questions

1. Probability is expressed as

- (a) Ratio
- (b) percentage
- (c) Proportion
- (d) all the above

Ans: all the above

2. Probability can take values from

- (a) $-\infty$ to $+\infty$
- (b) $-\infty$ to 1
- (c) 0 to +1
- (d) -1 to +1

Ans: 0 to +1

3. The probability of a sure event is One.

Ans: True

4. If A and B are mutually exclusive events, then $P(A \cup B) = \underline{P(A) + P(B)}$

5. An integer is chosen from 1 to 20. The probability that the number is divisible by 4 is $\frac{1}{4}$.

Ans: True

6. Mean of the Binomial Distribution is npq .

Ans: False

7. Define an independent event.

8. What is conditional probability?

9. State the addition and multiplication laws.

10. State the properties of Binomial distribution.

Lecture.7

Poisson Distributions - properties, Normal Distributions- properties

Theoretical Distributions

Theoretical distributions are

- | | | |
|--------------------------|---|-------------------------|
| 1. Binomial distribution | } | Discrete distribution |
| 2. Poisson distribution | | |
| 3. Normal distribution | → | Continuous distribution |

Discrete Probability distribution

Bernoulli distribution

A random variable x takes two values 0 and 1, with probabilities q and p ie., $p(x=1) = p$ and $p(x=0)=q$, $q=1-p$ is called a Bernoulli variate and is said to be Bernoulli distribution where p and q are probability of success and failure. It was given by Swiss mathematician James Bernoulli (1654-1705)

Example

- Tossing a coin(head or tail)
- Germination of seed(germinate or not)

Binomial distribution

Binomial distribution was discovered by James Bernoulli (1654-1705). Let a random experiment be performed repeatedly and the occurrence of an event in a trial be called as success and its non-occurrence is failure. Consider a set of n independent trails (n being finite), in which the probability p of success in any trail is constant for each trial. Then $q=1-p$ is the probability of failure in any trail.

The probability of x success and consequently n-x failures in n independent trails. But x successes in n trails can occur in ${}^n C_x$ ways. Probability for each of these ways is $p^x q^{n-x}$.

$$\begin{aligned} P(\text{sss...ff...fsf...f}) &= p(s)p(s)\dots p(f)p(f)\dots \\ &= p,p\dots q,q\dots \\ &= (p,p\dots p)(q,q\dots q) \\ &\quad (x \text{ times}) (n-x \text{ times}) \end{aligned}$$

Hence the probability of x success in n trials is given by

$${}^n C_x p^x q^{n-x}$$

Definition

A random variable x is said to follow binomial distribution if it assumes non-negative values and its probability mass function is given by

$$P(X=x) = p(x) = \begin{cases} {}^n C_x p^x q^{n-x}, & x=0,1,2\dots n \\ q=1-p \\ 0, & \text{otherwise} \end{cases}$$

The two independent constants n and p in the distribution are known as the parameters of the distribution.

Condition for Binomial distribution

We get the binomial distribution under the following experimentation conditions

1. The number of trial n is finite
2. The trials are independent of each other.

3. The probability of success p is constant for each trial.
4. Each trial must result in a success or failure.
5. The events are discrete events.

Properties

1. If p and q are equal, the given binomial distribution will be symmetrical. If p and q are not equal, the distribution will be skewed distribution.
2. Mean = $E(x) = np$
3. Variance = $V(x) = npq$ (mean > variance)

Application

1. Quality control measures and sampling process in industries to classify items as defectives or non-defective.
2. Medical applications such as success or failure, cure or no-cure.

Example 1

Eight coins are tossed simultaneously. Find the probability of getting atleast six heads.

Solution

Here number of trials, $n = 8$, p denotes the probability of getting a head.

$$\therefore p = \frac{1}{2} \text{ and } q = \frac{1}{2}$$

If the random variable X denotes the number of heads, then the probability of a success in n trials is given by

$$P(X = x) = {}^n C_x p^x q^{n-x}, \quad x = 0, 1, 2, \dots, n$$

$$= 8C_x \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{8-x} = 8C_x \left(\frac{1}{2}\right)^8$$

$$= \frac{1}{2^8} 8C_x$$

Probability of getting atleast six heads is given by

$$\begin{aligned} P(x \geq 6) &= P(x = 6) + P(x = 7) + P(x = 8) \\ &= \frac{1}{2^8} 8C_6 + \frac{1}{2^8} 8C_7 + \frac{1}{2^8} 8C_8 \\ &= \frac{1}{2^8} [8C_6 + 8C_7 + 8C_8] \\ &= \frac{1}{2^8} [28 + 8 + 1] = \frac{37}{256} \end{aligned}$$

Example 2 Ten coins are tossed simultaneously. Find the probability of getting (i) atleast seven heads (ii) exactly seven heads (iii) atleast seven heads

Solution

$$p = \text{Probability of getting a head} = \frac{1}{2}$$

$$q = \text{Probability of not getting a head} = \frac{1}{2}$$

The probability of getting x heads throwing 10 coins simultaneously is given by

$$\begin{aligned} P(X = x) &= {}^n C_x p^x q^{n-x}, \quad x = 0, 1, 2, \dots, n \\ &= {}^{10} C_x \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{10-x} = {}^{10} C_x \left(\frac{1}{2}\right)^{10} \\ &= \frac{1}{2^{10}} {}^{10} C_x \end{aligned}$$

i) Probability of getting atleast seven heads

$$\begin{aligned} P(x \geq 7) &= P(x = 7) + P(x = 8) + P(x = 9) + P(x = 10) \\ &= \frac{1}{2^{10}} [{}^{10} C_7 + {}^{10} C_8 + {}^{10} C_9 + {}^{10} C_{10}] \\ &= \frac{1}{1024} [120 + 45 + 10 + 1] = \frac{176}{1024} \end{aligned}$$

ii) Probability of getting exactly 7 heads

$$P(x = 7) = \frac{1}{2^{10}} {}^{10} C_7 = \frac{1}{2^{10}} (120) = \frac{120}{1024}$$

iii) Probability of getting atleast 7 heads

$$P(x \leq 7) = 1 - P(x > 7)$$

$$\begin{aligned}
 &= 1 \text{ symbol } \{P(x = 8) + P(x = 9) + P(x = 10)\} \\
 &= 1 - \frac{1}{2^{10}} [10C_8 + 10C_9 + 10C_{10}] \\
 &= 1 - \frac{1}{2^{10}} [45 + 10 + 1] \\
 &= 1 - \frac{56}{1024} \\
 &= \frac{968}{1024}
 \end{aligned}$$

Example 3: 20 wrist watches in a box of 100 are defective. If 10 watches are selected at random, find the probability that (i) 10 are defective (ii) 10 are good (iii) at least one watch is defective (iv) at most 3 are defective.

Solution

20 out of 100 wrist watches are defective

Probability of defective wrist watch, $p = \frac{20}{100} = \frac{1}{5}$

$q = 1 - p = \frac{4}{5}$

Since 10 watches are selected at random, $n = 10$

$P(X = x) = nC_x p^x q^{n-x}$, $x = 0, 1, 2, \dots, 10$

$$= 10C_x \left(\frac{1}{5}\right)^x \left(\frac{4}{5}\right)^{10-x}$$

i) Probability of selecting 10 defective watches

$$P(x = 10) = 10C_{10} \left(\frac{1}{5}\right)^{10} \left(\frac{4}{5}\right)^0 = 1 \cdot \frac{1}{5^{10}} \cdot 1 = \frac{1}{5^{10}}$$

ii) Probability of selecting 10 good watches (i.e. no defective)

$$\begin{aligned}
 P(x = 0) &= 10C_0 \left(\frac{1}{5}\right)^0 \left(\frac{4}{5}\right)^{10} \\
 &= 1 \cdot 1 \left(\frac{4}{5}\right)^{10} = \left(\frac{4}{5}\right)^{10}
 \end{aligned}$$

iii) Probability of selecting at least one defective watch

$$P(x \geq 1) = 1 - P(x < 1)$$

$$= 1 - P(x = 0)$$

$$= 1 - {}^{10}C_0 \left(\frac{1}{5}\right)^0 \left(\frac{4}{5}\right)^{10}$$

$$= 1 - \left(\frac{4}{5}\right)^{10}$$

iv) Probability of selecting at most 3 defective watches

$$P(x \leq 3) = P(x = 0) + P(x = 1) + P(x = 2) + P(x = 3)$$

$$= {}^{10}C_0 \left(\frac{1}{5}\right)^0 \left(\frac{4}{5}\right)^{10} + {}^{10}C_1 \left(\frac{1}{5}\right)^1 \left(\frac{4}{5}\right)^9 + {}^{10}C_2 \left(\frac{1}{5}\right)^2 \left(\frac{4}{5}\right)^8 + {}^{10}C_3 \left(\frac{1}{5}\right)^3 \left(\frac{4}{5}\right)^7$$

$$= 1.1 \left(\frac{4}{5}\right)^{10} + 10 \left(\frac{1}{5}\right)^1 \left(\frac{4}{5}\right)^9 + \frac{10 \cdot 9}{1 \cdot 2} \left(\frac{1}{5}\right)^2 \left(\frac{4}{5}\right)^8 + \frac{10 \cdot 9 \cdot 8}{1 \cdot 2 \cdot 3} \left(\frac{1}{5}\right)^3 \left(\frac{4}{5}\right)^7$$

$$= 1. (0.107) + 10 (0.026) + 45 (0.0062) + 120 (0.0016)$$

$$= 0.859 \text{ (approx)}$$

Poisson distribution

The Poisson distribution, named after Simeon Denis Poisson (1781-1840). Poisson distribution is a discrete distribution. It describes random events that occurs rarely over a unit of time or space.

It differs from the binomial distribution in the sense that we count the number of success and number of failures, while in Poisson distribution, the average number of success in given unit of time or space.

Definition

The probability that exactly x events will occur in a given time is as follows

$$P(x) = \frac{e^{-\lambda} \lambda^x}{x!}, x=0,1,2,\dots$$

called as probability mass function of Poisson distribution.

where λ is the average number of occurrences per unit of time

$$\lambda = np$$

Condition for Poisson distribution

Poisson distribution is the limiting case of binomial distribution under the following assumptions.

1. The number of trials n should be indefinitely large i.e., $n \rightarrow \infty$
2. The probability of success p for each trial is indefinitely small.
3. $np = \lambda$, should be finite where λ is constant.

Properties

1. Poisson distribution is defined by single parameter λ .
2. Mean = λ
3. Variance = λ . Mean and Variance are equal.

Application

1. It is used in quality control statistics to count the number of defects of an item.
2. In biology, to count the number of bacteria.
3. In determining the number of deaths in a district in a given period, by rare disease.
4. The number of error per page in typed material.
5. The number of plants infected with a particular disease in a plot of field.
6. Number of weeds in particular species in different plots of a field.

Example 4: Suppose on an average 1 house in 1000 in a certain district has a fire during a year. If there are 2000 houses in that district, what is the probability that exactly 5 houses will have a fire during the year? [given that $e^{-2} = 0.13534$]

Solution:

$$\text{Mean, } \bar{x} = np, n = 2000 \text{ and } p = \frac{1}{1000}$$

$$= 2000 \times \frac{1}{1000}$$

$$\lambda = 2$$

The Poisson distribution is

$$P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$$P(X = 5) = \frac{e^{-2} 2^5}{5!}$$

$$= \frac{(0.13534) \times 32}{120}$$

$$= \mathbf{0.036}$$

Example 5

If 2% of electric bulbs manufactured by a certain company are defective. Find the probability that in a sample of 200 bulbs i) less than 2 bulbs ii) more than 3 bulbs are defective. [e-4 = 0.0183]

Solution

$$\text{The probability of a defective bulb} = p = \frac{2}{100} = 0.02$$

Given that n = 200 since p is small and n is large

We use the Poisson distribution

$$\text{mean, } m = np = 200 \times 0.02 = 4$$

$$\text{Now, Poisson Probability function, } P(X = x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

i) Probability of less than 2 bulbs are defective

$$\begin{aligned}
 &= P(X < 2) \\
 &= P(x = 0) + P(x = 1) \\
 &= e^{-4} + e^{-4} (4) \\
 &= e^{-4} (1 + 4) = 0.0183 \times 5 \\
 &= 0.0915
 \end{aligned}$$

ii) Probability of getting more than 3 defective bulbs

$$\begin{aligned}
 P(x > 3) &= 1 - P(x \leq 3) \\
 &= 1 - \{P(x = 0) + P(x = 1) + P(x = 2) + P(x = 3)\} \\
 &= 1 - e^{-4} \left\{ 1 + 4 + \frac{4^2}{2!} + \frac{4^3}{3!} \right\} \\
 &= 1 - \{0.0183 \times (1 + 4 + 8 + 10.67)\} \\
 &= 0.567
 \end{aligned}$$

Normal distribution

Continuous Probability distribution is normal distribution. It is also known as error law or Normal law or Laplacian law or Gaussian distribution. Many of the sampling distribution like student-t, f distribution and χ^2 distribution.

Definition

A continuous random variable x is said to be a normal distribution with parameters μ and σ^2 , if the density function is given by the probability law

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}; -\infty < x < \infty, -\infty < \mu < \infty, \sigma > 0$$

Note

The mean m and standard deviation s are called the parameters of Normal distribution. The normal distribution is expressed by $X \sim N(m, s^2)$

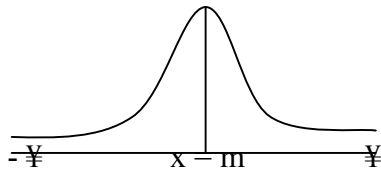
Condition of Normal Distribution

i) Normal distribution is a limiting form of the binomial distribution under the following conditions.

- a) n , the number of trials is indefinitely large i.e., $n \rightarrow \infty$ and
 - b) Neither p nor q is very small.
- ii) Normal distribution can also be obtained as a limiting form of Poisson distribution with parameter $m \rightarrow \infty$
- iii) Constants of normal distribution are mean = m , variation = s^2 , Standard deviation = s .

Normal probability curve

The curve representing the normal distribution is called the normal probability curve. The curve is symmetrical about the mean (m), bell-shaped and the two tails on the right and left sides of the mean extends to the infinity. The shape of the curve is shown in the following figure.



Properties of normal distribution

1. The normal curve is bell shaped and is symmetric at $x = m$.
2. Mean, median, and mode of the distribution are coincide
i.e., Mean = Median = Mode = m
3. It has only one mode at $x = m$ (i.e., unimodal)
4. The points of inflection are at $x = m \pm s$
5. The maximum ordinate occurs at $x = m$ and its value is $= \frac{1}{\sigma\sqrt{2\pi}}$
6. Area Property $P(m - s < x < m + s) = 0.6826$
 $P(m - 2s < x < m + 2s) = 0.9544$
 $P(m - 3s < x < m + 3s) = 0.9973$

Standard Normal distribution

Let X be random variable which follows normal distribution with mean μ and variance σ^2 . The standard normal variate is defined as $Z = \frac{X - \mu}{\sigma}$ which follows standard normal distribution with mean 0 and standard deviation 1 i.e., $Z \sim N(0,1)$. The

standard normal distribution is given by $\phi(Z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2}$; $-\infty < z < \infty$

The advantage of the above function is that it doesn't contain any parameter. This enables us to compute the area under the normal probability curve.

Note

Property of $\phi(Z)$

1. $\phi(-Z) = 1 - \phi(Z)$
2. $P(a \leq Z \leq b) = \phi(b) - \phi(a)$

Example 6: In a normal distribution whose mean is 12 and standard deviation is 2. Find the probability for the interval from $x = 9.6$ to $x = 13.8$

Solution

Given that $Z \sim N(12, 4)$

$$\begin{aligned} P(9.6 \leq Z \leq 13.8) &= P\left(\frac{9.6-12}{2} \leq Z \leq \frac{13.8-12}{2}\right) \\ &= P(-1.2 \leq Z \leq 0) + P(0 \leq Z \leq 0.9) \\ &= P(0 \leq Z \leq 1.2) + P(0 \leq Z \leq 0.9) \quad [\text{by using symmetric property}] \\ &= 0.3849 + 0.3159 \\ &= 0.7008 \end{aligned}$$

When it is converted to percentage (ie) 70% of the observations are covered between 9.6 to 13.8.

Example 7: For a normal distribution whose mean is 2 and standard deviation 3. Find the value of the variate such that the probability of the variate from the mean to the value is 0.4115

Solution:

Given that $Z \sim N(2, 9)$

To find X_1 :

We have $P(2 \leq Z \leq X_1) = 0.4115$

$$P\left(\frac{2-2}{3} \leq \frac{X-\mu}{\sigma} \leq \frac{X_1-2}{3}\right) = 0.4115$$

$$P(0 \leq Z \leq Z_1) = 0.4115 \text{ where } Z_1 = \frac{X_1 - 2}{3}$$

[From the normal table where 0.4115 lies is the value of Z_1]

From the normal table we have $Z_1 = 1.35$

$$\therefore 1.35 = \frac{X_1 - 2}{3}$$

$$\Rightarrow 3(1.35) + 2 = X_1$$

$$= X_1 = 6.05$$

(i.e) 41 % of the observation converged between 2 and 6.05

Questions

1. For a Poisson distribution

- | | |
|---------------------|---------------------|
| (a) mean > variance | (b) mean = variance |
| (c) mean < variance | (d) mean < variance |

Ans: mean = variance

2. In normal distribution, skewness is

- | | |
|----------------------|-------------------|
| (a) one | (b) zero |
| (c) greater than one | (d) less than one |

Ans: zero

3. Poisson distribution is a distribution for rare events

Ans: True

4. The total area under normal probability curve is one.

Ans: True

5. Poisson distribution is for continuous variable.

Ans: False

6. In a symmetrical curve mean, median and mode will coincide.

Ans: True

7. Give any two examples of Poisson distribution

8. The variance of a Poisson distribution is 0.5. Find $P(x = 3)$.

$[e^{-0.5} = 0.6065]$

9. The customer accounts of a certain departmental store have an average balance of Rs.1200 and a standard deviation of Rs.400. Assuming that the account balances are normally distributed. (i) what percentage of the accounts is over Rs.1500? (ii) What percentage of the accounts is between Rs.1000 and Rs.1500? iii) What percentage of the accounts is below Rs.1500?

10. State the Properties of normal distribution

Lecture.8

Sampling-basic concepts- sampling vs complete enumeration parameter and statistic-sampling methods-simple random sampling and stratified random sampling

Sampling

Population (Universe)

Population means aggregate of all possible units. It need not be human population. It may be population of plants, population of insects, population of fruits, etc.

Finite population

When the number of observation can be counted and is definite, it is known as finite population

- No. of plants in a plot.
- No. of farmers in a village.
- All the fields under a specified crop.

Infinite population

When the number of units in a population is innumerably large, that we cannot count all of them, it is known as infinite population.

- The plant population in a region.
- The population of insects in a region.

Frame

A list of all units of a population is known as frame.

Parameter

A summary measure that describes any given characteristic of the population is known as parameter. Population are described in terms of certain measures like mean, standard deviation etc. These measures of the population are called parameter and are usually denoted by Greek letters. For example, population mean is denoted by μ , standard deviation by σ and variance by σ^2 .

Sample

A portion or small number of unit of the total population is known as sample.

- All the farmers in a village(population) and a few farmers(sample)
- All plants in a plot is a population of plants.
- A small number of plants selected out of that population is a sample of plants.

Statistic

A summary measure that describes the characteristic of the sample is known as statistic. Thus sample mean, sample standard deviation etc is statistic. The statistic is usually denoted by roman letter.

\bar{x} - sample mean

s – standard deviation

The statistic is a random variable because it varies from sample to sample.

Sampling

The method of selecting samples from a population is known as sampling.

Sampling technique

There are two ways in which the information is collected during statistical survey.

They are

1. Census survey
2. Sampling survey

Census

It is also known as population survey and complete enumeration survey. Under census survey the information are collected from each and every unit of the population or universe.

Sample survey

A sample is a part of the population. Information are collected from only a few units of a population and not from all the units. Such a survey is known as sample survey.

Sampling technique is universal in nature, consciously or unconsciously it is adopted in every day life.

For eg.

1. A handful of rice is examined before buying a sack.
2. We taste one or two fruits before buying a bunch of grapes.
3. To measure root length of plants only a portion of plants are selected from a plot.

Need for sampling

The sampling methods have been extensively used for a variety of purposes and in great diversity of situations.

In practice it may not be possible to collect information on all units of a population due to various reasons such as

1. Lack of resources in terms of money, personnel and equipment.
2. The experimentation may be destructive in nature. Eg- finding out the germination percentage of seed material or in evaluating the efficiency of an insecticide the experimentation is destructive.
3. The data may be wasteful if they are not collected within a time limit. The census survey will take longer time as compared to the sample survey. Hence for getting quick results sampling is preferred. Moreover a sample survey will be less costly than complete enumeration.
4. Sampling remains the only way when population contains infinitely many number of units.
5. Greater accuracy.

Sampling methods

The various methods of sampling can be grouped under

- 1) Probability sampling or random sampling

2) Non-probability sampling or non random sampling

Random sampling

Under this method, every unit of the population at any stage has equal chance (or) each unit is drawn with known probability. It helps to estimate the mean, variance etc of the population.

Under probability sampling there are two procedures

1. Sampling with replacement (SWR)
2. Sampling without replacement (SWOR)

When the successive draws are made with placing back the units selected in the preceding draws, it is known as sampling with replacement. When such replacement is not made it is known as sampling without replacement.

When the population is finite sampling with replacement is adopted otherwise SWOR is adopted.

Mainly there are many kinds of random sampling. Some of them are.

1. Simple Random Sampling
2. Systematic Random Sampling
3. Stratified Random Sampling
4. Cluster Sampling

Simple Random sampling (SRS)

The basic probability sampling method is the simple random sampling. It is the simplest of all the probability sampling methods. It is used when the population is homogeneous.

When the units of the sample are drawn independently with equal probabilities. The sampling method is known as Simple Random Sampling (SRS). Thus if the population consists of N units, the probability of selecting any unit is $1/N$.

A theoretical definition of SRS is as follows

Suppose we draw a sample of size n from a population of size N . There are NC_n possible samples of size n . If all possible samples have an equal probability $1/NC_n$ of being drawn, the sampling is said to be simple random sampling.

There are two methods in SRS

1. Lottery method
2. Random no. table method

Lottery method

This is the most popular method and simplest method. In this method all the items of the universe are numbered on separate slips of paper of same size, shape and color. They are folded and mixed up in a drum or a box or a container. A blindfold selection is made. Required number of slips is selected for the desired sample size. The selection of items thus depends on chance.

For example, if we want to select 5 plants out of 50 plants in a plot, we number the 50 plants first. We write the numbers from 1-50 on slips of the same size, roll them and mix them. Then we make a blindfold selection of 5 plants. This method is also called unrestricted random sampling because units are selected from the population without any restriction. This method is mostly used in lottery draws. If the population is infinite, this method is inapplicable. There is a lot of possibility of personal prejudice if the size and shape of the slips are not identical.

Random number table method

As the lottery method cannot be used when the population is infinite, the alternative method is using a table of random numbers.

There are several standard tables of random numbers. But the credit for this technique goes to Prof. LHC. Tippett (1927). The random number table consists of 10,400 four-figured numbers. There are various other random numbers. They are fishers and Yates (1938) comprising of 15,000 digits arranged in twos. Kendall and B.B Smith (1939) consisting of 1, 00,000 numbers grouped in 25,000 sets of 4 digit random numbers, Rand corporation (1955) consisting of 2, 00,000 random numbers of 5 digits each etc.,

Merits

1. There is less chance for personal bias.
2. Sampling error can be measured.
3. This method is economical as it saves time, money and labour.

Demerits

1. It cannot be applied if the population is heterogeneous.
2. This requires a complete list of the population but such up-to-date lists are not available in many enquires.
3. If the size of the sample is small, then it will not be a representative of the population.

Stratified Sampling

When the population is heterogeneous with respect to the characteristic in which we are interested, we adopt stratified sampling.

When the heterogeneous population is divided into homogenous sub-population, the sub-populations are called strata. From each stratum a separate sample is selected using simple random sampling. This sampling method is known as stratified sampling.

We may stratify by size of farm, type of crop, soil type, etc.

The number of units to be selected may be uniform in all strata (or) may vary from stratum to stratum.

There are four types of allocation of strata

1. Equal allocation
2. Proportional allocation
3. Neyman's allocation
4. Optimum allocation

If the number of units to be selected is uniform in all strata it is known as equal allocation of samples.

If the number of units to be selected from a stratum is proportional to the size of the stratum, it is known as proportional allocation of samples.

When the cost per unit varies from stratum to stratum, it is known as optimum allocation.

When the costs for different strata are equal, it is known as Neyman's allocation.

Merits

1. It is more representative.
2. It ensures greater accuracy.
3. It is easy to administrate as the universe is sub-divided.

Demerits

1. To divide the population into homogeneous strata, it requires more money, time and statistical experience which is a difficult one.
2. If proper stratification is not done, the sample will have an effect of bias.

Questions

1. If each and every unit of population has equal chance of being included in the sample, it is known as

- (a) Restricted sampling
- (b) Purposive sampling
- (c) Simple random sampling
- (d) None of the above

Ans: Simple random sampling

2. In a population of size 10 the possible number of samples of size 2 will be

- (a) 45
- (b) 40
- (c) 54
- (d) None of the above

Ans: 45

3. A population consisting of an unlimited number of units is called an infinite population.

Ans: True

4. If all the units of a population are surveyed it is called census.

Ans: True

5. Random numbers are used for selecting the samples in simple random sampling method.

Ans: True

6. The list of all units in a population is called as Frame.

Ans: True

7. What is sampling?

8. Explain the Lottery method.

9. Explain the method of selection of samples in simple random sampling.

10. Explain the method of selection of samples in Stratified random sampling.

Lecture.9

Test of significance – Basic concepts – null hypothesis – alternative hypothesis – level of significance – Standard error and its importance – steps in testing

Test of Significance

Objective

To familiarize the students about the concept of testing of any hypothesis, the different terminologies used in testing and application of different types of tests.

Sampling Distribution

By drawing all possible samples of same size from a population we can calculate the statistic, for example, \bar{x} for all samples. Based on this we can construct a frequency distribution and the probability distribution of \bar{x} . Such probability distribution of a statistic is known as a sampling distribution of that statistic. In practice, the sampling distributions can be obtained theoretically from the properties of random samples.

Standard Error

As in the case of population distribution the characteristic of the sampling distributions are also described by some measurements like mean & standard deviation. Since a statistic is a random variable, the mean of the sampling distribution of a statistic is called the expected value of the statistic. The SD of the sampling distributions of the statistic is called standard error of the Statistic. The square of the standard error is known as the variance of the statistic. It may be noted that the standard deviation is for units whereas the standard error is for the statistic.

Theory of Testing Hypothesis

Hypothesis

Hypothesis is a statement or assumption that is yet to be proved.

Statistical Hypothesis

When the assumption or statement that occurs under certain conditions is formulated as scientific hypothesis, we can construct criteria by which a scientific hypothesis is either rejected or provisionally accepted. For this purpose, the scientific hypothesis is translated into statistical language. If the hypothesis is given in a statistical language it is called a statistical hypothesis.

For eg:-

The yield of a new paddy variety will be 3500 kg per hectare – scientific hypothesis.

In Statistical language it may be stated as the random variable (yield of paddy) is distributed normally with mean 3500 kg/ha.

Simple Hypothesis

When a hypothesis specifies all the parameters of a probability distribution, it is known as simple hypothesis. The hypothesis specifies all the parameters, i.e μ and σ of a normal distribution.

Eg:-

The random variable x is distributed normally with mean $\mu=0$ & $SD=1$ is a simple hypothesis. The hypothesis specifies all the parameters (μ & σ) of a normal distributions.

Composite Hypothesis

If the hypothesis specifies only some of the parameters of the probability distribution, it is known as composite hypothesis. In the above example if only the μ is specified or only the σ is specified it is a composite hypothesis.

Null Hypothesis - H_0

Consider for example, the hypothesis may be put in a form ‘paddy variety A will give the same yield per hectare as that of variety B’ or there is no difference between the average yields of paddy varieties A and B. These hypotheses are in definite terms. Thus these hypothesis form a basis to work with. Such a working hypothesis is known as null hypothesis. It is called null hypothesis because it nullifies the original hypothesis, that variety A will give more yield than variety B.

The null hypothesis is stated as ‘there is no difference between the effect of two treatments or there is no association between two attributes (ie) the two attributes are independent. Null hypothesis is denoted by H_0 .

Eg:-

There is no significant difference between the yields of two paddy varieties (or) they give same yield per unit area. Symbolically, $H_0: \mu_1=\mu_2$.

Alternative Hypothesis

When the original hypothesis is $\mu_1 > \mu_2$ stated as an alternative to the null hypothesis is known as alternative hypothesis. Any hypothesis which is complementary to null hypothesis is called alternative hypothesis, usually denoted by H_1 .

Eg:-

There is a significance difference between the yields of two paddy varieties. Symbolically,

$$H_1: \mu_1 \neq \mu_2 \text{ (two sided or directionless alternative)}$$

If the statement is that A gives significantly less yield than B (or) A gives significantly more yield than B. Symbolically,

$$H_1: \mu_1 < \mu_2 \text{ (one sided alternative-left tailed)}$$

$$H_1: \mu_1 > \mu_2 \text{ (one sided alternative-right tailed)}$$

Testing of Hypothesis

Once the hypothesis is formulated we have to make a decision on it. A statistical procedure by which we decide to accept or reject a statistical hypothesis is called testing of hypothesis.

Sampling Error

From sample data, the statistic is computed and the parameter is estimated through the statistic. The difference between the parameter and the statistic is known as the sampling error.

Test of Significance

Based on the sampling error the sampling distributions are derived. The observed results are then compared with the expected results on the basis of sampling distribution. If the difference between the observed and expected results is more than specified quantity of the standard error of the statistic, it is said to be significant at a specified probability level. The process up to this stage is known as test of significance.

Decision Errors

By performing a test we make a decision on the hypothesis by accepting or rejecting the null hypothesis H_0 . In the process we may make a correct decision on H_0 or commit one of two kinds of error.

- We may reject H_0 based on sample data when in fact it is true. This error in decisions is known as Type I error.

- We may accept H_0 based on sample data when in fact it is not true. It is known as Type II error.

	Accept H_0	Reject H_0
H_0 is true	Correct Decision	Type I error
H_0 is false	Type II error	Correct Decision

The relationship between type I & type II errors is that if one increases the other will decrease. The probability of type I error is denoted by α . The probability of type II error is denoted by β . The correct decision of rejecting the null hypothesis when it is false is known as the power of the test. The probability of the power is given by $1-\beta$.

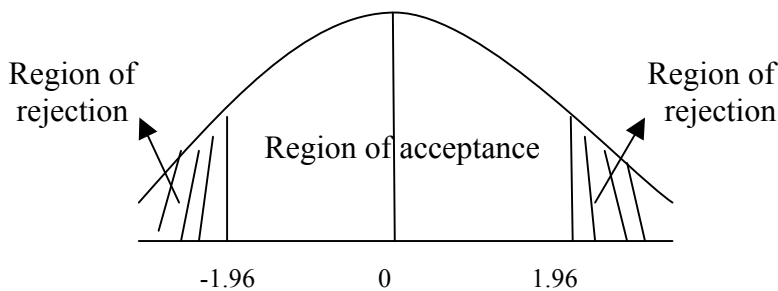
Critical Region

The testing of statistical hypothesis involves the choice of a region on the sampling distribution of statistic. If the statistic falls within this region, the null hypothesis is rejected; otherwise it is accepted. This region is called critical region.

Let the null hypothesis be $H_0: \mu_1 = \mu_2$ and its alternative be $H_1: \mu_1 \neq \mu_2$. Suppose H_0 is true. Based on sample data it may be observed that statistic $(\bar{x}_1 - \bar{x}_2)$ follows a normal distribution given by

$$Z = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{SE(\bar{x}_1 - \bar{x}_2)}$$

We know that 95% values of the statistic from repeated samples will fall in the range $(\bar{x}_1 - \bar{x}_2) \pm 1.96 \text{ times } SE(\bar{x}_1 - \bar{x}_2)$. This is represented by a diagram.



The border line value ± 1.96 is the critical value or tabular value of Z. The area beyond the critical values (shaded area) is known as critical region or region of rejection. The remaining area is known as region of acceptance.

If the statistic falls in the critical region we reject the null hypothesis and, if it falls in the region of acceptance we accept the null hypothesis.

In other words if the calculated value of a test statistic (Z, t, χ^2 etc) is more than the critical value in magnitude it is said to be significant and we reject H_0 and otherwise we accept H_0 . The critical values for the t and χ^2 are given in the form of readymade tables. Since the critical values are given in the form of table it is commonly referred as table value. The table value depends on the level of significance and degrees of freedom.

Example: $Z_{cal} < Z_{tab}$ -We accept the H_0 and conclude that there is no significant difference between the means

Test Statistic

The sampling distribution of a statistic like Z, t, & χ^2 are known as test statistic. Generally, in case of quantitative data

$$\text{Test statistic} = \frac{\text{Statistic} - \text{Parameter}}{\text{Standard Error(Statistic)}}$$

Note

The choice of the test statistic depends on the nature of the variable (ie) qualitative or quantitative, the statistic involved (i.e) mean or variance and the sample size, (i.e) large or small.

Level of Significance

The probability that the statistic will fall in the critical region is $\frac{\alpha}{2} + \frac{\alpha}{2} = \alpha$. This α is nothing but the probability of committing type I error. Technically the probability of committing type I error is known as level of Significance.

One and two tailed test

The nature of the alternative hypothesis determines the position of the critical region. For example, if H_1 is $\mu_1 \neq \mu_2$ it does not show the direction and hence the critical region falls on either end of the sampling distribution. If H_1 is $\mu_1 < \mu_2$ or $\mu_1 > \mu_2$ the direction is known. In the

first case the critical region falls on the left of the distribution whereas in the second case it falls on the right side.

One tailed test – When the critical region falls on one end of the sampling distribution, it is called one tailed test.

Two tailed test – When the critical region falls on either end of the sampling distribution, it is called two tailed test.

For example, consider the mean yield of new paddy variety (μ_1) is compared with that of a ruling variety (μ_2). Unless the new variety is more promising than the ruling variety in terms of yield we are not going to accept the new variety. In this case $H_1 : \mu_1 > \mu_2$ for which one tailed test is used. If both the varieties are new our interest will be to choose the best of the two. In this case $H_1 : \mu_1 \neq \mu_2$ for which we use two tailed test.

Degrees of freedom

The number of degrees of freedom is the number of observations that are free to vary after certain restrictions have been placed on the data. If there are n observations in the sample, for each restriction imposed upon the original observation the number of degrees of freedom is reduced by one.

The number of independent variables which make up the statistic is known as the degrees of freedom and is denoted by γ (Nu)

Steps in testing of hypothesis

The process of testing a hypothesis involves following steps.

1. Formulation of null & alternative hypothesis.
2. Specification of level of significance.
3. Selection of test statistic and its computation.
4. Finding out the critical value from tables using the level of significance, sampling distribution and its degrees of freedom.
5. Determination of the significance of the test statistic.
6. Decision about the null hypothesis based on the significance of the test statistic.
7. Writing the conclusion in such a way that it answers the question on hand.

Large sample theory

The sample size n is greater than 30 ($n \geq 30$) it is known as large sample. For large samples the sampling distributions of statistic are normal (Z test). A study of sampling distribution of statistic for large sample is known as large sample theory.

Small sample theory

If the sample size n is less than 30 ($n < 30$), it is known as small sample. For small samples the sampling distributions are t , F and χ^2 distribution. A study of sampling distributions for small samples is known as small sample theory.

Test of Significance

The theory of test of significance consists of various test statistics. The theory has been developed under two broad headings

1. Test of significance for large sample
Large sample test or Asymptotic test or Z test ($n \geq 30$)
2. Test of significance for small samples ($n < 30$)
Small sample test or Exact test- t , F and χ^2 .

It may be noted that small sample tests can be used in case of large samples also.

Large sample test

Large sample tests are

1. Sampling from attributes
2. Sampling from variables

Sampling from attributes

There are two types of tests for attributes

1. Test for single proportion
2. Test for equality of two proportions

Test for single proportion

In a sample of large size n , we may examine whether the sample would have come from a population having a specified proportion $P = P_0$. For testing

We may proceed as follows

1. Null Hypothesis (H₀)

H₀: The given sample would have come from a population with specified proportion P=P₀

2. Alternative Hypothesis(H₁)

H₁ : The given sample may not be from a population with specified proportion

P≠P₀ (Two Sided)

P>P₀(One sided-right sided)

P<P₀(One sided-left sided)

3. Test statistic

$$Z = \frac{|p - P|}{\sqrt{\frac{PQ}{n}}}$$

It follows a standard normal distribution with μ=0 and σ²=1

4. Level of Significance

The level of significance may be fixed at either 5% or 1%

5. Expected vale or critical value

In case of test statistic Z, the expected value is

$$Z_e = \left. \begin{array}{l} 1.96 \text{ at } 5\% \text{ level} \\ 2.58 \text{ at } 1\% \text{ level} \end{array} \right\} \longrightarrow \text{Two tailed test}$$

$$Z_e = \left. \begin{array}{l} 1.65 \text{ at } 5\% \text{ level} \\ 2.33 \text{ at } 1\% \text{ level} \end{array} \right\} \longrightarrow \text{One tailed test}$$

6. Inference

If the observed value of the test statistic Z₀ exceeds the table value Z_e we reject the Null Hypothesis H₀ otherwise accept it.

Test for equality of two proportions

Given two sets of sample data of large size n₁ and n₂ from attributes. We may examine whether the two samples come from the populations having the same proportion. We may proceed as follows:

1. Null Hypothesis (Ho)

Ho: The given two sample would have come from a population having the same proportion

$$P_1 = P_2$$

2. Alternative Hypothesis (H1)

H1 : The given two sample may not be from a population with specified proportion

$$P_1 \neq P_2 \text{ (Two Sided)}$$

$$P_1 > P_2 \text{ (One sided-right sided)}$$

$$P_1 < P_2 \text{ (One sided-left sided)}$$

3. Test statistic

$$Z = \frac{|(p_1 - p_2) - (P_1 - P_2)|}{\sqrt{\frac{P_1 Q_1}{n_1} + \frac{P_2 Q_2}{n_2}}}$$

When P₁ and P₂ are not known, then

$$Z = \frac{|p_1 - p_2|}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} \quad \text{for heterogeneous population}$$

Where q₁ = 1-p₁ and q₂ = 1-p₂

$$Z = \frac{|p_1 - p_2|}{\sqrt{pq \left(\frac{1}{n} + \frac{1}{n_2} \right)}} \quad \text{for homogeneous population}$$

p = combined or pooled estimate.

$$p = \frac{n_1 p_1 + n_2 p_2}{n_1 + n_2}$$

4. Level of Significance

The level may be fixed at either 5% or 1%

5. Expected value

The expected value is given by

$$Z_e = \quad 1.96 \text{ at 5\% level} \quad]$$

$Z_e =$	2.58 at 1% level	→	Two tailed test
	1.65 at 5% level	}	→ One tailed test
	2.33 at 1% level		

6. Inference

If the observed value of the test statistic Z exceeds the table value Z_e we may reject the Null Hypothesis H_0 otherwise accept it.

Sampling from variable

In sampling for variables, the test are as follows

1. Test for single Mean
2. Test for single Standard Deviation
3. Test for equality of two Means
4. Test for equality of two Standard Deviation

Test for single Mean

In a sample of large size n , we examine whether the sample would have come from a population having a specified mean

1. Null Hypothesis (H_0)

H_0 : There is no significance difference between the sample mean ie., $\mu = \mu_0$

or

The given sample would have come from a population having a specified mean

ie., $\mu = \mu_0$

2. Alternative Hypothesis(H_1)

H_1 : There is significance difference between the sample mean

ie., $\mu \neq \mu_0$ or $\mu > \mu_0$ or $\mu < \mu_0$

3. Test statistic

$$Z = \frac{|\bar{x} - \mu|}{\frac{\sigma}{\sqrt{n}}}$$

When population variance is not known, it may be replaced by its estimate

$$Z = \frac{|\bar{x} - \mu|}{\frac{s}{\sqrt{n}}}$$

where $s = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}}$

4. Level of Significance

The level may be fixed at either 5% or 1%

5. Expected value

The expected value is given by

$Z_e =$	1.96 at 5% level	}	→ Two tailed test
	2.58 at 1% level		

$Z_e =$	1.65 at 5% level	}	→ One tailed test
	2.33 at 1% level		

6. Inference

If the observed value of the test statistic Z exceeds the table value Z_e we may reject the Null Hypothesis H_0 otherwise accept it.

Test for equality of two Means

Given two sets of sample data of large size n_1 and n_2 from variables. We may examine whether the two samples come from the populations having the same mean. We may proceed as follows

1. Null Hypothesis (H_0)

H_0 : There is no significance difference between the sample mean i.e., $\mu = \mu_0$

or

The given sample would have come from a population having a specified mean

i.e., $\mu_1 = \mu_2$

2. Alternative Hypothesis (H_1)

H_1 : There is significance difference between the sample mean i.e., $\mu \neq \mu_0$

i.e., $\mu_1 \neq \mu_2$ or $\mu_1 < \mu_2$ or $\mu_1 > \mu_2$

3. Test statistic

When the population variances are known and unequal (i.e) $\sigma_1^2 \neq \sigma_2^2$

$$Z = \frac{|\bar{x}_1 - \bar{x}_2 - (\mu_1 - \mu_2)|}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

When $\sigma_1^2 = \sigma_2^2$,

$$Z = \frac{|\bar{x}_1 - \bar{x}_2|}{\sigma \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

where $\sigma = \frac{n_1\sigma_1^2 + n_2\sigma_2^2}{n_1 + n_2}$

The equality of variances can be tested by using F test.

When population variance is unknown, they may be replaced by their estimates s_1^2 and s_2^2

$$Z = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad \text{when } s_1^2 \neq s_2^2$$

when $s_1^2 = s_2^2$

$$Z = \frac{|\bar{x}_1 - \bar{x}_2|}{s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$\text{where } s^2 = \frac{n_1 s_1^2 + n_2 s_2^2}{n_1 + n_2}$$

4. Level of Significance

The level may be fixed at either 5% or 1%

5. Expected value

The expected value is given by

$Z_e =$	1.96 at 5% level	}	→ Two tailed test
	2.58 at 1% level		
$Z_e =$	1.65 at 5% level	}	→ One tailed test
	2.33 at 1% level		

6. Inference

If the observed value of the test statistic Z exceeds the table value Z_e we may reject the Null Hypothesis H_0 otherwise accept it.

Questions

1. A hypothesis may be classified as

- | | |
|------------|-------------------|
| (a) Simple | (b) Composite |
| (c) Null | (d) All the above |

Ans: All the above

2. Area of the critical region depends on

- | | |
|-----------------------------|----------------------------|
| (a) Size of type I error | (b) Size of type II error |
| (c) Value of the statistics | (d) Number of observations |

Ans: Size of type I error

3. Large sample test can be applied when the sample size exceeds 30.

Ans: True

4. If the calculated test statistic is greater than the critical value, the null hypothesis is rejected.

Ans: True

5. The standard error of mean is given by $\frac{\sigma}{\sqrt{n}}$

Ans: True

6. If the alternative hypothesis is $\mu_1 \neq \mu_2$ then the test is known as one tailed test.

Ans: False

7. Define standard error.

8. Define Type I and Type II error.

9. Describe the procedure of comparing two group means.

10. Describe the procedure of comparing two proportions.

Lecture.10

T-test – definition – assumptions – test for equality of two means-independent and paired t test

Student’s t test

When the sample size is smaller, the ratio $Z = \frac{|\bar{x} - \mu|}{\frac{s}{\sqrt{n}}}$ will follow t distribution

and not the standard normal distribution. Hence the test statistic is given as $t = \frac{|\bar{x} - \mu|}{\frac{s}{\sqrt{n}}}$

which follows normal distribution with mean 0 and unit standard deviation. This follows a t distribution with (n-1) degrees of freedom which can be written as $t_{(n-1)}$ d.f.

This fact was brought out by Sir William Gosset and Prof. R.A Fisher. Sir William Gosset published his discovery in 1905 under the pen name Student and later on developed and extended by Prof. R.A Fisher. He gave a test known as t-test.

Applications (or) uses

1. To test the single mean in single sample case.
2. To test the equality of two means in double sample case.
 - (i) Independent samples(Independent t test)
 - (ii) Dependent samples (Paired t test)
3. To test the significance of observed correlation coefficient.
4. To test the significance of observed partial correlation coefficient.
5. To test the significance of observed regression coefficient.

Test for single Mean

1. Form the null hypothesis

$$H_0: \mu = \mu_0$$

(i.e) There is no significance difference between the sample mean and the population mean

2. Form the Alternate hypothesis

$$H_1: \mu \neq \mu_0 \text{ (or } \mu > \mu_0 \text{ or } \mu < \mu_0)$$

ie., There is significance difference between the sample mean and the population mean

3. Level of Significance

The level may be fixed at either 5% or 1%

4. Test statistic

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \text{ which follows t distribution with } (n-1) \text{ degrees of freedom}$$

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

6. Find the table value of t corresponding to (n-1) d.f. and the specified level of significance.

7. Inference

If $t < t_{tab}$ we accept the null hypothesis H_0 . We conclude that there is no significant difference sample mean and population mean

(or) if $t > t_{tab}$ we reject the null hypothesis H_0 . (ie) we accept the alternative hypothesis and conclude that there is significant difference between the sample mean and the population mean.

Example 1

Based on field experiments, a new variety of green gram is expected to given a yield of 12.0 quintals per hectare. The variety was tested on 10 randomly selected farmer's fields. The yield (quintals/hectare) were recorded as 14.3,12.6,13.7,10.9,13.7,12.0,11.4,12.0,12.6,13.1. Do the results conform to the expectation?

Solution

Null hypothesis $H_0: \mu=12.0$

(i.e) the average yield of the new variety of green gram is 12.0 quintals/hectare.

Alternative Hypothesis: $H_1: \mu \neq 12.0$

(i.e) the average yield is not 12.0 quintals/hectare, it may be less or more than 12 quintals / hectare

Level of significance: 5 %

Test statistic:

$$t = \left| \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \right|$$

From the given data

$$\sum x = 126.3 \quad \sum x^2 = 1605.77$$

$$\bar{x} = \frac{\sum x}{n} = \frac{126.3}{10} = 12.63$$

$$s = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}} = \sqrt{\frac{1605.77 - \frac{1595.169}{9}}{9}} = \sqrt{\frac{10.601}{9}}$$

$$= 1.0853$$

$$\frac{s}{\sqrt{n}} = \frac{1.0853}{\sqrt{10}} = 0.3432$$

Now $t = \left| \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \right|$

$$t = \frac{12.63 - 12}{0.3432} = 1.836$$

Table value for t corresponding to 5% level of significance and 9 d.f. is 2.262 (two tailed test)

Inference

$$t < t_{\text{tab}}$$

We accept the null hypothesis H_0

We conclude that the new variety of green gram will give an average yield of 12 quintals/hectare.

Note

Before applying t test in case of two samples the equality of their variances has to be tested by using F-test

$$F = \frac{s_1^2}{s_2^2} \sim F_{(n_1-1, n_2-1)} \text{ d.f. if } s_1^2 > s_2^2$$

or

$$F = \frac{s_2^2}{s_1^2} \sim F_{(n_2-1, n_1-1)} \text{ d.f. if } s_2^2 > s_1^2$$

where s_1^2 is the variance of the first sample whose size is n_1 .

s_2^2 is the variance of the second sample whose size is n_2 .

It may be noted that the numerator is always the greater variance. The critical value for F is read from the F table corresponding to a specified d.f. and level of significance

Inference

$$F < F_{\text{tab}}$$

We accept the null hypothesis H_0 .(i.e) the variances are equal otherwise the variances are unequal.

Test for equality of two Means (Independent Samples)

Given two sets of sample observation $x_{11}, x_{12}, x_{13} \dots x_{1n}$, and $x_{21}, x_{22}, x_{23} \dots x_{2n}$ of sizes n_1 and n_2 respectively from the normal population.

1. Using F-Test , test their variances

(i) Variances are Equal

$$H_0: \mu_1 = \mu_2$$

$$H_1 \mu_1 \neq \mu_2 \text{ (or } \mu_1 < \mu_2 \text{ or } \mu_1 > \mu_2)$$

Test statistic

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

where the combined variance

$$s^2 = \frac{\left[\sum x_1^2 - \frac{(\sum x_1)^2}{n_1} \right] + \left[\sum x_2^2 - \frac{(\sum x_2)^2}{n_2} \right]}{n_1 + n_2 - 2}$$

The test statistic t follows a t distribution with (n1+n2-2) d.f.

(ii) Variances are unequal and n1=n2

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

It follows a t distribution with $\left(\frac{n_1 + n_2}{2} \right) - 1$ d.f.

(i) Variances are unequal and n1≠n2

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

This statistic follows neither t nor normal distribution but it follows Behrens-Fisher distribution. The Behrens – Fisher test is laborious one. An alternative simple method has been suggested by Cochran & Cox. In this method the critical value of t is altered as t_w (i.e) weighted t

$$t_w = \frac{t_1 \left(\frac{s_1^2}{n_1} \right) + t_2 \left(\frac{s_2^2}{n_2} \right)}{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

where t_1 is the critical value for t with (n1-1) d.f. at a dspecified level of significance and t_2 is the critical value for t with (n2-1) d.f. at a dspecified level of significance and

Example 2

In a fertilizer trial the grain yield of paddy (Kg/plot) was observed as follows

Under ammonium chloride 42,39,38,60 &41 kgs

Under urea 38, 42, 56, 64, 68, 69,& 62 kgs.

Find whether there is any difference between the sources of nitrogen?

Solution

Ho: $\mu_1 = \mu_2$ (i.e) there is no significant difference in effect between the sources of nitrogen.

H₁: $\mu_1 \neq \mu_2$ (i.e) there is a significant difference between the two sources

Level of significance = 5%

Before we go to test the means first we have to test their variances by using F-test.

F-test

Ho:., $\sigma_1^2 = \sigma_2^2$

H1:., $\sigma_1^2 \neq \sigma_2^2$

$$s_1^2 = \frac{\sum x_1^2 - \frac{(\sum x_1)^2}{n1}}{n1 - 1} = 82.5$$

$$s_2^2 = \frac{\sum x_2^2 - \frac{(\sum x_2)^2}{n2}}{n2 - 1} = 154.33$$

$\therefore F = \frac{s_2^2}{s_1^2} \sim F_{(n_2 - 1, n_1 - 1)}$ d.f if $s_2^2 > s_1^2$

$$F = \frac{154.33}{32.5} = 1.8707$$

$F_{tab}(6,4)$ d.f. = 6.16

$\Rightarrow F < F_{tab}$

We accept the null hypothesis H₀. (i.e) the variances are equal.

Use the test statistic

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

where

$$s^2 = \frac{\left[\sum x_1^2 - \frac{(\sum x_1)^2}{n_1} \right] + \left[\sum x_2^2 - \frac{(\sum x_2)^2}{n_2} \right]}{n_1 + n_2 - 2} = \frac{330 + 926}{10} = 125.6$$

$$t = \frac{|44 - 57|}{\sqrt{125.6 \left(\frac{1}{7} + \frac{1}{75} \right)}} = 1.98$$

The degrees of freedom is $5+7-2= 10$. For 5 % level of significance, table value of t is 2.228

Inference:

$$t < t_{\text{tab}}$$

We accept the null hypothesis H_0

We conclude that the two sources of nitrogen do not differ significantly with regard to the grain yield of paddy.

Example 3

The summary of the results of an yield trial on onion with two methods of propagation is given below. Determine whether the methods differ with regard to onion yield. The onion yield is given in Kg/plot.

Method I	Method II
$n_1=12$	$n_2=12$
$\bar{x}_1 = 25.25$	$\bar{x}_2 = 28.83$
$SS_1=186.25$	$SS_2=737.6667$
$s_1^2 = 16.9318$	$s_2^2 = 67.0606$

Solution

$H_0: \mu_1 = \mu_2$ (i.e) the two propagation methods do not differ with regard to onion yield.

$H_1: \mu_1 \neq \mu_2$ (i.e) the two propagation methods differ with regard to onion yield.

Level of significance = 5%

Before we go to test the means first we have to test their variability using F-test.

F-test

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$H_1: \sigma_1^2 \neq \sigma_2^2$$

$$s_1^2 = \frac{\sum x_1^2 - \frac{(\sum x_1)^2}{n1}}{n1 - 1} = 16.9318$$

$$s_2^2 = \frac{\sum x_2^2 - \frac{(\sum x_2)^2}{n2}}{n2 - 1} = 67.0606$$

$$\therefore F = \frac{s_2^2}{s_1^2} \sim F_{(n_2 - 1, n_1 - 1)} \text{ d.f. if } s_2^2 > s_1^2$$

$$F = \frac{67.0606}{16.9318} = 3.961$$

$$F_{\text{tab}}(11, 11) \text{ d.f.} = 2.82$$

$$\Rightarrow F > F_{\text{tab}}$$

We reject the null hypothesis H_0 . we conclude that the variances are unequal.

Here the variances are unequal with equal sample size then the test statistic is

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$$

where

$$s^2 = \frac{\left[\sum x_1^2 - \frac{(\sum x_1)^2}{n1} \right] + \left[\sum x_2^2 - \frac{(\sum x_2)^2}{n2} \right]}{n_1 + n_2 - 2}$$

$$s^2 = \frac{SS1 + SS2}{n_1 + n_2 - 2} = \frac{186.25 + 737.6667}{12 + 12 - 2} = 41.9962$$

$$t = \frac{25.25 - 28.83}{\sqrt{41.9962 \left(\frac{1}{12} + \frac{1}{12} \right)}} = \frac{3.58}{\sqrt{6.9994}} = 1.353$$

$$t = 1.353$$

The table value for $\left(\frac{12+12}{2} - 1 \right) = 11$ d.f. at 5% level of significance is 2.201

Inference:

$$t < t_{\text{tab}}$$

We accept the null hypothesis H_0

We conclude that the two propagation methods do not differ with regard to onion yield.

Example 4

The following data relate the rubber yield of two types of rubber plants, where the sample have been drawn independently. Test whether the two types of rubber plants differ in their yield.

Type I	6.21	5.70	6.04	4.47	5.22	4.45	4.84	5.84	5.88	5.82	6.09	5.59
	6.06	5.59	6.74	5.55								
Type II	4.28	7.71	6.48	7.71	7.37	7.20	7.06	6.40	8.93	5.91	5.51	6.36

Solution

$H_0: \mu_1 = \mu_2$ (i.e) there is no significant difference between the two rubber plants.

$H_1: \mu_1 \neq \mu_2$ (i.e) there is a significant difference between the two rubber plants.

Level of significance = 5%

Here

$n_1 = 16$	$n_2 = 12$
$\sum x_1 = 90.09$	$\sum x_2 = 80.92$
$\bar{x}_1 = 5.63$	$\bar{x}_2 = 6.7431$
$\sum x_1^2 = 513.085$	$\sum x_2^2 = 561.64$

Before we go to test the means first we have to test their variability using F-test.

F-test

$H_0: \sigma_1^2 = \sigma_2^2$

$H_1: \sigma_1^2 \neq \sigma_2^2$

$$s_1^2 = \frac{\sum x_1^2 - \frac{(\sum x_1)^2}{n_1}}{n_1 - 1} = 0.388$$

$$s_2^2 = \frac{\sum x_2^2 - \frac{(\sum x_2)^2}{n_2}}{n_2 - 1} = 1.452$$

$$\therefore F = \frac{s_2^2}{s_1^2} \sim F_{(n_2 - 1, n_1 - 1)} \text{ d.f. if } s_2^2 > s_1^2$$

$$F = \frac{1.452}{0.388} = 3.742$$

$$F_{\text{tab}}(11, 15) \text{ d.f.} = 2.51$$

$$\Rightarrow F > F_{\text{tab}}$$

We reject the null hypothesis H_0 . Hence, the variances are unequal.

Here the variances are unequal with unequal sample size then the test statistic is

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$t = \frac{(5.63 - 6.7431_2)}{\sqrt{\frac{0.388}{16} + \frac{1.452}{12}}} = 2.912$$

$$t_w = \frac{t_1 \left(\frac{S_1^2}{n_1} \right) + t_2 \left(\frac{S_2^2}{n_2} \right)}{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

$$t_1 = t_{(16-1)} \text{ d.f.} = 2.131$$

$$t_2 = t_{(12-1)} \text{ d.f.} = 2.201$$

$$t_w = \frac{2.131 \left(\frac{0.388}{16} \right) + 2.201 \left(\frac{1.452}{12} \right)}{\frac{0.388}{16} + \frac{1.425}{12}} = 2.187$$

Inference:

$$t > t_w$$

We reject the null hypothesis H_0 . We conclude that the second type of rubber plant yields more rubber than that of first type.

Equality of two means (Dependant samples)

Paired t test

In the t-test for difference between two means, the two samples were independent of each other. Let us now take particular situations where the samples are not independent.

In agricultural experiments it may not be possible to get required number of homogeneous experimental units. For example, required number of plots which are similar in all; characteristics may not be available. In such cases each plot may be divided into two equal parts and one treatment is applied to one part and second treatment to another part of the plot. The results of the experiment will result in two correlated samples. In some other situations two observations may be taken on the same experimental unit. For example, the soil properties before and after the application of industrial effluents may be observed on number of plots. This will result in paired observation. In such situations we apply paired t test.

Suppose the observation before treatment is denoted by x and the observation after treatment is denoted by y. for each experimental unit we get a pair of observation(x,y). In case of n experimental units we get n pairs of observations : (x1,y1), (x2,y2)...(xn,yn). In order to apply the paired t test we find out the differences (x1- y1), (x2-y2),...,(xn-yn) and denote them as d1, d2,...,dn. Now d1, d2...form a sample . we

apply the t test procedure for one sample (i.e) $t = \frac{|\bar{d}|}{\sqrt{s^2/n}}$

$$\text{where } \bar{d} = \frac{\sum di}{n}, s^2 = \frac{\sum di^2 - \frac{(\sum di)^2}{n}}{n-1}$$

the mean \bar{d} may be positive or negative. Hence we take the absolute value as $|\bar{d}|$. The test statistic t follows a t distribution with (n-1) d.f.

Example 5

In an experiment the plots where divided into two equal parts. One part received soil treatment A and the second part received soil treatment B. each plot was planted with sorghum. The sorghum yield (kg/plot) was absorbed. The results are given below. Test the effectiveness of soil treatments on sorghum yield.

Soil treatment A	49	53	51	52	47	50	52	53
Soil treatment B	52	55	52	53	50	54	54	53

Solution

H₀: $\mu_1 = \mu_2$, there is no significant difference between the effects of the two soil treatments

H₁: $\mu_1 \neq \mu_2$, there is significant difference between the effects of the two soil treatments

Level of significance = 5%

Test statistic

$$t = \frac{|\bar{d}|}{\sqrt{s^2 / n}}$$

x	y	d=x-y	d ²
49	52	-3	9
53	55	-2	4
51	52	-1	1
51	52	-1	1
47	50	-3	16
50	54	-4	16
52	54	-2	4
53	53	0	0
Total		-16	44

$$\bar{d} = \frac{\sum di}{n} = \frac{-16}{8} = -2,$$

$$s^2 = \frac{\sum di^2 - \frac{(\sum di)^2}{n}}{n-1} = 1.7143$$

$$t = \frac{|-2|}{\sqrt{1.7143/8}} = 4.32$$

Table value of t for 7 d.f. at 5% l.o.s is 2.365

Inference:

$$t > t_{tab}$$

We reject the null hypothesis H_0 . We conclude that there is a significant difference between the two soil treatments between A and B. Soil treatment B increases the yield of sorghum significantly,

Questions

1. The test statistic $F = \frac{s_1^2}{s_2^2}$ is used for testing

- (a) $H_0: \mu_1 = \mu_2$ (b) $H_0: \sigma_1^2 = \sigma_2^2$
- (c) $H_0: \sigma_1 = \sigma_2$ (d) $H_0: \sigma_2 = \sigma_1$

Ans: $H_0: \sigma_1^2 = \sigma_2^2$

2. In paired t test with n observations in each group the degrees of freedom is

- (a) n (b) n-1 (c) n-2 (d) n+1

Ans: n-1

3. Student t- test is applicable in case of small samples.

Ans: True

4.F test is also known as variance ratio test.

Ans: True

5. In case of comparing the equality of two variances the greater variance should be taken in the numerator.

Ans: True

6. While comparing the means of two independent samples the variances of the two samples will be always equal.

Ans: False

7. Define t statistic.

8. Define F statistic.

9. Explain the procedure of testing the equality of two variances.

10. How to compare the means of two independent small samples.

Lecture.11

Attributes- Contingency table – 2x2 contingency table – Test for independence of attributes – test for goodness of fit of mendalian ratio

Test based on χ^2 -distribution

In case of attributes we can not employ the parametric tests such as F and t. Instead we have to apply χ^2 test. When we want to test whether a set of observed values are in agreement with those expected on the basis of some theories or hypothesis. The χ^2 statistic provides a measure of agreement between such observed and expected frequencies.

The χ^2 test has a number of applications. It is used to

- (1) Test the independence of attributes
- (2) Test the goodness of fit
- (3) Test the homogeneity of variances
- (4) Test the homogeneity of correlation coefficients
- (5) Test the equaslity of several proportions.

In genetics it is applied to detect linkage.

Applications

χ^2 – test for goodness of fit

A very powerful test for testing the significance of the discrepancy between theory and experiment was given by Prof. Karl Pearson in 1900 and is known as “chi-square test of goodness of fit “.

If O_i , ($i=1,2,\dots,n$) is a set of observed (experimental frequencies) and E_i ($i=1,2,\dots,n$) is the corresponding set of expected (theoretical or hypothetical) frequencies, then,

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

It follows a χ^2 distribution with n-1 d.f. In case of χ^2 only one tailed test is used.

Example

In plant genetics, our interest may be to test whether the observed segregation ratios deviate significantly from the mendelian ratios. In such situations we want to test the agreement between the observed and theoretical frequency, such test is called as test of goodness of fit.

Conditions for the validity of χ^2 -test:

χ^2 -test is an approximate test for large values of ‘n’ for the validity of χ^2 -test of goodness of fit between theory and experiment, the following conditions must be satisfied.

1. The sample observations should be independent.
2. Constraints on the cell frequency, if any, should be linear.

Example: $\sum O_i = \sum E_i$.

3. N, the total frequency should be reasonably large, say greater than (>) 50.
4. No theoretical cell frequency should be less than (<)5. If any theoretical cell frequency is <5, then for the application of χ^2 - test, it is pooled with the preceding or succeeding frequency so that the pooled frequency is more than 5 and finally adjust for degree's of freedom lost in pooling.

Example1

The number of yeast cells counted in a haemocytometer is compared to the theoretical value is given below. Does the experimental result support the theory?

No. of Yeast cells in the square	Obseved Frequency	Expected Frequency
0	103	106
1	143	141
2	98	93
3	42	41
4	8	14
5	6	5

Solution

H₀: the experimental results support the theory

H₁: the experimental results does not support the theory.

Level of significance=5%

Test Statistic:

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

O _i	E _i	O _i .E _i	(O _i .E _i) ²	(O _i .E _i) ² /E _i
103	106	-3	9	0.0849
143	141	2	4	0.0284
98	93	5	25	0.2688
42	41	1	1	0.0244
8	14	-6	36	2.5714
6	5	1	1	0.2000
400	400			3.1779

∴ χ² = 3.1779

Table value

χ² (6-1=5 at 5 % l.os) = 11.070

Inference

χ² < χ²_{tab}

We accept the null hypothesis.

(i.e) there is a good correspondence between theory and experiment.

χ² test for independence of attributes

At times we may consider two characteristics on attributes simultaneously. Our interest will be to test the association between these two attributes

Example:- An entomologist may be interested to know the effectiveness of different concentrations of the chemical in killing the insects. The concentrations of chemical form one attribute. The state of insects ‘killed & not killed’ forms another attribute. The result of this experiment can be arranged in the form of a contingency table. In general one attribute may be divided into m classes as A₁, A₂,A_m and the other attribute may be divided into n classes as B₁, B₂,B_n. Then the contingency table will have m x n cells. It is termed as m x n contingency table

A	A1	A2	...	Aj	...	Am	Row Total
B							

B1	O11	O12	...	O1j		O1m	r1
B2	O21	O22	...	O2j		O2m	r2
⋮							
⋮							
Bi	Oij	Oi2	...	Oij		Oim	ri
⋮							
⋮							
Bn	On1	On2	...	Onj		Onm	rk
Column	c1	c2	...	cj	...	cm	$n = \sum ri = \sum cj$
Total							

where O_{ij} 's are observed frequencies.

The expected frequencies corresponding to O_{ij} is calculated as $\frac{r_i \cdot c_j}{n}$. The χ^2 is

computed as

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^m \frac{(o_{ij} - E_{ij})^2}{E_{ij}}$$

where

O_{ij} – observed frequencies

E_{ij} – Expected frequencies

n = number of rows

m = number of columns

It can be verified that $\sum O_{ij} = \sum E_{ij}$

This χ^2 is distributed as χ^2 with $(n-1)(m-1)$ d.f.

2x2 – contingency table

When the number of rows and number of columns are equal to 2 it is termed as 2 x 2 contingency table. It will be in the following form

	B ₁	B ₂	Row Total
A ₁	a	b	a+b r ₁
A ₂	c	d	c+d r ₂
Column Total	a+c	b+d	a+b+c+d =n
	c ₁	c ₂	

Where a, b, c and d are cell frequencies c1 and c2 are column totals, r1 and r2 are row totals and n is the total number of observations.

In case of 2 x 2 contingency table χ^2 can be directly found using the short cut formula,

$$\chi^2 = \frac{n(ad - bc)^2}{c1.c2.r1.r2}$$

The d.f associated with χ^2 is (2-1) (2-1) =1

Yates correction for continuity

If anyone of the cell frequency is < 5, we use Yates correction to make χ^2 as continuous. The yates correction is made by adding 0.5 to the least cell frequency and adjusting the other cell frequencies so that the column and row totals remain same . suppose, the first cell frequency is to be corrected then the contingency table will be as follows:

	B1	B2	Row Total
A1	a + 0.5	b - 0.5	a+b=r1
A2	c - 0.5	d + 0.5	c+d =r2
Column Total	a+c=c1	b+d=c2	n = a+b+c+d

Then use the χ^2 - statistic as

$$\chi^2 = \frac{n\left(|ad - bc| - \frac{n}{2}\right)^2}{c1.c2.r1.r2}$$

The d.f associated with χ^2 is (2-1) (2-1) =1

Exapmle 2

The severity of a disease and blood group were studied in a research project. The findings sre given in the following table, knowmn as the m xn contingency table. Can this severity of the condition and blood group are associated.

Severity of a disease classified by blood group in 1500 patients.

Condition	Blood Groups				Total
	O	A	B	AB	

Severe	51	40	10	9	110
Moderate	105	103	25	17	250
Mild	384	527	125	104	1140
Total	540	670	160	130	1500

Solution

H₀: The severity of the disease is not associated with blood group.

H₁: The severity of the disease is associated with blood group.

Calculation of Expected frequencies

Condition	Blood Groups				Total
	O	A	B	AB	
Severe	39.6	49.1	11.7	9.5	110
Moderate	90.0	111.7	26.7	21.7	250
Mild	410.4	509.2	121.6	98.8	1140
Total	540	670	160	130	1500

Test statistic:

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^m \frac{(o_{ij} - E_{ij})^2}{E_{ij}}$$

The d.f. associated with the χ^2 is (3-1)(4-1) = 6

Calculations

O _i	E _i	O _i ·E _i	(O _i ·E _i) ²	(O _i ·E _i) ² /E _i
51	39.6	11.4	129.96	3.2818
40	49.1	-9.1	82.81	1.6866
10	11.7	-1.7	2.89	0.2470
9	9.5	-0.5	0.25	0.0263
105	90.0	15	225.00	2.5000
103	111.7	-8.7	75.69	0.6776
25	26.7	-1.7	2.89	0.1082
17	21.7	-4.7	22.09	1.0180
384	410.4	-26.4	696.96	1.6982
527	509.2	17.8	316.84	0.6222
125	121.6	3.4	11.56	0.0951

104	98.8	5.2	27.04	0.2737
Total				12.2347

$$\therefore \chi^2 = 12.2347$$

Table value of χ^2 for 6 d.f. at 5% level of significance is 12.59

Inference

$$\chi^2 < \chi^2_{\text{tab}}$$

We accept the null hypothesis.

The severity of the disease has no association with blood group.

Example 3

In order to determine the possible effect of a chemical treatment on the rate of germination of cotton seeds a pot culture experiment was conducted. The results are given below

Chemical treatment and germination of cotton seeds

	Germinated	Not germinated	Total
Chemically Treated	118	22	140
Untreated	120	40	160
Total	238	62	300

Does the chemical treatment improve the germination rate of cotton seeds?

Solution

H_0 : The chemical treatment does not improve the germination rate of cotton seeds.

H_1 : The chemical treatment improves the germination rate of cotton seeds.

Level of significance = 1%

Test statistic

$$= \frac{n(ad - bc)^2}{(a + b)(c + d)(a + c)(b + d)} \text{ with 1 d.f.}$$

$$\chi^2 = \frac{300(118 \times 40 - 22 \times 120)^2}{140 \times 160 \times 62 \times 238} = 3.927$$

Table value

$$\chi^2 (1) \text{ d.f. at } 1\% \text{ L.O.S} = 6.635$$

Inference

$$\chi^2 < \chi^2_{\text{tab}}$$

We accept the null hypothesis.

The chemical treatment will not improve the germination rate of cotton seeds significantly.

Example 4

In an experiment on the effect of a growth regulator on fruit setting in muskmelon the following results were obtained. Test whether the fruit setting in muskmelon and the application of growth regulator are independent at 1% level.

	Fruit set	Fruit not set	Total
Treated	16	9	25
Control	4	21	25
Total	20	30	50

Solution

H₀: Fruit setting in muskmelon does not depend on the application of growth regulator.

H₁: Fruit setting in muskmelon depend on the application of growth regulator.

Level of significance = 1%

After Yates correction we have

	Fruit set	Fruit not set	Total
Treated	15.5	9.5	25
Control	4.5	20.5	25
Total	20	30	50

Tet statistic

$$\chi^2 = \frac{n \left(|ad - bc| - \frac{n}{2} \right)^2}{(a+b)(c+d)(a+c)(b+d)}$$

$$\chi^2 = \frac{50 \left[|15.5 \times 20.5 - 9.5 \times 4.5| - \frac{50}{2} \right]^2}{25 \times 25 \times 20 \times 30} = 8.33$$

Table value

χ^2 (1) d.f. at 1 % level of significance is 6.635

Inference

$$\chi^2 > \chi^2_{\text{tab}}$$

We reject the null hypothesis.

Fruit setting in muskmelon is influenced by the growth regulator. Application of growth regulator will increase fruit setting in musk melon.

Questions

1. The calculated value of χ^2 is
 (a) always positive (b) always negative
 (c) can be either positive or negative (d) none of these

Ans: always positive

2. Degrees of freedom for Chi-square in case of contingency table of order (4 × 3) are
 (a) 12 (b) 9 (c) 8 (d) 6

Ans: 6

3. One condition for application of χ^2 test is that no cell frequency should be less than five.

Ans: True

4. The distribution of the χ^2 depends on the degrees of freedom.

Ans: True

5. The greater the discrepancy between the observed and expected Frequency lesser the value of χ^2 .

Ans: False

6. When observed and expected frequencies completely coincide χ^2 will be zero.

Ans: True

7. What is a contingency table?

8. When and how to apply Yates correction?

9. Explain the χ^2 test of goodness of fit?

10. Explain how to test the independence of attributes?

Lecture.12**Correlation – definition – Scatter diagram -Pearson's correlation co-efficient –
properties of correlation coefficient****Correlation**

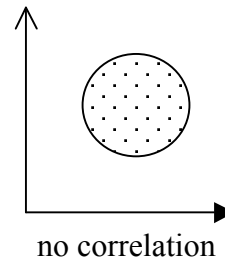
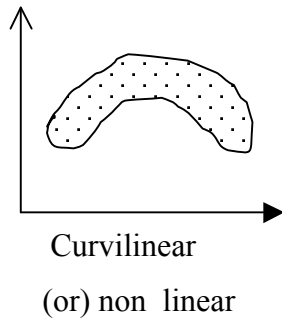
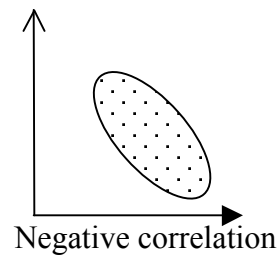
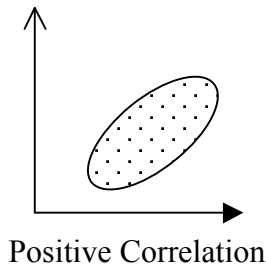
Correlation is the study of relationship between two or more variables. Whenever we conduct any experiment we gather information on more related variables. When there are two related variables their joint distribution is known as bivariate normal distribution and if there are more than two variables their joint distribution is known as multivariate normal distribution.

In case of bi-variate or multivariate normal distribution, we are interested in discovering and measuring the magnitude and direction of relationship between 2 or more variables. For this we use the tool known as correlation.

Suppose we have two continuous variables X and Y and if the change in X affects Y, the variables are said to be correlated. In other words, the systematic relationship between the variables is termed as correlation. When only 2 variables are involved the correlation is known as simple correlation and when more than 2 variables are involved the correlation is known as multiple correlation. When the variables move in the same direction, these variables are said to be correlated positively and if they move in the opposite direction they are said to be negatively correlated.

Scatter Diagram

To investigate whether there is any relation between the variables X and Y we use scatter diagram. Let $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be n pairs of observations. If the variables X and Y are plotted along the X-axis and Y-axis respectively in the x-y plane of a graph sheet the resultant diagram of dots is known as scatter diagram. From the scatter diagram we can say whether there is any correlation between x and y and whether it is positive or negative or the correlation is linear or curvilinear.



Pearsons Correlation coefficient

The measures of the degree of relationship between two continuous variables is called correlation coefficient. It is denoted by r (in case of sample) and ρ (in case of population). The correlation coefficient r is known as Pearson's correlation coefficient as it was discovered by Karl Pearson. It is also called as product moment correlation.

The correlation coefficient r is given as the ratio of covariance of the variables X and Y to the product of the standard deviation of X and Y .

Symbolically,

$$r = \frac{\frac{1}{n-1} \sum (x - \bar{x})(y - \bar{y})}{\sqrt{\frac{1}{n-1} \sum (x - \bar{x})^2} \sqrt{\frac{1}{n-1} \sum (y - \bar{y})^2}}$$

which can be simplified as

$$r = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sqrt{\sum x^2 - \frac{(\sum x)^2}{n}} \sqrt{\sum y^2 - \frac{(\sum y)^2}{n}}}$$

This correlation coefficient r is known as Pearson's Correlation coefficient. The numerator is termed as sum of product of X and Y and abbreviated as $SP(XY)$. In the denominator the first term is called sum of squares of X (i.e) $SS(X)$ and second term is called sum of squares of Y (i.e) $SS(Y)$

$$\therefore r = \frac{SP(XY)}{\sqrt{SS(X)}\sqrt{SS(Y)}}$$

The denominator in the above formula is always positive. The numerator may be positive or negative making r to be either positive or negative.

Assumptions in correlation analysis:

Correlation coefficient r is used under certain assumptions, they are

1. The variables under study are continuous random variables and they are normally distributed
2. The relationship between the variables is linear
3. Each pair of observations is unconnected with other pair (independent)

Properties

1. The correlation coefficient value ranges between -1 and $+1$.
2. The correlation coefficient is not affected by change of origin or scale or both.
3. If $r > 0$ it denotes positive correlation

$r < 0$ it denotes negative correlation between the two variables x and y .

$r = 0$ then the two variables x and y are not linearly correlated.(i.e)two variables are independent.

$r = +1$ then the correlation is perfect positive

$r = -1$ then the correlation is perfect negative.

Testing the significance of r

The significance of r can be tested by Student's t test. The test statistics is given by

$$t = \frac{|r|}{\sqrt{\frac{1-r^2}{n-2}}}$$

This t is distributed as Student's t distribution with (n-2) degrees of freedom.

The relationship between the variables is interpreted by the square of the correlation coefficient (r^2) which is called coefficient of determination. The value $1-r^2$ is called as coefficient of alienation. If r^2 is 0.72, it implies that on the basis of the samples 72% of the variation in one variable is caused by the variation of the other variable. The coefficient of determination is used to compare 2 correlation coefficients.

Problem

Compute Pearsons coefficient of correlation between plant height (cm) and yield (Kgs) as per the data given below:

Plant Height (cm)	39	65	62	90	82	75	25	98	36	78
Yield in Kgs	47	53	58	86	62	68	60	91	51	84

Solution

H_0 : The correlation coefficient r is not significant

H_1 : The correlation coefficient r is significant.

Level of significance 5%

From the data

n = 10

$$\sum x = 650 \quad \sum y = 660 \quad \sum xy = 45604 \quad \sum x^2 = 47648 \quad \sum y^2 = 45784$$

$$r = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sqrt{\sum x^2 - \frac{(\sum x)^2}{n}} \sqrt{\sum y^2 - \frac{(\sum y)^2}{n}}}$$

$$\begin{aligned}
 &= \frac{45604 - \frac{(650)(660)}{10}}{\sqrt{47648 - \frac{(650)^2}{10}} \sqrt{45784 - \frac{(660)^2}{10}}} \\
 &= \frac{45604 - 42900}{(73.47)(47.1)} = 0.7804
 \end{aligned}$$

Correlation coefficient is positively correlated.

Test Statistic

$$t = \frac{|r|}{\sqrt{\frac{1-r^2}{n-2}}} \sim (n-2) \text{ d.f.}$$

$$t = \frac{0.7804}{\sqrt{\frac{1-(0.7804)^2}{10-2}}} = 3.530$$

$$t_{\text{tab}} = t_{(10-2, 5\% \text{los})} = 2.306$$

Inference

$t > t_{\text{tab}}$, we reject null hypothesis.

\therefore The correlation coefficient r is significant. (i.e) there is a relation between plant height and yield.

Questions

1. Limits for correlation coefficient.

- (a) $-1 \leq r \leq 1$ (b) $0 \leq r \leq 1$
- (c) $-1 \leq r \leq 0$ (d) $1 \leq r \leq 2$

Ans: $-1 \leq r \leq 1$

2. The correlation coefficient is unaffected by change of

- (a) Origin (b) scale
- (c) Scale & origin (d) None of these

Ans: scale & origin

3. When $r = +1$, there is Perfect positive correlation.

Ans: True

4. Karl pearsons correlation coefficient is calculated only when the two variables are continuous.

Ans: True

5. The correlation between two variables is symmetric

Ans: True

6. The correlation between two variables is known as multiple correlation.

Ans: False

7. What is a scatter diagram? Mention its uses

8. Define correlation.

9. Explain the method how to calculate the Karl pearsons correlation coefficient?

10. Mention the properties of the correlation coefficient?

Lecture.13

Regression – definition – fitting of simple linear regression equation – testing the significance of the regression coefficient

Regression

Regression is the functional relationship between two variables and of the two variables one may represent cause and the other may represent effect. The variable representing cause is known as independent variable and is denoted by X. The variable X is also known as predictor variable or repressor. The variable representing effect is known as dependent variable and is denoted by Y. Y is also known as predicted variable. The relationship between the dependent and the independent variable may be expressed as a function and such functional relationship is termed as regression. When there are only two variables the functional relationship is known as simple regression and if the relation between the two variables is a straight line it is known as simple linear regression. When there are more than two variables and one of the variables is dependent upon others, the functional relationship is known as multiple regression. The regression line is of the form $y=a+bx$ where a is a constant or intercept and b is the regression coefficient or the slope. The values of 'a' and 'b' can be calculated by using the method of least squares. An alternate method of calculating the values of a and b are by using the formula:

The regression equation of y on x is given by $y = a + bx$

The regression coefficient of y on x is given by

$$b = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sum x^2 - \frac{(\sum x)^2}{n}}$$

and $a = \bar{y} - b \bar{x}$

The regression line indicates the average value of the dependent variable Y associated with a particular value of independent variable X.

Assumptions

1. The x's are non-random or fixed constants
2. At each fixed value of X the corresponding values of Y have a normal distribution about a mean.
3. For any given x, the variance of Y is same.
4. The values of y observed at different levels of x are completely independent.

Properties of Regression coefficients

1. The correlation coefficient is the geometric mean of the two regression coefficients
2. Regression coefficients are independent of change of origin but not of scale.
3. If one regression coefficient is greater than unit, then the other must be less than unit but not vice versa. ie. both the regression coefficients can be less than unity but both cannot be greater than unity, ie. if $b_1 > 1$ then $b_2 < 1$ and if $b_2 > 1$, then $b_1 < 1$.
4. Also if one regression coefficient is positive the other must be positive (in this case the correlation coefficient is the positive square root of the product of the two regression coefficients) and if one regression coefficient is negative the other must be negative (in this case the correlation coefficient is the negative square root of the product of the two regression coefficients). ie. if $b_1 > 0$, then $b_2 > 0$ and if $b_1 < 0$, then $b_2 < 0$.
5. If θ is the angle between the two regression lines then it is given by

$$\tan \theta = \frac{(1 - r^2)\sigma_x \sigma_y}{r(\sigma_x^2 + \sigma_y^2)}$$

Testing the significance of regression co-efficient

To test the significance of the regression coefficient we can apply either a t test or analysis of variance (F test). The ANOVA table for testing the regression coefficient will be as follows:

Sources of variation	d.f.	SS	MS	F
Due to regression	1	SS(b)	S_b^2	S_b^2 / S_e^2

Deviation from regression	n-2	SS(Y)-SS(b)	S_e^2	
Total	n-1	SS(Y)		

In case of t test the test statistic is given by

$$t = b / SE (b) \text{ where } SE (b) = s_e^2 / SS(X)$$

The regression analysis is useful in predicting the value of one variable from the given values of another variable. Another use of regression analysis is to find out the causal relationship between variables.

Uses of Regression

The regression analysis is useful in predicting the value of one variable from the given value of another variable. Such predictions are useful when it is very difficult or expensive to measure the dependent variable, Y. The other use of the regression analysis is to find out the causal relationship between variables. Suppose we manipulate the variable X and obtain a significant regression of variables Y on the variable X. Thus we can say that there is a causal relationship between the variable X and Y. The causal relationship between nitrogen content of soil and growth rate in a plant, or the dose of an insecticide and mortality of the insect population may be established in this way.

Example 1

From a paddy field, 36 plants were selected at random. The length of panicles(x) and the number of grains per panicle (y) of the selected plants were recorded. The results are given below. Fit a regression line y on x. Also test the significance (or) regression coefficient.

The length of panicles in cm (x) and the number of grains per panicle (y) of paddy plants.

S.No.	Y	X	S.No.	Y	X	S.No.	Y	X
1	95	22.4	13	143	24.5	25	112	22.9
2	109	23.3	14	127	23.6	26	131	23.9
3	133	24.1	15	92	21.1	27	147	24.8
4	132	24.3	16	88	21.4	28	90	21.2
5	136	23.5	17	99	23.4	29	110	22.2
6	116	22.3	18	129	23.4	30	106	22.7
7	126	23.9	19	91	21.6	31	127	23.0

8	124	24.0	20	103	21.4	32	145	24.0
9	137	24.9	21	114	23.3	33	85	20.6
10	90	20.0	22	124	24.4	34	94	21.0
11	107	19.8	23	143	24.4	35	142	24.0
12	108	22.0	24	108	22.5	36	111	23.1

Null Hypothesis H_0 : regression coefficient is not significant.

Alternative Hypothesis H_1 : regression coefficient is significant.

$$\sum y = 4174 \quad \sum y^2 = 496258 \quad \bar{y} = \frac{\sum y}{n} = 115.94$$

$$\sum x = 822.9 \quad \sum x^2 = 18876.83 \quad \bar{x} = \frac{\sum x}{n} = 22.86$$

$$\sum xy = 96183.4$$

$$SS(Y) = \sum y^2 - \frac{(\sum y)^2}{n} = 496258 - \frac{(4174)^2}{36} = 12305.8889$$

$$SS(X) = \sum x^2 - \frac{(\sum x)^2}{n} = 18876.83 - \frac{(822.9)^2}{36} = 66.7075$$

The regression line y on x is $\bar{y} = a + b \bar{x}$

$$b_1 = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sum x^2 - \frac{(\sum x)^2}{n}} = \frac{96183.4 - \frac{(822.9)(4174)}{36}}{66.7075} = 11.5837$$

$$\bar{y} = a + b \bar{x}$$

$$115.94 = a + (11.5837)(22.86)$$

$$a = 115.94 - 264.8034$$

$$a = -148.8633$$

The fitted regression line is $y = -148.8633 + 11.5837x$

$$SS(b) = \frac{\left(\sum xy - \frac{\sum x \sum y}{n} \right)^2}{\sum x^2 - \frac{(\sum x)^2}{n}} = \frac{(722.7167)^2}{66.7075} = 8950.8841$$

Anova Table

Sources of Variation	d.f.	SS	MSS	F
Regression	1	8950.8841	8950.8841	90.7093
Error	36-2=34	3355.0048	98.6766	
Total	35	12305.8889		

For t-test

$$t = \frac{b}{SE(b)} \sim t_{(n-2)} d.f$$

$$SE(b) = \sqrt{\frac{Se^2}{SS(X)}} = \sqrt{\frac{98.6776}{66.7075}} = 1.2162$$

$$t = \frac{11.5837}{1.2162} = 9.5245$$

Table Value:

$t_{(n-2)}$ d.f.= t_{34} d.f at 5% level=2.032

$t > t_{tab}$. we reject H_0 .

Hence t is significant.

Questions

1. When the correlation coefficient $r = +1$, then the two regression lines

- a) are perpendicular to each other
- b) coincide
- c) are parallel to each other
- d) none of these

Ans: coincide

2. If one regression coefficient is greater than unity then the other must be
- a) greater than unity
 - b) equal to unity
 - c) less than unity
 - d) none of these

Ans: less than unity

3. If the correlation between the two variables is positive the regression coefficient will be positive.

Ans: True

4. The Dependent variable is also called as predicted variable.

Ans: True

5. Correlation coefficient is the geometric mean of two regression coefficients.

Ans: True

6. Regression gives the functional relationship between two variables.

Ans: True

7. What is meant by Cause and effect?

8. State the properties of regression coefficient.

9. From the following data, find the regression equation

$$\sum X = 21, \sum Y = 20, \sum X^2 = 91, \sum XY = 74, n = 7$$

10. Explain how to fit the regression equation of y on x and test the significance of the regression coefficient.

Lecture.14

Design of experiments – basic concepts – treatment – experimental unit – experimental error - basic principle – replication, randomization and local control.

Design of Experiments

Choice of treatments, method of assigning treatments to experimental units and arrangement of experimental units in different patterns are known as designing an experiment. We study the effect of changes in one variable on another variable. For example how the application of various doses of fertilizer affects the grain yield. Variable whose change we wish to study is known as **response variable**. Variable whose effect on the response variable we wish to study is known as **factor**.

Treatment: Objects of comparison in an experiment are defined as treatments. Examples are Varieties tried in a trail and different chemicals.

Experimental unit: The object to which treatments are applied or basic objects on which the experiment is conducted is known as experimental unit.

Example: piece of land, an animal, etc

Experimental error: Response from all experimental units receiving the same treatment may not be same even under similar conditions. These variations in responses may be due to various reasons. Other factors like heterogeneity of soil, climatic factors and genetic differences, etc also may cause variations (known as extraneous factors). The variations in response caused by extraneous factors are known as **experimental error**.

Our aim of designing an experiment will be to minimize the experimental error.

Basic principles

To reduce the experimental error we adopt certain principles known as basic principles of experimental design.

The basic principles are 1) Replication, 2) Randomization and 3) Local control

Replication

Repeated application of the treatments is known as replication.

When the treatment is applied only once we have no means of knowing about the variation in the results of a treatment. Only when we repeat several times we can estimate the experimental error.

With the help of experimental error we can determine whether the obtained differences between treatment means are real or not. When the number of replications is increased, experimental error reduces.

Randomization

When all the treatments have equal chance of being allocated to different experimental units it is known as randomization.

If our conclusions are to be valid, treatment means and differences among treatment means should be estimated without any bias. For this purpose we use the technique of randomization.

Local Control

Experimental error is based on the variations from experimental unit to experimental unit. This suggests that if we group the homogenous experimental units into blocks, the experimental error will be reduced considerably. Grouping of homogenous experimental units into blocks is known as local control of error.

In order to have valid estimate of experimental error the principles of replication and randomization are used.

In order to reduce the experimental error, the principles of replication and local control are used.

In general to have precise, valid and accurate result we adopt the basic principles.

Questions

1. For valid conclusions we should have

- | | |
|-----------------------|---------------------|
| (a) Unbiased estimate | (b) biased estimate |
| (c) random estimate | (d) none of these |

Ans: Unbiased estimate

2. Response variable is also called as
- (a) Independent variable
 - (b) dependent variable
 - (c) treatment
 - (d) error

Ans: dependent variable

3. The genetic differences of varieties are termed as extraneous factors.

Ans: True

4. Repetition of the treatment is known as replication.

Ans: True

5. Replication will increase the error.

Ans: False

6. Basic principles are adopted to reduce the experimental error.

Ans: True

7. What is experimental error?

8. Define treatment and experimental unit.

9. What is meant by designing an experiment?

10. Explain the basic principles and its uses?

Lecture.15

Completely randomized design – description – layout – analysis – advantages and disadvantages

Completely Randomized Design (CRD)

CRD is the basic single factor design. In this design the treatments are assigned completely at random so that each experimental unit has the same chance of receiving any one treatment. But CRD is appropriate only when the experimental material is homogeneous. As there is generally large variation among experimental plots due to many factors CRD is not preferred in field experiments.

In laboratory experiments and greenhouse studies it is easy to achieve homogeneity of experimental materials and therefore CRD is most useful in such experiments.

Layout of a CRD

Completely randomized Design is the one in which all the experimental units are taken in a single group which are homogeneous as far as possible.

The randomization procedure for allotting the treatments to various units will be as follows.

Step 1: Determine the total number of experimental units.

Step 2: Assign a plot number to each of the experimental units starting from left to right for all rows.

Step 3: Assign the treatments to the experimental units by using random numbers.

The statistical model for CRD with one observation per unit

$$Y_{ij} = \mu + t_i + e_{ij}$$

μ = overall mean effect

t_i = true effect of the i^{th} treatment

e_{ij} = error term of the j^{th} unit receiving i^{th} treatment

The arrangement of data in CRD is as follows:

	Treatments				
	T ₁	T ₂	T _i	T _K	
	y ₁₁	y ₂₁	y _{i1}	Y _{K1}	
	y ₁₂	y ₂₂	y _{i2}	Y _{K2}	
	y _{1r1}	y _{2r2}	y _{iri}	Y _{k rk}	
Total	Y₁	Y₂	Y_i	T_k	GT

(GT – Grand total)

The null hypothesis will be

H₀: μ₁ = μ₂ = = μ_k or There is no significant difference between the treatments

And the alternative hypothesis is

H₁: μ₁ ≠ μ₂ ≠ ≠ μ_k. There is significant difference between the treatments

The different steps in forming the analysis of variance table for a CRD are:

$$1. \quad C.F = \frac{(GT)^2}{n}$$

n= Total number of observations

$$2. \quad \text{Total SS} = \text{TSS} = \sum_{i=1}^k \sum_{j=1}^r y_{ij}^2 - C.F$$

$$3. \quad \text{Treatment SS} = \text{TrSS} = \frac{Y_1^2}{r_1} + \frac{Y_2^2}{r_2} + \dots + \frac{Y_k^2}{r_k} - C.F$$

$$= \sum_{i=1}^k \frac{Y_i^2}{r_i} - C.F$$

$$4. \quad \text{Error SS} = \text{ESS} = \sum_{i=1}^k \sum_{j=1}^r y_{ij}^2 - \sum_{i=1}^k \frac{Y_i^2}{r_i}$$

$$= \text{TSS} - \text{TrSS}$$

5. Form the following ANOVA table and calculate F value.

Source of variation	d.f.	SS	MS	F
Treatments	t-1	TrSS	TrMS = $\frac{TrSS}{t-1}$	$\frac{TrMS}{EMS}$
Error	n-t	ESS	EMS = $\frac{ESS}{n-t}$	
Total	n-1	TSS		

6. Compare the calculated F with the critical value of F corresponding to treatment degrees of freedom and error degrees of freedom so that acceptance or rejection of the null hypothesis can be determined.

7. If null hypothesis is rejected that indicates there is significant differences between the different treatments.

8. Calculate C D value.

$$C.D. = SE(d). t$$

$$\text{where S.E(d)} = \sqrt{EMS\left(\frac{1}{r_i} + \frac{1}{r_j}\right)}$$

r_i = number of replications for treatment i

r_j = number of replications for treatment j and

t is the critical t value for error degrees of freedom at specified level of significance, either 5% or 1%.

Advantages of a CRD

1. Its layout is very easy.
2. There is complete flexibility in this design i.e. any number of treatments and replications for each treatment can be tried.
3. Whole experimental material can be utilized in this design.
4. This design yields maximum degrees of freedom for experimental error.
5. The analysis of data is simplest as compared to any other design.
6. Even if some values are missing the analysis can be done.

Disadvantages of a CRD

1. It is difficult to find homogeneous experimental units in all respects and hence CRD is seldom suitable for field experiments as compared to other experimental designs.
2. It is less accurate than other designs.

Questions

1. CRD can be used with
(a) Equal replication (b) unequal replication
(c) Equal and unequal replication (d) single replication

Ans: Equal and unequal replication

2. When there are 5 treatments each replicated 4 times the total number of experimental plots will be
(a) 5 (b) 4 (c) 9 (d) 20

Ans: 20

3. In CRD the error degrees of freedom is $rt-1$.

Ans: True

4. CRD can be adopted only when the experimental material is homogenous.

Ans: True

5. CRD is a single factor experiment.

Ans: True

6. In CRD the total sum of squares is divided into treatment sum of squares, Replication sum of squares and error sum of squares.

Ans: False

7. Mention any two advantages of CRD?

8. When the treatments are large in a CRD what will happen to the precision of the experiment?

9. Explain the Layout of the CRD?

10. Explain the Layout of the CRD?

Lecture.16

Randomized blocks design – description – layout – analysis – advantages and disadvantages

Randomized Blocks Design (RBD)

When the experimental material is heterogeneous, the experimental material is grouped into homogenous sub-groups called blocks. As each block consists of the entire set of treatments a block is equivalent to a replication.

If the fertility gradient runs in one direction say from north to south or east to west then the blocks are formed in the opposite direction. Such an arrangement of grouping the heterogeneous units into homogenous blocks is known as randomized blocks design. Each block consists of as many experimental units as the number of treatments. The treatments are allocated randomly to the experimental units within each block independently such that each treatment occurs once. The number of blocks is chosen to be equal to the number of replications for the treatments.

The analysis of variance model for RBD is

$$Y_{ij} = \mu + t_i + r_j + e_{ij}$$

where

μ = the overall mean

t_i = the i^{th} treatment effect

r_j = the j^{th} replication effect

e_{ij} = the error term for i^{th} treatment and j^{th} replication

Analysis of RBD

The results of RBD can be arranged in a two way table according to the replications (blocks) and treatments.

There will be $r \times t$ observations in total where r stands for number of replications and t for number of treatments. .

The data are arranged in a two way table form by representing treatments in rows and replications in columns.

Treatment	Replication					Total
	1	2	3	r	
1	y11	y12	y13	y1r	T1
2	y21	y22	y23	y2r	T2
3	y31	y32	y33	y3r	T3
t	yt1	yt2	yt3	ytr	Tt
Total	R1	R2	R3		Rr	G.T

In this design the total variance is divided into three sources of variation viz., between replications, between treatments and error

$$CF = \frac{(GT)^2}{rt}$$

$$\text{Total SS} = TSS = \sum \sum y_{ij}^2 - CF$$

$$\text{Replication SS} = RSS = \frac{1}{t} \sum R_j^2 - CF$$

$$\text{Treatments SS} = TrSS = \frac{1}{r} \sum T_i^2 - CF$$

$$\text{Error SS} = ESS = \text{Total SS} - \text{Replication SS} - \text{Treatment SS}$$

The skeleton ANOVA table for RBD with t treatments and r replications

Sources of variation	d.f.	SS	MS	F Value
Replication	r-1	RSS	RMS	RM S/ EM S
Treatment	t-1	TrSS	TrMS	TrMS/EMS
Error	(r-1)(t-1)	ESS	EMS	
Total	rt - 1	TSS		

$$CD = SE(d) \cdot t \quad \text{where } S.E(d) = \sqrt{\frac{2EMS}{r}}$$

t = critical value of t for a specified level of significance and error degrees of freedom

Based on the CD value the bar chart can be drawn.

From the bar chart conclusion can be written.

Advantages of RBD

The precision is more in RBD. The amount of information obtained in RBD is more as compared to CRD. RBD is more flexible. Statistical analysis is simple and easy. Even if some values are missing, still the analysis can be done by using missing plot technique.

Disadvantages of RBD

When the number of treatments is increased, the block size will increase. If the block size is large maintaining homogeneity is difficult and hence when more number of treatments is present this design may not be suitable.

Questions

1. RBD can be used with
- | | |
|-----------------------------------|-------------------------|
| (a) Equal replication | (b) unequal replication |
| (c) Equal and unequal replication | (d) single replication |

Ans: Equal replication

2. When there are 5 treatments each replicated 4 times the total number of experimental plots will be
- (a) 5 (b) 4 (c) 9 (d) 20

Ans: 20

3. In RBD the error degrees of freedom is $(r-1)(t-1)$.

Ans: True

4. RBD can be adopted when the experimental material is heterogeneous.

Ans: True

5. In RBD the blocking is done in one direction.

Ans: True

6. In RBD the total sum of squares is divided into treatment sum of squares, Replication sum of squares and error sum of squares.

Ans: True

7. Mention any two advantages of RBD?

8. Furnish the ANOVA model for RBD

9. Explain the Layout of the RBD?

10. Explain the computational procedure of RBD?

Lecture.17

Latin square design – description – layout – analysis – advantages and disadvantages.

Latin Square Design

When the experimental material is divided into rows and columns and the treatments are allocated such that each treatment occurs only once in each row and each column, the design is known as L S D.

In LSD the treatments are usually denoted by A B C D etc.

For a 5 x 5 LSD the arrangements may be

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>B</i>	<i>A</i>	<i>E</i>	<i>C</i>	<i>D</i>
<i>C</i>	<i>D</i>	<i>A</i>	<i>E</i>	<i>B</i>
<i>D</i>	<i>E</i>	<i>B</i>	<i>A</i>	<i>C</i>
<i>E</i>	<i>C</i>	<i>D</i>	<i>B</i>	<i>A</i>
Square 1				

	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>B</i>	<i>A</i>	<i>D</i>	<i>E</i>	<i>C</i>
<i>C</i>	<i>E</i>	<i>A</i>	<i>B</i>	<i>D</i>
<i>D</i>	<i>C</i>	<i>E</i>	<i>A</i>	<i>B</i>
<i>E</i>	<i>D</i>	<i>B</i>	<i>C</i>	<i>A</i>
Square 2				

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>A</i>
<i>C</i>	<i>D</i>	<i>E</i>	<i>A</i>	<i>B</i>
<i>D</i>	<i>E</i>	<i>A</i>	<i>B</i>	<i>C</i>
<i>E</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>
Square 3				

Analysis

The ANOVA model for LSD is

$$Y_{ijk} = \mu + r_i + c_j + t_k + e_{ijk}$$

r_i is the i^{th} row effect

c_j is the j^{th} column effect

t_k is the k^{th} treatment effect and

e_{ijk} is the error term

The analysis of variance table for LSD is as follows:

Sources of Variation	d.f.	S S	M S	F
Rows	t-1	RSS	RMS	RMS/EMS
Columns	t-1	CSS	CMS	CMS/EMS
Treatments	t-1	TrSS	TrMS	TrMS/EMS
Error	(t-1)(t-2)	ESS	EMS	
Total	t^2-1	TSS		

F table value

$F_{[t-1,(t-1)(t-2)]}$ degrees of freedom at 5% or 1% level of significance

Steps to calculate the above Sum of Squares are as follows:

$$\text{Correction Factor (CF)} = \frac{(GT)^2}{(t)^2}$$

$$\text{Total Sum of Squares (TSS)} = \sum (y_{ijk})^2 - CF$$

$$\text{Row sum of squares (RSS)} = \frac{1}{t} \sum_{i=1}^t (R_i)^2 - CF$$

$$\text{Column sum of squares (CSS)} = \frac{1}{t} \sum_{j=1}^t (C_j)^2 - CF$$

$$\text{Treatment sum of squares (TrSS)} = \frac{1}{t} \sum_{k=1}^t (T_k)^2 - CF$$

$$\text{Error Sum of Squares} = \text{TSS} - \text{RSS} - \text{CSS} - \text{TrSS}$$

These results can be summarized in the form of analysis of variance table.

Calculation of SE, SE (d) and CD values

$$SE = \sqrt{\frac{EMS}{r}}$$

where r is the number of rows

$$SE(d) = \sqrt{2} SE .$$

$$CD = SE (d). t$$

where t = table value of t for a specified level of significance and error degrees of freedom

Using CD value the bar chart can be drawn and the conclusion may be written.

Advantages

- LSD is more efficient than RBD or CRD. This is because of double grouping that will result in small experimental error.
- When missing values are present, missing plot technique can be used and analysed.

Disadvantages

- This design is not as flexible as RBD or CRD as the number of treatments is limited to the number of rows and columns. LSD is seldom used when the number of treatments is more than 12. LSD is not suitable for treatments less than five.

Because of the limitations on the number of treatments, LSD is not widely used in agricultural experiments.

Note: The number of sources of variation is two for CRD, three for RBD and four for LSD.

Questions

1. In a Latin Square design the number of rows will be equal to
- | | |
|------------------------|--|
| a) No. of columns | b) No. of Treatments |
| c) No. of Replications | d) No. of Columns & Number of Treatments |

Ans: No. of Columns & Number of Treatments

2. In a Latin Square design with 5 treatments the number of experimental units will be equal to
- | | | | |
|-------|-------|-------|-------|
| a) 25 | b) 20 | c) 24 | d) 36 |
|-------|-------|-------|-------|

Ans: 25

3. If the number of experimental units is 36 then the number of rows will be equal to 6.

Ans: True

4. The error degrees of freedom in LSD with t treatments will be $(t-1)(t-2)$.

Ans: True

5. If the experimental material is homogeneous then LSD can be adopted.

Ans: False

6. In a LSD each treatment should occur only once in each row and each column.

Ans: True

7. Furnish the ANOVA model for LSD.
8. What is a Latin Square Design?
9. State the advantages and disadvantages of LSD.
10. Explain the computational procedure of LSD?

Lecture.18

Factorial experiments – factor and levels – types – symmetrical and asymmetrical – simple, main and interaction effects – advantages and disadvantages

Factorial Experiments

When two or more number of factors are investigated simultaneously in a single experiment such experiments are called as factorial experiments.

Terminologies

1. **Factor:** Factor refers to a set of related treatments. We may apply of different doses of nitrogen to a crop. Hence nitrogen irrespective of doses is a factor.
2. **Levels of a factor:** Different states or components making up a factor are known as the levels of that factor. eg different doses of nitrogen.

Types of factorial Experiment

A factorial experiment is named based on the number of factors and levels of factors. For example, when there are 3 factors each at 2 levels the experiment is known as 2 X 2 X 2 or 2³ factorial experiments.

If there are 2 factors each at 3 levels then it is known as 3 X 3 or 3² factorial experiment.

- In general if there are n factors each with p levels then it is known as pⁿ factorial experiment.
- For varying number of levels the arrangement is described by the product. For example, an experiment with 3 factors each at 2 levels, 3 levels and 4 levels respectively then it is known as 2 X 3 X 4 factorial experiment.
- If all the factors have the same number of levels the experiment is known as symmetrical factorial otherwise it is called as mixed factorial.
- Factors are represented by capital letters. Treatment combinations are usually by small letters.
- For example, if there are 2 varieties v0 and v1 and 2 dates of sowing d0 and d1 the treatment combinations will be
- v0d0, v1d0, v1d0 and v1d1.

Simple and Main Effects

Simple effect of a factor is the difference between its responses for a fixed level of other factors.

Main effect is defined as the average of the simple effects.

Interaction is defined as the dependence of factors in their responses. Interaction is measured as the mean of the differences between simple effects.

Advantages

1. In such type of experiments we study the individual effects of each factor and their interactions.
2. In factorial experiments a wide range of factor combinations are used.
3. Factorial approach will result in considerable saving of the experimental resources, experimental material and time.

Disadvantages

1. When number of factors or levels of factors or both are increased, the number of treatment combinations increases. Consequently block size increases. If block size increases it may be difficult to maintain homogeneity of experimental material. This will lead to increase in experimental error and loss of precision in the experiment.
2. All treatment combinations are to be included for the experiment irrespective of its importance and hence this results in wastage of experimental material and time.
3. When many treatment combinations are included the execution of the experiment and statistical analysis become difficult.

Questions

1. In a factorial experiment the minimum number of factors will be
a) One b) two c) three d) none of the above

Ans: two

2. With factor A which has two levels and factor B which has three levels the number of combinations will be
a) 5 b) 6 c) 2 d) 3

Ans: 6

3. In a factorial experiments Factor refers to a set of related treatments.

Ans: True

4. A three 3 x3 factorial experiment can also be written as 3^2 factorial experiment.

Ans: True

5. In a factorial experiment factors are always represented by small letters.

Ans: False

6. If all the factors have the same number of levels the experiment is known as symmetrical factorial.

Ans: True

7. What is a main effect?

8. What is an interaction effect?

9. What are the advantages and disadvantages of Factorial experiments?

10. Write down the treatment combinations for a factorial experiment with varieties as one factor and n levels as the second factor. The levels for the varieties are 4 (V_0, V_1, V_2 & V_3) for n levels 4 (n_0, n_1, n_2 & n_3)

Lecture.19

2² Factorial Experiments in RBD – lay out – analysis

2² Factorial Experiments in RBD

2² factorial experiment means two factors each at two levels. Suppose the two factors are A and B and both are tried with two levels the total number of treatment combinations will be four i.e. a0b0, a0b1, a1b0 and a1b1.

The allotment of these four treatment combinations will be as allotted in RBD. That is each block is divided into four experimental units. By using the random numbers these four combinations are allotted at random for each block separately.

The analysis of variance table for two factors A with a levels and B with b levels with r replications tried in RBD will be as follows:

Sources of Variation	d.f.	SS	MS	F
Replications	r-1	RSS	RMS	
Factor A	a-1	ASS	AMS	AMS / EMS
Factor B	b-1	BSS	BMS	BMS / EMS
AB (interaction)	(a-1)(b-1)	ABSS	ABMS	ABMS / EMS
Error	(r-1)(ab-1)	ESS	EMS	
Total	rab-1	TSS		

As in the previous designs calculate the replication totals to calculate the RSS, TSS in the usual way. To calculate ASS, BSS and ABSS, form a two way table A X B by taking the levels of A in rows and levels of B in the columns. To get the values in this table the missing factor is replication. That is by adding over replication we can form this table.

$$CF = \frac{(G T)^2}{r \times a \times b}$$

$$RSS = \frac{\sum_{i=1}^r R_i^2}{a \times b} - CF$$

A X B Two way table

B A	b ₀	b ₁	Total
a ₀	a ₀ b ₀	a ₀ b ₁	A ₀
a ₁	a ₁ b ₀	a ₁ b ₁	A ₁
Total	B ₀	B ₁	Grand Total

$$ASS = \frac{A_0^2 + A_1^2}{b \times r} - CF$$

$$BSS = \frac{B_0^2 + B_1^2}{a \times r} - CF$$

$$ABSS = \frac{(a_0 b_0)^2 + (a_0 b_1)^2 + (a_1 b_0)^2 + (a_1 b_1)^2}{r} - CF - ASS - BSS$$

$$ESS = TSS - RSS - ASS - BSS - ABSS$$

By substituting the above values in the ANOVA table corresponding to the columns sum of squares, the mean squares and F value can be calculated.

Questions

1. 2² factorial experiment means two factors each at
 a) two levels b) three levels c) four levels d) one level

Ans: two levels

2. If the total number of combinations are four then each block is divided into _____ experimental units

- a) Two b) three c) four d) none of the above

Ans: four

3. The two factors are A and B and both are tried with three levels the total number of combinations will be nine.

Ans: True

4. The error degrees of freedom for an experiment with one factor at 2 levels and another at three levels and 3 replications will be 10.

Ans: True

5. In a two factor experiment the highest order interaction will be the two factor interaction.

Ans: True

6. In a factorial experiment with two factors the treatment sum of squares will be split up into factor 1, factor 2 and factor 1x factor 2 interaction sum of squares.

Ans: True

7. How to calculate the replication sum of squares in a two factor experiment.

8. In a factorial experiment what is the total number of experimental units when there are 3 replications, 4 levels for factor A and 3 levels for factor B

9. Furnish the ANOVA table for a two factor experiment with factor A at a levels, Factor B at b levels and the number of replication r.

10. Explain the procedure of forming a A X B two way table and calculating the factor ASS, BSS and A X BSS.

Lecture.20

2³ factorial experiments in RBD – lay out – analysis

2³ Factorial Experiment in RBD

2³ factorial experiment mean three factors each at two levels. Suppose the three factors are A, B and C are tried with two levels the total number of combinations will be eight i.e. a0b0c0, a0b0c1, a0b1c0, a0b1c1, a1b0c0, a1b0c1, a1b1c0 and a1b1c1.

The allotment of these eight treatment combinations will be as allotted in RBD. That is each block is divided into eight experimental units. By using the random numbers these eight combinations are allotted at random for each block separately.

The analysis of variance table for three factors A with a levels, B with b levels and C with c levels with r replications tried in RBD will be as follows:

Sources of Variation	d.f.	SS	MS	F
Replications	r-1	RSS	RMS	
Factor A	a-1	ASS	AMS	AMS / EMS
Factor B	b-1	BSS	BMS	BMS / EMS
Factor C	c-1	CSS	CMS	CMS / EMS
AB	(a-1)(b-1)	ABSS	ABMS	ABMS / EMS
AC	(a-1)(c-1)	ACSS	ACMS	ACMS / EMS
BC	(b-1)(c-1)	BCSS	BCMS	BCMS / EMS
ABC	(a-1)(b-1)(c-1)	ABCSS	ABCMS	ABCMS / EMS
Error	(r-1)(abc-1)	ESS	EMS	
Total	rabc-1	TSS		

Analysis

1. Arrange the results as per treatment combinations and replications.

Treatment combination	Replication				Treatment Total
	R1	R2	R3	...	
a0b0c0					T1
a0b0c1					T2
a0b1c0					T3
a0b1c1					T4
a1b0c0					T5
a1b0c1					T6
a1b1c0					T7
a1b1c1					T8

As in the previous designs calculate the replication totals to calculate the CF, RSS, TSS, overall TrSS in the usual way. To calculate ASS, BSS, CSS, ABSS, ACSS, BCSS and ABCSS, form three two way tables A X B, AXC and BXC.

AXB two way table can be formed by taking the levels of A in rows and levels of B in the columns. To get the values in this table the missing factor is replication. That is by adding over replication we can form this table.

A X B Two way table

B A	b ₀	b ₁	Total
a ₀	a ₀ b ₀	a ₀ b ₁	A ₀
a ₁	a ₁ b ₀	a ₁ b ₁	A ₁

Total	B ₀	B ₁	Grand Total
-------	----------------	----------------	-------------

$$ASS = \frac{A_0^2 + A_1^2}{b \times c \times r} - CF$$

$$BSS = \frac{B_0^2 + B_1^2}{a \times c \times r} - CF$$

$$ABSS = \frac{(a_0 b_0)^2 + (a_0 b_1)^2 + (a_1 b_0)^2 + (a_1 b_1)^2}{c \times r} - CF - ASS - BSS$$

A X C two way table can be formed by taking the levels of A in rows and levels of C in the columns

A X C Two way table

C \ A	c ₀	c ₁	Total
a ₀	a ₀ c ₀	a ₀ c ₁	A ₀
a ₁	a ₁ c ₀	a ₁ c ₁	A ₁
Total	C ₀	C ₁	Grand Total

$$CSS = \frac{C_0^2 + C_1^2}{a \times b \times r} - CF$$

$$ACSS = \frac{(a_0 c_0)^2 + (a_0 c_1)^2 + (a_1 c_0)^2 + (a_1 c_1)^2}{b \times r} - CF - ASS - CSS$$

B X C two way table can be formed by taking the levels of B in rows and levels of C in the columns

B X C Two way table

C \ B	c ₀	c ₁	Total
-------	----------------	----------------	-------

b ₀	b ₀ c ₀	b ₀ c ₁	B ₀
b ₁	b ₁ c ₀	b ₁ c ₁	B ₁
Total	C ₀	C ₁	Grand Total

$$BCSS = \frac{(b_0c_0)^2 + (b_0c_1)^2 + (b_1c_0)^2 + (b_1c_1)^2}{a \times r} - CF - BSS - CSS$$

$$ABCSS = \frac{(a_0b_0c_0)^2 + (a_0b_0c_1)^2 + (a_0b_1c_0)^2 + (a_0b_1c_1)^2 + (a_1b_0c_0)^2 + (a_1b_0c_1)^2 + (a_1b_1c_0)^2 + (a_1b_1c_1)^2}{r}$$

$$-CF-ASS-BSS-CSS-ABSS-ACSS-BCSS$$

$$ESS = TSS-RSS- ASS-BSS-CSS-ABSS-ACSS-BCSS-ABCSS$$

By substituting the above values in the ANOVA table corresponding to the columns sum of squares, the mean squares and F value can be calculated.

Questions

1. 2³ factorial experiment means two factors each at

- a) two levels b) three levels c) four levels d) one level

Ans: three levels

2. If the total number of combinations are eight then each block is divided into _____ experimental units

- a) Two b) three c) four d) eight

Ans: eight

3. If the three factors are A, B and C are tried with three levels the total number of combinations will be twenty seven.

Ans: True

4. The error degrees of freedom for an experiment with one factor at 2 levels, second factor at 3 levels and the third at 3 levels and 3 replications will be 34.

Ans: True

5. In a three factor experiment the highest order interaction will be the two factor interaction.

Ans: False

6. In a factorial experiment with three factors the treatment sum of squares will be split up into factor 1, factor 2 and factor 1 x factor 2 interaction sum of squares.

Ans: False

7. How to calculate the three factor interaction sum of squares in a three factor experiment.

8. In a factorial experiment what is the total number of experimental units when there are 3 replications, 4 levels for factor A and 3 levels for factor B and C.

9. Furnish the ANOVA table for a three factor experiment with factor A at a levels, Factor B at b levels and factor C at c levels and the number of replication r.

10. Explain the procedure of forming a B X C two way table and calculating the factor BSS, CSS and B X CSS.

Lecture.21

Split plot design – layout – ANOVA Table

Split-plot Design

In field experiments certain factors may require larger plots than for others. For example, experiments on irrigation, tillage, etc requires larger areas. On the other hand experiments on fertilizers, etc may not require larger areas. To accommodate factors which require different sizes of experimental plots in the same experiment, split plot design has been evolved.

In this design, larger plots are taken for the factor which requires larger plots. Next each of the larger plots is split into smaller plots to accommodate the other factor. The different treatments are allotted at random to their respective plots. Such arrangement is called split plot design.

In split plot design the larger plots are called main plots and smaller plots within the larger plots are called as sub plots. The factor levels allotted to the main plots are main plot treatments and the factor levels allotted to sub plots are called as sub plot treatments.

Layout and analysis of variance table

First the main plot treatment and sub plot treatment are usually decided based on the needed precision. The factor for which greater precision is required is assigned to the sub plots.

The replication is then divided into number of main plots equivalent to main plot treatments. Each main plot is divided into subplots depending on the number of sub plot treatments. The main plot treatments are allocated at random to the main plots as in the case of RBD. Within each main plot the sub plot treatments are allocated at random as in the case of RBD. Thus randomization is done in two stages. The same procedure is followed for all the replications independently.

The analysis of variance will have two parts, which correspond to the main plots and sub-plots. For the main plot analysis, replication X main plot treatments table is

formed. From this two-way table sum of squares for replication, main plot treatments and error (a) are computed. For the analysis of sub-plot treatments, main plot X sub-plot treatments table is formed. From this table the sums of squares for sub-plot treatments and interaction between main plot and sub-plot treatments are computed. Error (b) sum of squares is found out by residual method. The analysis of variance table for a split plot design with m main plot treatments and s sub-plot treatments is given below.

Analysis of variance for split plot with factor A with m levels in main plots and factor B with s levels in sub-plots will be as follows:

Sources of Variation	d.f.	SS	MS	F
Replication	$r-1$	RSS	RMS	RMS/EMS (a)
A	$m-1$	ASS	AMS	AMS/EMS (a)
Error (a)	$(r-1)(m-1)$	ESS (a)	EMS (a)	
B	$s-1$	BSS	BMS	BMS/EMS (b)
AB	$(m-1)(s-1)$	ABSS	ABMS	ABMS/EMS (b)
Error (b)	$m(r-1)(s-1)$	ESS (b)	EMS (b)	
Total	$rms - 1$	TSS		

Analysis

Arrange the results as follows

Treatment Combination	Replication				Total
	R1	R2	R3	...	
A0B0	a0b0	a0b0	a0b0	...	T00
A0B1	a0b1	a0b1	a0b1	...	T01
A0B2	a0b2	a0b2	a0b2	...	T02
Sub Total	A01	A02	A03	...	T0
A1B0	a1b0	a1b0	a1b0	...	T10
A1B1	a1b1	a1b1	a1b1	...	T11
A1B2	a1b2	a1b2	a1b2	...	T12
Sub Total	A11	A12	A13	...	T1
.
.
.
Total	R1	R2	R3	...	G.T

$$\text{Compute CF} = \frac{(G T)^2}{r \times m \times s}$$

$$\text{TSS} = [(a_0b_0)^2 + (a_0b_1)^2 + (a_0b_2)^2 + \dots] - \text{CF}$$

Form A x R Table and calculate RSS, ASS and Error (a) SS

Treatment	Replication				Total
	R1	R2	R3	...	
A0	A01	A02	A03	...	T0
A1	A11	A12	A13	...	T1
A2	A21	A22	A23	...	T2
.
.
.
Total	R1	R2	R3	...	GT

$$\text{RSS} = \left(\frac{R1^2 + R2^2 + R3^2 + \dots}{m \cdot s} \right) - \text{CF}$$

$$\text{ASS} = \left(\frac{T0^2 + T1^2 + T2^2 + \dots}{r \cdot s} \right) - \text{CF}$$

$$\text{A x R table SS} = \left(\frac{A01^2 + A02^2 + A03^2 + \dots}{b} \right) - \text{CF}$$

Error (a) SS = A x R TSS - RASS - ASS.

Form A x B Table and calculate BSS, Ax B SSS and Error (b) SS

Treatment	Replication				Total
	B0	B1	B2	...	
A0	T00	T01	T02	...	T0
A1	T10	T11	T12	...	T1
A2	T20	T21	T22	...	T2
.
.
.
Total	C0	C1	C2	...	GT

$$\text{BSS} = \left(\frac{C0^2 + C1^2 + C2^2 + \dots}{r \cdot m} \right) - \text{CF}$$

$$A \times B \text{ table SS} = \left(\frac{T_0^2 + T_1^2 + T_3^2 + \dots}{r} \right) - CF$$

ABSS= A x B Table SS – ASS- ABSS

Error (b) SS= Table SS-ASS-BSS-ABSS –Error (a) SS.

Then complete the ANOVA table.

Questions

1. To accommodate factors which require different sizes of experimental plots in the same experiment _____ design has been evolved

- a) Split plot b) CRD c) RBD d) LSD

Ans: Split plot

2. The number of error terms in a split plot design is

- a) One b) two c) three d) none of these

Ans: two

3. The plot size for the subplot treatment will be small when compared to the plot size of the main plot treatments.

Ans: True

4. The precision for the sub plot treatments is more when compared to the main plot treatments.

Ans: True

5. The plot sizes for the main plot treatment and subplot treatment are not same.

Ans: True

6. The plot size for subplot and interaction are same in split plot design.

Ans: True

7. When will you adopt split plot design?

8. What is error (a) in a split plot design?

9. Furnish the skeleton ANOVA table with 3 replications, 4 main plot treatments and 3 subplot treatments.

10. Furnish the layout of a split plot design.

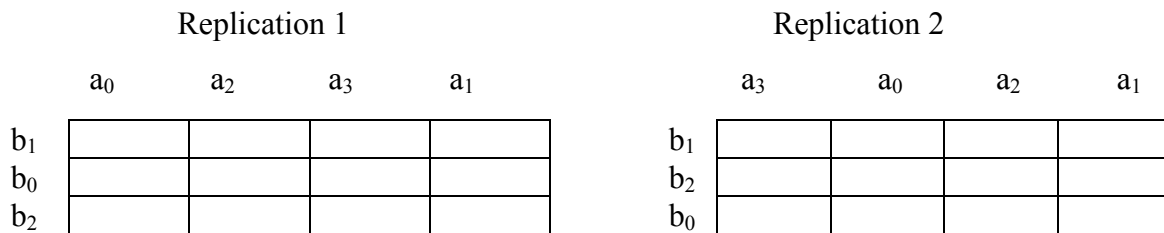
Lecture.22

Strip plot design – layout – ANOVA Table

Strip Plot Design

This design is also known as **split block design**. When there are two factors in an experiment and both the factors require large plot sizes it is difficult to carryout the experiment in split plot design. Also the precision for measuring the interaction effect between the two factors is higher than that for measuring the main effect of either one of the two factors. Strip plot design is suitable for such experiments.

In strip plot design each block or replication is divided into number of vertical and horizontal strips depending on the levels of the respective factors.



In this design there are plot sizes.

1. Vertical strip plot for the first factor – vertical factor
2. Horizontal strip plot for the second factor – horizontal factor
3. Interaction plot for the interaction between 2 factors

The vertical strip and the horizontal strip are always perpendicular to each other. The interaction plot is the smallest and provides information on the interaction of the 2 factors. Thus we say that interaction is tested with more precision in strip plot design.

Analysis

The analysis is carried out in 3 parts.

1. Vertical strip analysis
2. Horizontal strip analysis

3. Interaction analysis

Suppose that A and B are the vertical and horizontal strips respectively. The following two way tables, viz., A X Rep table, B X Rep table and A X B table are formed. From A X Rep table, SS for Rep, A and Error (a) are computed. From B X Rep table, SS for B and Error (b) are computed. From A X B table, A X B SS is calculated.

When there are r replications, a levels for factor A and b levels for factor B, then the ANOVA table is

X	d.f.	SS	MS	F
Replication	(r-1)	RSS	RMS	RMS/EMS (a)
A	(a-1)	ASS	AMS	AMS/EMS (a)
Error (a)	(r-1) (a-1)	ESS (a)	EMS (a)	
B	(b-1)	BSS	BMS	BMS/EMS (b)
Error (b)	(r-1) (b-1)	ESS (b)	EMS (b)	
AB	(a-1) (b-1)	ABSS	ABMS	ABMS/EMS (c)
Error (c)	(r-1) (a-1) (b-1)	E SS (c)	EMS (c)	
Total		(rab - 1)	TSS	

Analysis

Arrange the results as follows:

Treatment Combination	Replication				Total
	R1	R2	R3	...	
A0B0	a0b0	a0b0	a0b0	...	T00
A0B1	a0b1	a0b1	a0b1	...	T01
A0B2	a0b2	a0b2	a0b2	...	T02
Sub Total	A01	A02	A03	...	T0
A1B0	a1b0	a1b0	a1b0	...	T10
A1B1	a1b1	a1b1	a1b1	...	T11
A1B2	a1b2	a1b2	a1b2	...	T12
Sub Total	A11	A12	A13	...	T1
.
.
.
Total	R1	R2	R3	...	G.T

$$\text{Compute CF} = \frac{(G T)^2}{r \times m \times s}$$

TSS = [(a0b0)² +

Treatment	Replication				Total
	R1	R2	R3	...	
B0	B01	B02	B03	...	T0
B1	B11	B12	B13	...	T1
B2	B21	B22	B23	...	T2
.
.
.
Total	R1	R2	R3	...	GT

(a0b1)²+(a0b2)²+...]-CF

1) Vertical Strip Analysis

Form A x R Table and calculate RSS, ASS and Error(a) SS

Treatment	Replication				Total
	R1	R2	R3	...	
A0	A01	A02	A03	...	T0
A1	A11	A12	A13	...	T1
A2	A21	A22	A23	...	T2
.
.
.
Total	R1	R2	R3	...	GT

$$RSS = \left(\frac{R1^2 + R2^2 + R3^2 + \dots}{m. s} \right) - CF$$

$$ASS = \left(\frac{T0^2 + T1^2 + T3^2 + \dots}{r. s} \right) - CF$$

$$A \times R \text{ table SS} = \left(\frac{A01^2 + A02^2 + A03^2 + \dots}{b} \right) - CF$$

Error (a) SS= A x R TSS-RASS-ASS.

2) Horizontal Strip Analysis

Form B x R Table and calculate RSS, BSS and Error(b) SS

$$3) \text{ BSS} = \left(\frac{T_0^2 + T_1^2 + T_3^2 + \dots}{r \cdot s} \right) - CF$$

$$4) \text{ B x R table SS} = \left(\frac{B_01^2 + B_02^2 + B_03^2 + \dots}{a} \right) - CF$$

$$5) \text{ Error (b) SS} = \text{B x R TSS} - \text{RSS} - \text{BSS}$$

3) Interaction Analysis

Form A x B Table and calculate BSS, Ax B SSS and Error (b) SS

Treatment	Replication				Total
	B0	B1	B2	...	
A0	T00	T01	T02	...	T0
A1	T10	T11	T12	...	T1
A2	T20	T21	T22	...	T2
.
.
.
Total	C0	C1	C2	...	GT

$$\text{A x B table SS} = \left(\frac{T_0o^2 + T_01^2 + T_03^2 + \dots}{r} \right) - CF$$

$$\text{ABSS} = \text{A x B Table SS} - \text{ASS} - \text{ABSS}$$

$$\text{Error (c) SS} = \text{TSS} - \text{ASS} - \text{BSS} - \text{ABSS} - \text{Error (a) SS} - \text{Error (a) SS}$$

Then complete the ANOVA table.

Questions

1. To accommodate factors which require different sizes of experimental plots in the same experiment _____ design has been evolved

- a) Strip plot b) CRD c) RBD d) LSD

Ans: Strip plot

2. The number of error terms in a strip plot design is

- a) One b) two c) three d) none of these

Ans: three

4. The plot size for the treatments allotted in vertical strips will not be equal when compared to the treatments allotted in horizontal strips.

Ans: True

5. The degrees of freedom for Error (b) in a strip plot design is $(r-1)(a-1)$.

Ans: False

6. The analysis of a strip plot design is carried out in three stages, viz, horizontal strip analysis, vertical strip analysis and interaction analysis.

Ans: True

7. When will you adopt strip plot design?

8. What is error (c) in a strip plot design?

9. Furnish the skeleton ANOVA table with 3 replications, 3 treatments in horizontal strip and 3 treatments in vertical strip.

10. Furnish the layout of a strip plot design.

Lecture.23

Long term experiments – ANOVA table – guard rows – optimum plot size – determination methods.

Long Term Experiments

A long term experiment is an experimental procedure that runs through a long period of time, in order to test a hypothesis or observe a phenomenon that takes place at an extremely slow rate. Several agricultural field experiments have run for more than 100 years. Experiments that are conducted at several sites or repeated over different seasons can also be classified as long term experiments. Performance of crops varies considerably from location to location as well as season to season. This is because of the influence of environmental factors such as rainfall, temperature etc. In order to determine the effects, the experiments have to be repeated at different locations and seasons. With such repetition of experiments practical recommendations may be made with greater confidence especially with new crop varieties or new techniques are introduced. Here we discuss the experiments that are conducted over different locations or different seasons.

Layout of experiment

Once the locations or seasons are decided upon the next step is to select the appropriate design of experiment. The individual experiments may be designed as CRD, RBD, split plot etc. The same design is adopted for all the locations or seasons. However randomization of treatments should be done afresh for each experiment.

Analysis

The results of repeated experiments are analysed using combined analysis of variance method.

The combined analysis is aimed at

1. to test whether there are significant differences between the treatments at various environments or loc or seasons etc.

- test the consistency of the treatment at different environments. i.e. to test the presence or absence of interaction of the treatment with environments.

The presence of interaction will indicate that the responses change with environment.

In the first stage of the combined analysis the results of the individual locations are analysed based on the basic experimental design tried. In the second stage of the analysis various SS are computed by combining all the data.

If the basic design adopted is RBD with t treatments and r replications and p locations the ANOVA table will be

Sources of Variation	Degrees of Freedom	Sum of Squares	Mean Squares	F-ratio
Replication within locations	$p(r-1)$	RSS	RMS	
Locations	$p-1$	LSS	LMS	
Treatments	$t-1$	TrSS	TrMS	$\frac{\text{TrMS}}{\text{LXTMS}}$
Location x Treatments	$(p-1)(t-1)$	LXTSS	LXTMS	$\frac{\text{LXTMS}}{\text{EMS}}$
Combined error	$p(r-1)(t-1)$	ESS	EMS	
Total	$rtp-1$	TSS		

But before proceeding with the combined analysis it is necessary to test whether the EMS of the individual experiments are homogenous and the heterogeneity of EMS can be tested by either Bartlett's test or Hartley's test.

When the EMS are homogenous the analysis is done as follows:

Rep within location SS = Sum of replication SS of all locations

Pooled error SS = sum of error SS of all locations

The treatment X location two-way table is formed. From this two way table treatment SS, locations SS and treatment X location SS are computed.

The significance of treatment X location interaction is tested and if it is found to be significant then the interaction mean square is used for calculating the F value for treatments.

Optimum plot size

Size and shape of experimental units will affect the accuracy of the experimental units. Select a plot with optimum plot size for this purpose. **Minimum size of experimental plot for a given degree of precision is known as optimum plot size.** Optimum plot size depends on crop, available land area, number of treatments etc.

To determine the optimum plot size two methods are available. They are (1) Maximum curvature method and (2) Fairfield Smith's variance law. For determining the optimum plot size in either method data are to be collected by conducting an Uniformity trial.

An uniformity trial is a trial conducted over an experimental material by selecting a particular variety of a crop and for the entire experimental unit uniform treatments are given. At harvest, the experimental unit is divided into small basic units (depending on the crop) and yield recorded. Then to find the optimum plot size, the basic units are combined by adding the basic units in rows or columns. But while combining rows or columns no row or column should be left out. Then for the new units formed we calculate coefficient of variation and based on the CV values the optimum plot size is determined.

Questions

1. A long term experiment is an experiment conducted in
a) one season b) more than one season c) more than one year d) **both b and c**

Ans: both b and c

2. The homogeneity of the error variances of the individual seasons or locations is tested by
a) t test b) F test c) Bartlett's test d) none of these

Ans: Bartlett's test

3. The significant interaction indicated that the responses change with the environment.

Ans: True

4. Minimum size of experimental plot for a given degree of precision is known as optimum plot size.

Ans: True

5. The designs adopted in two seasons for the same experiment need not be the same design.

Ans: False

6. The combined analysis is used to test whether there are significant differences between the treatments at various environments or loc or seasons etc.

Ans: True

7. What is a uniformity trail?

8. Mention the methods of determining the optimum plot size.

9. How to determine the optimum plot size?

10. Furnish the ANOVA table of an experiment conducted in RBD in s seasons.

Exercise.1

Diagrammatic and graphic representation – simple, multiple, component and percentage bar diagram – pie chart – histogram. Frequency polygon, frequency curve

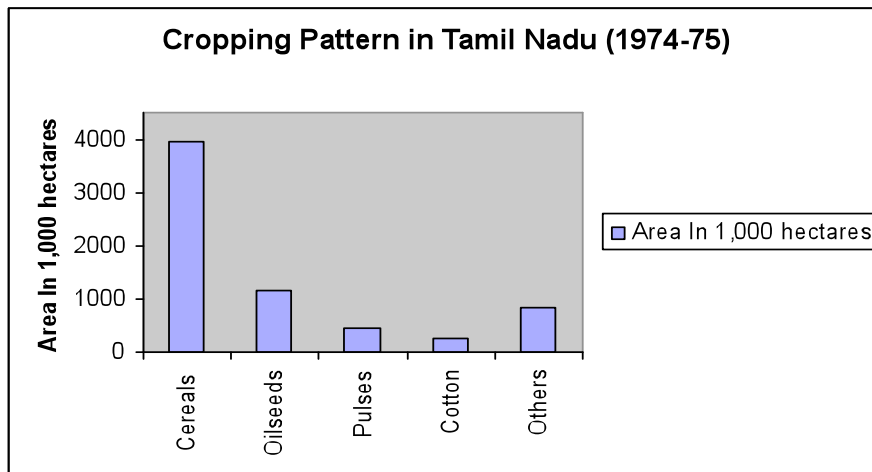
Simple Bar diagram

Example

The cropping pattern in Tamil Nadu in the year 1974-75 was as follows.

Crops	Area In 1,000 hectares
Cereals	3940
Oilseeds	1165
Pulses	464
Cotton	249
Others	822

The simple bar diagram for this data is given below.

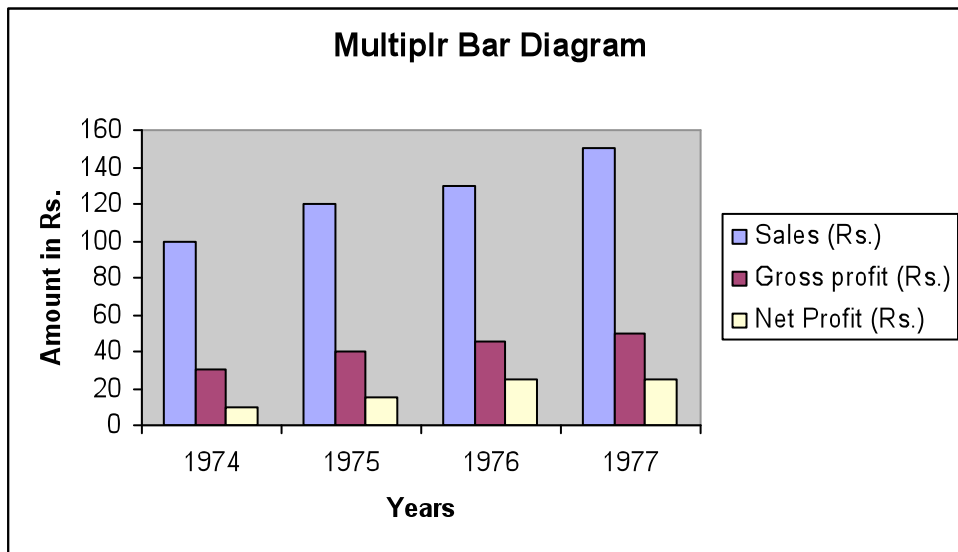


Multiple bar diagram

Example 1

Draw a multiple bar diagram for the following data

Year	Sales (Rs.)	Gross Profit (Rs.)	Net Profit (Rs.)
1974	100	30	10
1975	120	40	15
1976	130	45	25
1977	150	50	25
Total	500	165	75

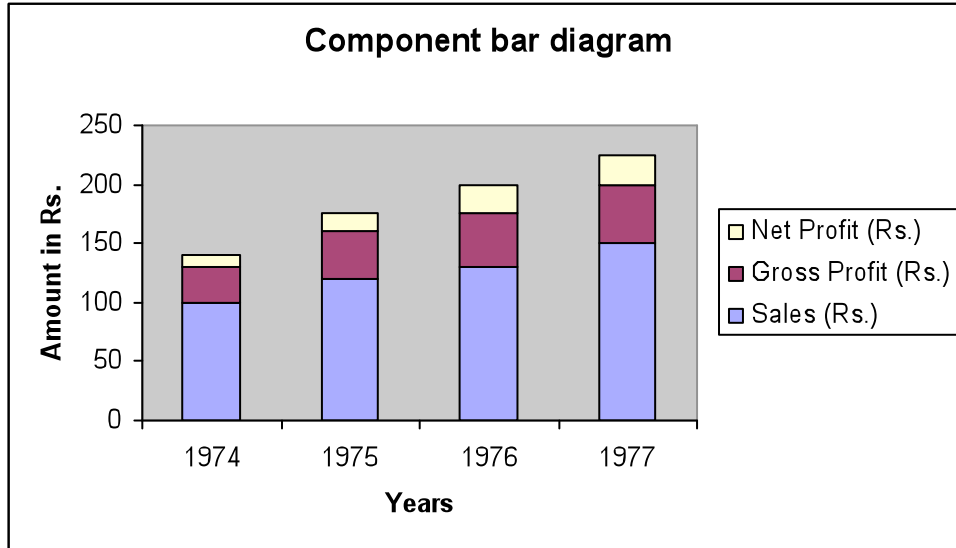


Component bar diagram

Example 2

Draw a component bar diagram for the following data

Year	Sales (Rs.)	Gross Profit (Rs.)	Net Profit (Rs.)
1974	100	30	10
1975	120	40	15
1976	130	45	25
1977	150	50	25



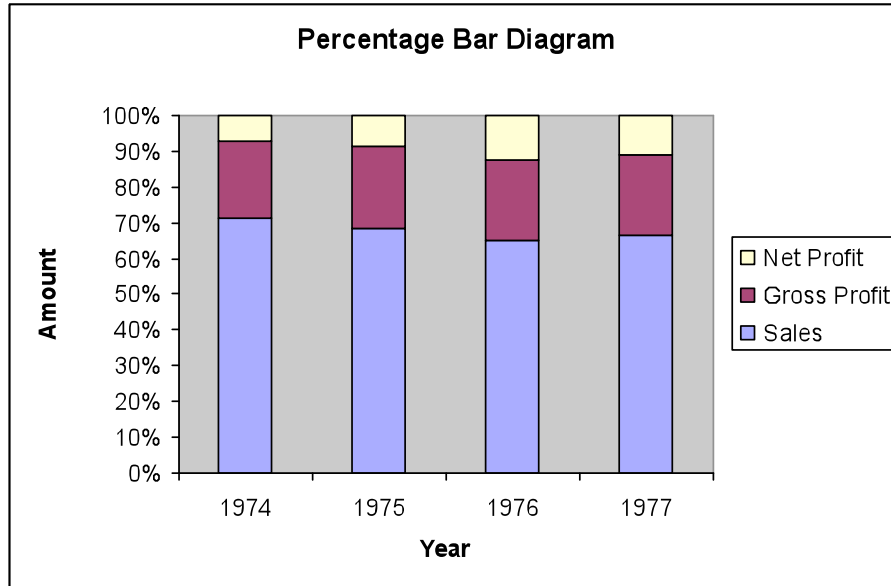
Percentage bar diagram

Example 3

Draw a Percentage bar diagram for the following data

Using the formula $\text{Percentage} = \frac{\text{Actual value}}{\text{Total of the actual value}} \times 100$, the above table is converted.

Year	Sales (Rs.)	Gross Profit (Rs.)	Net Profit (Rs.)
1974	71.43	21.43	7.14
1975	68.57	22.86	8.57
1976	65	22.5	12.5
1977	66.67	22.22	11.11



Pie chart / Pie Diagram

Example 4

Given the population of 1991 of four southern states of India. Construct a pie diagram for the following data.

State	Population
Andhra Pradesh	663
Karnataka	448
Kerala	290
Tamil Nadu	556
Total	1957

Using the formula

$$\text{Angle} = \frac{\text{Actual value}}{\text{Total of the actual value}} \times 360^\circ$$

(or)

$$\text{Angle} = \frac{\text{Percentage}}{100} \times 360^\circ$$

The table value becomes

State	Population
Andhra Pradesh	121.96
Karnataka	82.41
Kerala	53.35
Tamil Nadu	102.28

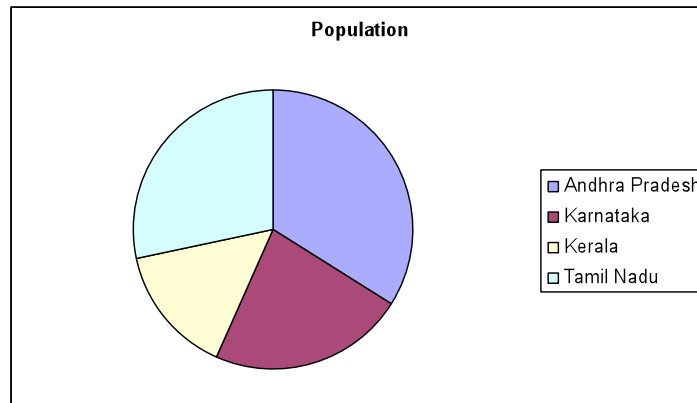
$$\text{Radius} = \pi r^2$$

Here $\pi r^2 = 1957$

$$r^2 = \frac{1957}{\pi} = \frac{1957}{3.14} = 623.24$$

$$r = 24.96$$

$$r = 25 \text{ (approx)}$$



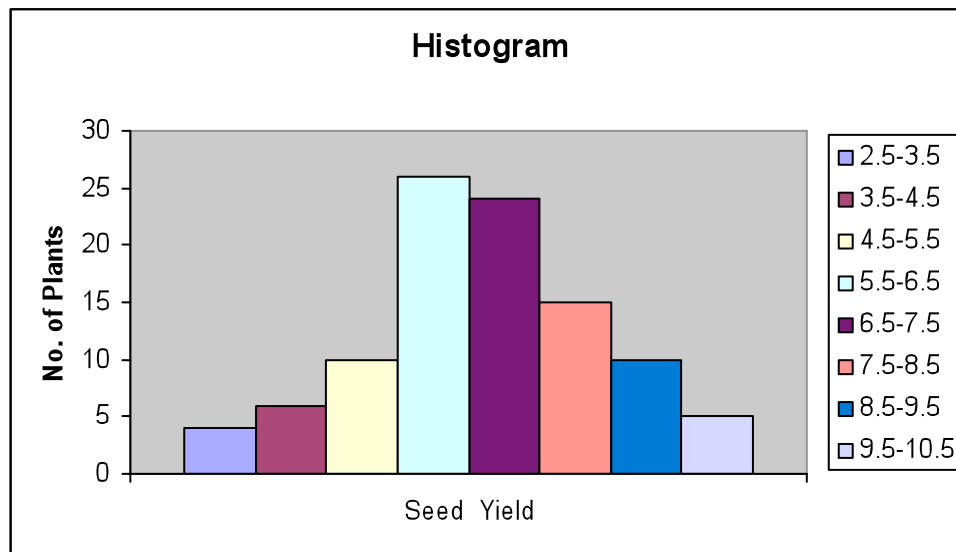
Histogram

Example 5

Draw a histogram for the following data

Seed Yield	No. of Plants
2.5-3.5	4
3.5-4.5	6
4.5-5.5	10

5.5-6.5	26
6.5-7.5	24
7.5-8.5	15
8.5-9.5	10
9.5-10.5	5



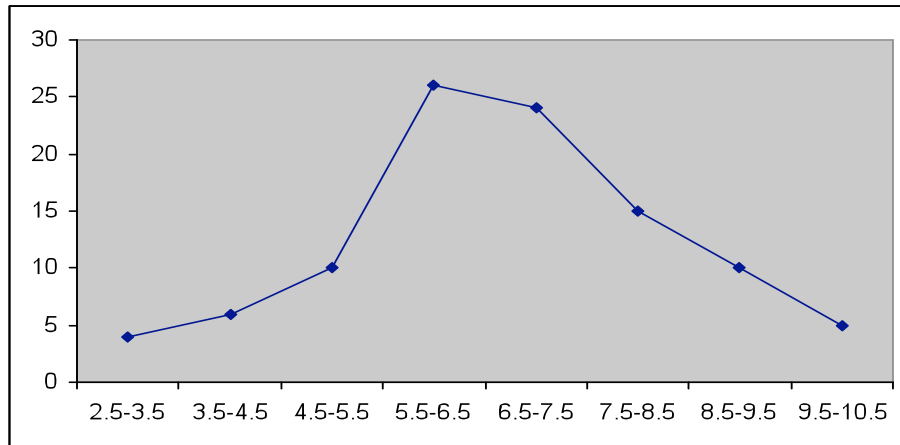
Frequency Polygon

Example 6

Draw frequency polygon for the following data

Seed Yield	No. of Plants
2.5-3.5	4
3.5-4.5	6
4.5-5.5	10
5.5-6.5	26
6.5-7.5	24

7.5-8.5	15
8.5-9.5	10
9.5-10.5	5

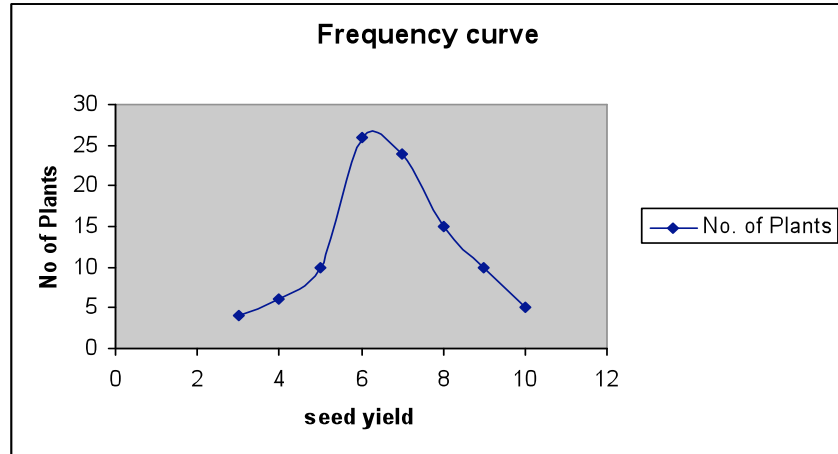


Frequency curve

Example 7

Draw frequency curve for the following data

Seed Yield	No. of Plants
2.5-3.5	4
3.5-4.5	6
4.5-5.5	10
5.5-6.5	26
6.5-7.5	24
7.5-8.5	15
8.5-9.5	10
9.5-10.5	5

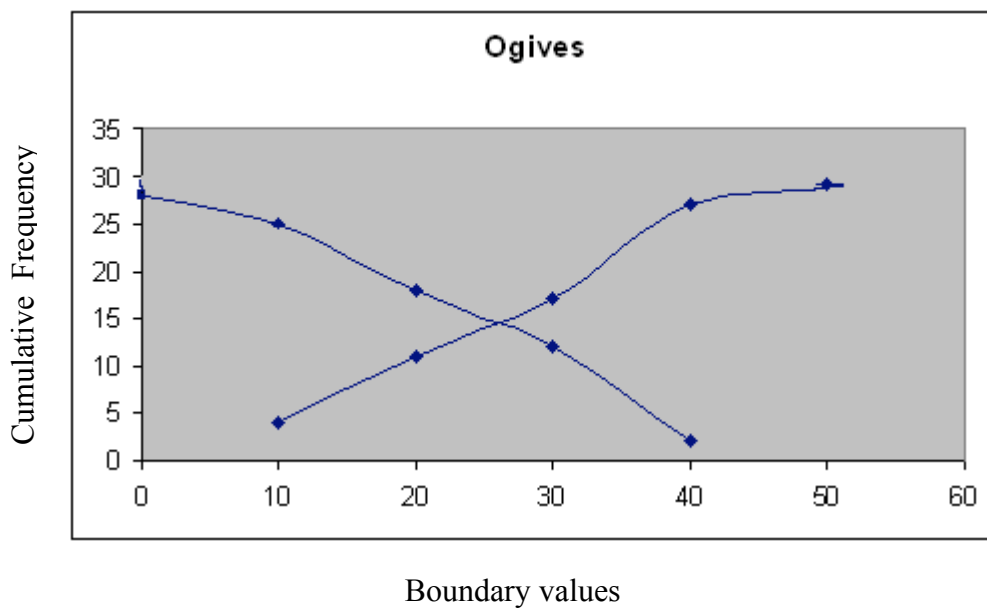


Ogives

Example 8

Draw ogives for the following data

Continuous Interval	Mid Point	Frequency	< cumulative Frequency	> cumulative frequency
0-10	5	4	4	29
10-20	15	7	11	25
20-30	25	6	17	18
30-40	35	10	27	12
40-50	45	2	29	2



Learning Exercise

1) The mean yields of green gram (Kg/hectare) under different weedicide treatment were as follows

Weedicide	Yield(Kg/ha)
Oxadiazon	1382
Fluchloralin	1117
Isopofuron	1066
Unweeded Control	767

Draw Simple Bar Diagram

2) The cropping pattern of Tamil Nadu in 3 different years was as follows.

Crops	Area		
	2002	2003	2004
Cereals	3600	3650	3950
Oilseeds	1000	1150	1100
Pulses	400	450	460
Cotton	200	230	240
Others	800	820	820

Draw multiple bar diagram, Component bar diagram, Percentage bar diagram and Pie chart.

3) The yields of a crop sorghum from 100 experimental plots is given below. Construct histogram, frequency polygon, frequency curve and ogives.

Grain Yield	No. of Plants
--------------------	----------------------

Statistics

65-85	3
85-105	5
105-125	7
125-145	20
145-165	24
165-185	26
185-205	12
205-225	02
225-245	01

Exercise.2

Measures of central tendency – mean median, mode, geometric mean, harmonic mean for raw data

Arithmetic mean or mean

Arithmetic mean or simply the mean of a variable is defined as the sum of the observations divided by the number of observations. If the variable x assumes n values $x_1, x_2 \dots x_n$ then the mean, \bar{x} , is given by

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

This formula is for the ungrouped or raw data.

Example 1

Calculate the mean for 2, 4, 6, 8, 10

Solution

$$\bar{x} = \frac{2 + 4 + 6 + 8 + 10}{5} = \frac{30}{5} = 6$$

Short-Cut method

Under this method an assumed or an arbitrary average (indicated by A) is used as the basis of calculation of deviations from individual values. The formula is

$$\bar{x} = A + \frac{\sum d}{n}$$

Where, A = the assumed mean or any value in x

d = the deviation of each value from the assumed mean

Example 2

A student's marks in 5 subjects are 75, 68, 80, 92, 56. Find his average mark.

Solution

X	d=x-A
75	7
68	0
80	12
92	24
56	-12
Total	31

$$\begin{aligned} \bar{x} &= A + \frac{\sum d}{n} \\ &= 68 + \frac{31}{5} \\ &= 68 + 6.2 \\ &= 74.2 \end{aligned}$$

Median

The median is the middle most item that divides the group into two equal parts, one part comprising all values greater, and the other, all values less than median.

Ungrouped or Raw data

Arrange the given values in the ascending or decreasing order. If the number of values is odd, median is the middle value

If the number of values are even, median is the mean of middle two values.

By formula

When n is odd Median = Md = $\left(\frac{n+1}{2}\right)^{th}$ value

When n is even Average of $\left(\frac{n}{2}\right)$ and $\left(\frac{n}{2} + 1\right)^{th}$ value

Example 3

If the weights of sorghum ear heads are 45, 60, 48, 100, 65 gms. Calculate the median

Solution

Here $n = 5$

First arrange it in ascending order

45, 48, 60, 65, 100

$$\begin{aligned} \text{Median} &= \left(\frac{n+1}{2} \right)^{\text{th}} \text{ value} \\ &= \left(\frac{5+1}{2} \right) = 3^{\text{rd}} \text{ value} = 60 \end{aligned}$$

Example 4

If the sorghum ear- heads are 5, 48, 60, 65, 65, 100 gms. Calculate the median.

Solution

Here $n = 6$

$$\text{Median} = \text{Average of } \left(\frac{n}{2} \right) \text{ and } \left(\frac{n}{2} + 1 \right)^{\text{th}} \text{ value}$$

$$\left(\frac{n}{2} \right) = \frac{6}{2} = 3^{\text{rd}} \text{ value} = 60 \text{ and } \left(\frac{n}{2} + 1 \right) = \frac{6}{2} + 1 = 4^{\text{th}} \text{ value} = 65$$

$$\text{Median} = \frac{60 + 65}{2} = 62.5 \text{ g}$$

Mode

The mode refers to that value in a distribution, which occur most frequently.

Computation of the mode

Ungrouped or Raw Data

For ungrouped data or a series of individual observations, mode is often found by mere inspection.

Example 5

Find the mode for the following seed weight

2, 7, 10, 15, 10, 17, 8, 10, 2 gms

Mode = 10

In some cases the mode may be absent while in some cases there may be more than one mode.

Example 6

1. 12, 10, 15, 24, 30 (no mode)

2. 7, 10, 15, 12, 7, 14, 24, 10, 7, 20, 10

the modes are 7 and 10

Geometric mean

The geometric mean of a series containing n observations is the nth root of the product of the values.

If x_1, x_2, \dots, x_n are observations then

$$G.M = \sqrt[n]{x_1, x_2, \dots, x_n}$$

$$= (x_1, x_2, \dots, x_n)^{1/n}$$

$$\text{Log GM} = \frac{1}{n} \log(x_1, x_2, \dots, x_n)$$

$$= \frac{1}{n} (\log x_1 + \log x_2 + \dots + \log x_n)$$

$$= \frac{\sum \log x_i}{n}$$

$$GM = \text{Antilog } \frac{\sum \log x_i}{n}$$

For grouped data

$$GM = \text{Antilog } \left[\frac{\sum f \log x_i}{N} \right]$$

Example 7

If the weights of sorghum ear heads are 45, 60, 48, 100, 65 gms. Find the Geometric mean for the following data

Weight of ear head x (g)	Log x
45	1.653
60	1.778
48	1.681
100	2
65	1.813
Total	8.925

Solution

Here n = 5

$$GM = \text{Antilog } \frac{\sum \log x_i}{n}$$

$$= \text{Antilog } \frac{8.925}{5}$$

$$= \text{Antilog } 1.785$$

$$= 60.95$$

Harmonic mean (H.M)

Harmonic mean of a set of observations is defined as the reciprocal of the arithmetic average of the reciprocal of the given values. If x_1, x_2, \dots, x_n are n observations,

$$H.M = \frac{n}{\sum_{i=1}^n \left(\frac{1}{x_i} \right)}$$

Example 8

From the given data calculate H.M 5, 10, 17, 24, 30

X	$\frac{1}{x}$
5	0.2000
10	0.1000
17	0.0588
24	0.0417
30	0.4338

$$H.M = \frac{5}{0.4338} = 11.526$$

Learning Exercise

The weight of 15 earheads of sorghum are 100, 102, 118, 124, 126, 98, 100, 100, 118, 95, 113, 115, 123, 121, 117. Find

- (i) Average of weight
- (ii) Median
- (iii) Mode
- (iv) Harmonic mean**
- (v) Geometric mean

Practical.3

Measures of central tendency – mean, median, mode, geometric mean and harmonic mean for grouped data

Arithmetic mean or mean

Grouped Data

The mean for grouped data is obtained from the following formula:

$$\bar{x} = \frac{\sum fx}{N}$$

Where x = the mid-point of individual class

f = the frequency of individual class

N = the sum of the frequencies or total frequencies.

Short-cut method

$$\bar{x} = A + \frac{\sum fd}{N} \times c$$

Where $d = \frac{x - A}{c}$

A = any value in x

N = total frequency

c = width of the class interval

Example 1

Given the following frequency distribution, calculate the arithmetic mean

Marks	:	64	63	62	61	60	59
Number of	}						
Students		:	8	18	12	9	7

Solution

X	f	fx	d=x-A	fd
64	8	512	2	16
63	18	1134	1	18
62	12	744	0	0
61	9	549	-1	-9
60	7	420	-2	-14
59	6	354	-3	-18
	60	3713		-7

Direct method

$$\bar{x} = \frac{\sum fx}{N}$$

$$\bar{x} = \frac{3713}{60} = 61.88$$

Short-cut method

$$\bar{x} = A + \frac{\sum fd}{N} \times c$$

Here A = 62

$$\bar{x} = 62 - \frac{7}{60} \times 1 = 61.88$$

Example 2

For the frequency distribution of seed yield of sesamum given in table calculate the mean yield per plot.

Yield per plot in(in g)	64.5-84.5	84.5-104.5	104.5-124.5	124.5-144.5
No of plots	3	5	7	20

Solution

Yield (in g)	No of Plots (f)	Mid X	$d = \frac{x - A}{c}$	fd
64.5-84.5	3	74.5	-1	-3
84.5-104.5	5	94.5	0	0
104.5-124.5	7	114.5	1	7
124.5-144.5	20	134.5	2	40
Total	35			44

$$A=94.5$$

The mean yield per plot is

$$\bar{x} = A + \frac{\sum fd}{N} \times c$$

$$\bar{x} = 94.5 + \frac{44}{35} \times 20 = 119.64 \text{ g}$$

Median

Grouped data

In a grouped distribution, values are associated with frequencies. Grouping can be in the form of a discrete frequency distribution or a continuous frequency distribution. Whatever may be the type of distribution, cumulative frequencies have to be calculated to know the total number of items.

Cumulative frequency: (cf)

Cumulative frequency of each class is the sum of the frequency of the class and the frequencies of the pervious classes, ie adding the frequencies successively, so that the last cumulative frequency gives the total number of items.

Discrete Series

Step1: Find cumulative frequencies.

Step2: Find $\left(\frac{N}{2} + 1\right)$

Step3: See in the cumulative frequencies the value just greater than $\left(\frac{N}{2} + 1\right)$

Step4: Then the corresponding value of x is median.

Example 3

The following data pertains to the number of members in a family. Find median size of the family.

Number of members x	1	2	3	4	5	6	7	8	9	10	11	12
Frequency f	1	3	5	6	10	13	9	5	3	2	2	1

Solution

X	f	cf
1	1	1
2	3	4
3	5	9
4	6	15
5	10	25
6	13	38
7	9	47
8	5	52
9	3	55
10	2	57
11	2	59
12	1	60
	60	

$$\begin{aligned} \text{Median} &= \text{size of } \left(\frac{N+1}{2} \right)^{\text{th}} \text{ item} \\ &= \text{size of } \left(\frac{60+1}{2} \right)^{\text{th}} \text{ item} \\ &= 30.5^{\text{th}} \text{ item} \end{aligned}$$

The cumulative frequency just greater than 30.5 is 38. and the value of x corresponding to 38 is 6. Hence the median size is 6 members per family.

Continuous Series

The steps given below are followed for the calculation of median in continuous series.

Step1: Find cumulative frequencies.

Step2: Find $\left(\frac{N}{2} \right)$

Step3: See in the cumulative frequency the value first greater than $\left(\frac{N}{2} \right)$, Then the corresponding class interval is called the Median class. Then apply the formula

$$\text{Median} = l + \frac{\frac{N}{2} - m}{f} \times c$$

Where l = Lower limit of the median class
 m = cumulative frequency preceding the median
 c = width of the median class
 f = frequency in the median class.
 N = Total frequency.

Example 4

For the frequency distribution of weights of sorghum ear-heads given in table below. Calculate the median.

Weights of ear heads (in g)	No of ear heads (f)	Cumulative frequency (m)
60-80	22	22
80-100	38	60
100-120	45	105
120-140	35	140
140-160	20	160
Total	160	

Solution

$$\text{Median} = l + \frac{\frac{N}{2} - m}{f} \times c$$

$$\left(\frac{N}{2}\right) = \left(\frac{160}{2}\right) = 80$$

Here $l = 100$, $N=160$, $f = 45$, $c = 20$, $m = 60$

$$\text{Median} = 100 + \frac{80 - 60}{45} \times 20 = 108.8 \text{ gms}$$

Geometric mean

Grouped Data

For grouped data

$$\text{GM} = \text{Antilog} \left[\frac{\sum f \log x_i}{N} \right]$$

Example 5

Find the Geometric mean for the following

Weight of sorghum (x)	No. of ear head(f)
50	4
65	6
75	16
80	8
95	7
100	4

Solution

Weight of sorghum (x)	No. of ear head(f)	Log x	flog x
50	5	1.699	8.495
63	10	10.799	17.99
65	5	1.813	9.065
130	15	2.114	31.71
135	15	2.130	31.95
Total	50	9.555	99.21

Here N= 50

$$\begin{aligned}
 \text{GM} &= \text{Antilog} \left[\frac{\sum f \log x_i}{N} \right] \\
 &= \text{Antilog} \left[\frac{99.21}{50} \right] \\
 &= \text{Antilog } 1.9842 = 96.43
 \end{aligned}$$

Continuous distribution

Example 6

For the frequency distribution of weights of sorghum ear-heads given in table below. Calculate the Geometric mean

Weights of ear heads (in g)	No of ear heads (f)
60-80	22
80-100	38
100-120	45
120-140	35
140-160	20
Total	160

Solution

Weights of ear heads (in g)	No of ear heads (f)	Mid x	Log x	f log x
60-80	22	70	1.845	40 59
80-100	38	90	1.954	74.25
100-120	45	110	2.041	91.85
120-140	35	130	2.114	73.99
140-160	20	150	2.176	43.52
Total	160			324.2

Here N = 160

$$GM = \text{Antilog} \left[\frac{\sum f \log x_i}{N} \right]$$

$$= \text{Antilog} \left[\frac{324.2}{160} \right]$$

$$= \text{Antilog} [2.02625]$$

$$= 106.23$$

Harmonic mean

For a frequency distribution

$$H.M = \frac{N}{\sum_{i=1}^n f\left(\frac{1}{x_i}\right)}$$

Example 7

The marks secured by some students of a class are given below. Calculate the harmonic mean.

Marks	20	21	22	23	24	25
Number of Students	4	2	7	1	3	1

Solution

Marks X	No of Students f	$\frac{1}{x}$	$f\left(\frac{1}{x}\right)$
20	4	0.0500	0.2000
21	2	0.0476	0.0952
22	7	0.0454	0.3178
23	1	0.0435	0.0435
24	3	0.0417	0.1251
25	1	0.0400	0.0400
	18		0.8216

$$H.M = \frac{N}{\sum_{i=1}^n f\left(\frac{1}{x_i}\right)} = \frac{18}{0.1968} = 21.91$$

Learning Exercise

For the following frequency distribution find the

- (i) Mean
- (ii) Median
- (iii) Mode
- (iv) Harmonic mean
- (iv) Geometric mean

Weight of earheads in gms	No. of earhead
40 - 60	6
60 - 80	8
80 – 100	35
100 -120	55
120 -140	30
140 – 160	15
160 – 180	12
180 – 200	9

Practical.4

Measures of dispersion – variance, standard deviation and coefficient of variation for raw data

Variance

The square of the standard deviation is called variance

(i.e) variance = \sqrt{SD}

Standard Deviation

It is defined as the positive square-root of the arithmetic mean of the Square of the deviations of the given observation from their arithmetic mean.

The standard deviation is denoted by S in case of sample and Greek letter s (sigma) in case of population.

The formula for calculating standard deviation is as follows

$$S = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n - 1}} \text{ for raw data}$$

And for grouped data the formulas are

$$S = \sqrt{\frac{\sum fx^2}{N} - \left(\frac{\sum fx}{N}\right)^2} \text{ for discrete data}$$

$$S = C \times \sqrt{\frac{\sum fd^2}{N} - \left(\frac{\sum fd}{N}\right)^2} \text{ for continuous data}$$

Where $d = \frac{x - A}{C}$

C = class interval

Example 1

Raw Data

The weights of 5 ear-heads of sorghum are 100, 102, 118, 124, 126 gms. Find the standard deviation.

Solution

x	x ²
100	10000
102	10404
118	13924
124	15376
126	15876
Total	570 65580

$$\text{Standard deviation } S = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}}$$

$$= \sqrt{\frac{65580 - \frac{(570)^2}{5}}{5-1}} = \sqrt{150} = 12.25 \text{ gms}$$

$$\text{Variance} = \sqrt{12.25} = 3.5$$

Coefficient of Variation

The Standard deviation is an absolute measure of dispersion. It is expressed in terms of units in which the original figures are collected and stated. The standard deviation of heights of students cannot be compared with the standard deviation of weights of students, as both are expressed in different units, i.e heights in centimeter and weights in kilograms. Therefore the standard deviation must be converted into a relative measure of dispersion for the purpose of comparison. The relative measure is known as

the coefficient of variation. The coefficient of variation is obtained by dividing the standard deviation by the mean and expressed in percentage. Symbolically, Coefficient of

$$\text{variation (C.V)} = \frac{SD}{mean} \times 100$$

If we want to compare the variability of two or more series, we can use C.V. The series or groups of data for which the C.V. is greater indicate that the group is more variable, less stable, less uniform, less consistent or less homogeneous. If the C.V. is less, it indicates that the group is less variable or more stable or more uniform or more consistent or more homogeneous.

Example 2

Consider the measurement on yield and plant height of a paddy variety. The mean and standard deviation for yield are 50 kg and 10 kg respectively. The mean and standard deviation for plant height are 55 am and 5 cm respectively.

Here the measurements for yield and plant height are in different units. Hence the variabilities can be compared only by using coefficient of variation.

For yield, $CV = \frac{10}{50} \times 100 = 20\%$

For plant height, $CV = \frac{5}{55} \times 100 = 9.1\%$

The yield is subject to more variation than the plant height.

Learning Exercise

1.	The weights of 8 earheads of sorghum are 14, 29, 9, 15, 20, 17, 12, and 11. Find Standard Deviation and Variance and coefficient of variation.											
2.	Find out which of the following batsmen is more consistent in scoring.											
	Batsman A	5	7	16	27	39	53	56	61	80	101	105
	Batsman B	0	4	16	21	41	43	57	78	83	93	95

Exercise.5
Measures of dispersion – variance, standard deviation and coefficient of variation
for grouped data

Standard deviation and Variance

Example 1

Discrete distribution

The frequency distributions of seed yield of 50 sesamum plants are given below. Find the standard deviation.

Seed yield in gms (x)	3	4	5	6	7
Frequency (f)	4	6	15	165	10

Solution

Seed yield in gms (x)	f	fx	fx ²
3	4	12	36
4	6	24	96
5	15	75	375
6	15	90	540
7	10	70	490
Total	50	271	1537

Here N = 50

$$\begin{aligned} \text{Standard deviation } S &= \sqrt{\frac{\sum fx^2}{N} - \left(\frac{\sum fx}{N}\right)^2} \\ &= \sqrt{\frac{1537}{50} - \left(\frac{271}{50}\right)^2} \\ &= \sqrt{30.74 - 29.3764} \end{aligned}$$

$$= 1.1677 \text{ gms}$$

$$\text{Variance} = \sqrt{1.1677} = 1.081$$

Example 2

Continuous distribution

The Frequency distributions of seed yield of 50 sesamum plants are given below. Find the standard deviation.

Seed yield in gms (x)	2.5-3.5	3.5-4.5	4.5-5.5	5.5-6.5	6.5-7.5
No. of plants (f)	4	6	15	165	10

Solution

Seed yield in gms (x)	No. of Plants f	Mid x	$d = \frac{x - A}{C}$	df	$d^2 f$
2.5-3.5	4	3	-2	-8	16
3.5-4.5	6	4	-1	-6	6
4.5-5.5	15	5	0	0	0
5.5-6.5	15	6	1	15	15
6.5-7.5	10	7	2	20	40
Total	50	25	0	21	77

A=Assumed mean = 5

N=50, C=1

$$S = C \times \sqrt{\frac{\sum fd^2}{N} - \left(\frac{\sum fd}{N}\right)^2}$$

$$= 1 \times \sqrt{\frac{77}{50} - \left(\frac{21}{50}\right)^2}$$

$$= \sqrt{1.54 - 0.1764}$$

$$= \sqrt{1.3636} = 1.1677$$

Variance = $\sqrt{1.1677} = 1.081$

Coefficient of variation

Example 3

Consider the measurement on yield and plant height of a paddy variety. The mean and standard deviation for yield are 50 kg and 10 kg respectively. The mean and standard deviation for plant height are 55 am and 5 cm respectively.

Here the measurements for yield and plant height are in different units. Hence the variabilities can be compared only by using coefficient of variation.

For yield, $CV = \frac{10}{50} \times 100 = 20\%$

For plant height, $CV = \frac{5}{55} \times 100 = 9.1\%$

The yield is subject to more variation than the plant height.

Learning Exercise

1.	From the data given below, find which series is more consistent						
	Variable	10-20	20-30	30-40	40-50	50-60	60-70
	Series A	10	16	30	40	26	18
	Series B	22	18	32	34	18	16
2.	The yield of a crop sorghum from 31 experimental plots are given below. Find the Range, Standard deviation, Variance, Coefficient of variation.						
		Grain Yield			No. of Plots		

	130	3
	135	4
	140	6
	145	6
	146	3
	148	5
	149	2
	150	1
	157	1

3. The following table gives the protein intake of 400 families. Find the Range, Standard deviation, Variance, Coefficient of variation.

Protein intake / Consumption unit Per day in grams	No. of Families
15 - 25	30
25 – 35	40
35 – 45	100
45 – 55	110
55 – 65	80
65 –75	30
75 - 85	10

Exercise.6

Selection of simple random sampling using lottery method and random numbers

Simple Random sampling (SRS)

The basic probability sampling method is the simple random sampling. It is the simplest of all the probability sampling methods. It is used when the population is homogeneous.

When the units of the sample are drawn independently with equal probabilities. The sampling method is known as Simple Random Sampling (SRS). Thus if the population consists of N units, the probability of selecting any unit is $1/N$.

A theoretical definition of SRS is as follows

Suppose drawn sample of size n from a population of size N . There are ${}^N C_n$ possible sample of size n . If all possible samples have an equal probability $1/{}^N C_n$ of being drawn, the sampling is said to be simple random sampling.

There are two methods in SRS

1. Lottery method
2. Random no. table method

Lottery method

This is most popular method and simplest method. In this method all the items of the universe are numbered on separate slips of paper of same size, shape and color. They are folded and mixed up in a drum or a box or a container. A blindfold selection is made. Required numbers of slips are selected for the desired sample size. The selection of items thus depends on chance.

For example if we select 5 students out of 50 students on slips of the same size and mix them, then we make a blindfold selection of 5 students. This method is also called unrestricted random sampling because units are selected from the population without any restriction. This method is mostly used in lottery draws. If the universe is infinite, this method is inapplicable. There is a lot of possibility of personal prejudice if the size and shape of the slips are not identical.

Random number table method

As the lottery method cannot be used when the population is infinite, the alternative method is using of table of random numbers.

There are several standard tables of random numbers. But the credit for this technique goes to Prof. LHC. Tippett (1927). The random number table consists of 10,400 four-figured numbers. There are various other random numbers. They are fishers and Yates (1938) comprising of 15,000 digits arranged in twos. Kendall and B.B Smith (1939) consisting of 1, 00,000 digits grouped in 25,000 sets of 4 digit random numbers, Rand corporation (1955) consisting of 2, 00,000 random numbers of 5 digits each etc.,

Learning Exercise

The following data refers to the Kapas yield of 96 plants.

82	102	88	93	97	38	103	92
102	62	63	72	64	68	59	69
73	65	46	79	87	84	29	52
28	36	37	53	49	51	30	37
56	66	42	37	35	97	32	35
89	99	54	72	26	67	18	27
60	72	33	42	52	82	14	22
57	73	63	61	63	92	40	58
62	61	43	25	42	36	17	30
75	87	47	56	76	36	35	44
56	51	111	73	93	58	49	89
50	80	54	55	91	12	82	76

Statistics

Select a sample of 25 plants by using simple random sampling method. Also calculate the mean of the 25 samples and verify whether the mean is equal to the mean of the 96 plants.

Exercise.7
Students's t test – paired and independent t test

Test for single Mean (n<30)

1. Form the null hypothesis

$$H_0: \mu = \mu_0$$

(i.e) There is no significance difference between the sample mean and the population mean i.e., $\mu = \mu_0$

2. Form the Alternate hypothesis

$$H_1: \mu \neq \mu_0 \text{ (or } \mu > \mu_0 \text{ or } \mu < \mu_0)$$

i.e., There is significance difference between the sample mean and the population mean

3. Level of Significance

The level may be fixed at either 5% or 1%

4. Test statistic

$$t_{cal} = \left| \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \right| \sim t_{(n-1)} \text{ d.f}$$

Where $\bar{x} = \frac{\sum x_i}{n}$

Where $s = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}}$

5. Find the table value

$$t_{tab} = t_{(0.05, (n-1))} \text{ d.f}$$

1. Inference

If $t_{cal} < t_{tab}$

- (i) We accept the null hypothesis H_0
- (ii) There is no significant difference

(or) if $t_{cal} > t_{tab}$

- (i) We reject the null hypothesis H_0 (ie) we accept the alternative hypothesis
- (ii) There is significant difference between the sample mean and the population mean.

Example 1

Based on field experiments, a new variety green gram is expected to give a yield of 12.0 quintals per hectare. The variety was tested on 10 randomly selected farmers' fields. The yield (quintals/hectare) were recorded as 14.3, 12.6, 13.7, 10.9, 13.7, 12.0, 11.4, 12.0, 12.6, 13.1. Do the results conform to the expectation?

Solution

Null hypothesis $H_0: \mu = 12.0$

(i.e) the average yield of the new variety of green gram is 12.0 quintals/hectare.

Alternative Hypothesis: $H_1: \mu \neq 12.0$

(i.e) the average yield is not 12.0 quintals/hectare

Level of significance: 5 %

Test statistic

$$t_{cal} = \frac{\left| \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \right| \sim t_{(n-1)} \text{ d.f.}}$$

From the given data

$$\sum x = 126.3 \quad \sum x^2 = 1605.77$$

$$\bar{x} = \frac{\sum x}{n} = \frac{126.3}{10} = 12.63$$

$$s = \sqrt{\frac{\sum x^2 - \frac{(\sum x)^2}{n}}{n-1}} = \sqrt{\frac{1605.77 - \frac{1595.169}{9}}{9}} = \sqrt{\frac{10.601}{9}}$$

$$= 1.0853$$

$$\frac{s}{\sqrt{n}} = \frac{1.0853}{\sqrt{10}} = 0.3432$$

$$\text{Now } t_{cal} = \frac{\left| \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \right| \sim t_{(n-1)} \text{ d.f.}}$$

$$= t_{cal} = \frac{12.63 - 12}{0.3432} = 1.836$$

Table value

$$t_{(0.05,9)} = 2.262 \quad (\text{two tailed test})$$

Inference

$$t_{cal} < t_{tab}$$

We accept the null hypothesis H_0

We conclude that the new variety of green gram will give an average yield of 12 quintals/hectare.

Note

F-test is used to test the equality of two means

$$F = \frac{S_1^2}{S_2^2} \sim F_{(n_1 - 1, n_2 - 1)} \text{ d.f. if } S_1^2 > S_2^2$$

where S_1^2 is the variance of the first sample whose size is n_1 .

S_2^2 is the variance of the second sample whose size is n_2 .

Otherwise

$$F = \frac{S_2^2}{S_1^2} \sim F_{(n_2 - 1, n_1 - 1)} \text{ d.f. if } S_2^2 > S_1^2$$

Inference

$$F_{cal} < F_{tab}$$

We accept the null hypothesis H_0 . (i.e) the variances are equal.

Test for equality of two means (Independent Samples)

Given two sets of sample observation $x_{11}, x_{12}, x_{13} \dots x_{1n}$. Similarly $x_{21}, x_{22}, x_{23} \dots x_{2n}$ of sizes n_1 and n_2 from the normal population.

1. Using F-Test , test their variances

(i) Variances are Equal:

$$H_0: \mu_1 = \mu_2$$

$$H_1 \mu_1 \neq \mu_2 \text{ (or } \mu_1 < \mu_2 \text{ or } \mu_1 > \mu_2)$$

Test statistic

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{(n_1+n_2-2)} d.f$$

Where

$$S^2 = \frac{\left[\sum x_1^2 - \frac{(\sum x_1)^2}{n_1} \right] + \left[\sum x_2^2 - \frac{(\sum x_2)^2}{n_2} \right]}{n_1 + n_2 - 2}$$

Variances are equal

(a) When the samples have unequal variances and equal number of observations ($n_1=n_2$), the test statistic is

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{\left(\frac{n_1+n_2}{2}-1\right)} d.f$$

(b) When the samples have unequal variances and unequal number of observations ($n_1 \neq n_2$), the test statistic is

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}}$$

This statistic follows neither t nor normal distribution but it follows Behrens-Fisher test. The Behrens – Fisher test is laborious one. An alternative simple method has been suggested by Cochran & Cox. In this method the critical value of t is altered as t_w (i.e) weighted t test

$$t_w = \frac{t_1 \left(\frac{S_1^2}{n_1} \right) + t_2 \left(\frac{S_2^2}{n_2} \right)}{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

Where $t_1 = t_{(n_1-1)}$ d.f

$t_2 = t_{(n_2-1)}$ d.f

Example 2

A group of 5 patients treated with medicine A is of weight 42,39,38,60 &41 kgs. Second group of 7 patients from the same hospital treated with medicine B is of weight 38, 42, 56, 64, 68, 69, & 62 kgs. Find whether there is any difference between medicines?

Solution

Ho:., $\mu_1 = \mu_2$ (i.e) there is no significant difference between the medicines A and B as regards on increase in weight.

H₁ $\mu_1 \neq \mu_2$ (i.e) there is a significant difference between the medicines A and B

Level of significance = 5%

Before we go to test the means first we have to test their variability using F-test.

F-test

Ho:., $\sigma_1^2 = \sigma_2^2$

H₁:., $\sigma_1^2 \neq \sigma_2^2$

$$S_1^2 = \frac{\sum x_1^2 - \frac{(\sum x_1)^2}{n_1}}{n_1 - 1} = 82.5$$

$$S_2^2 = \frac{\sum x_2^2 - \frac{(\sum x_2)^2}{n_2}}{n_2 - 1} = 154.33$$

$$\therefore F = \frac{S_2^2}{S_1^2} \sim F_{(n_2 - 1, n_1 - 1)} \text{ d.f if } S_2^2 > S_1^2$$

$$F_{cal} = \frac{154.33}{32.5} = 1.8707$$

F_{tab}(6,4) d.f=6.16

$$\Rightarrow F_{cal} < F_{tab}$$

We accept the null hypothesis H₀.(i.e) the variances are equal.

Test statistic

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{(n_1+n_2-2)} d.f$$

Where

$$S^2 = \frac{\left[\sum x_1^2 - \frac{(\sum x_1)^2}{n_1} \right] + \left[\sum x_2^2 - \frac{(\sum x_2)^2}{n_2} \right]}{n_1 + n_2 - 2} = \frac{330 + 926}{10} = 125.6$$

$$t = \frac{|44 - 57|}{\sqrt{125.6 \left(\frac{1}{7} + \frac{1}{75} \right)}} = 1.98$$

Table value

$t_{\text{tab}[(5+7-2)=10]} d.f$ at 5% l.o.s = 2.228

Inference:

$$t_{\text{cal}} < t_{\text{tab}}$$

We accept the null hypothesis H_0

We conclude that the medicines A and also B do not differ significantly.

Example 3

The summary of the results of an yield trial on onion with two methods of propagation is given below. Determine whether the methods differ with regard to onion yield. The onion yield is given in Kg/plot.

Method I	Method II
n1=12	n2=12
$\bar{x}_1 = 25.25$	$\bar{x}_2 = 28.83$
SS1=186.25	SS2=737.6667
$S_1^2 = 16.9318$	$S_2^2 = 67.0606$

Solution

Ho:., $\mu_1 = \mu_2$ (i.e) the two propagation method do not differ with regard to onion yield.

H₁ $\mu_1 \neq \mu_2$ (i.e) the two propagation method differ with regard to onion yield.

Level of significance = 5%

Before we go to test the means first we have to test their variability using F-test.

F-test

Ho:., $\sigma_1^2 = \sigma_2^2$

H1:., $\sigma_1^2 \neq \sigma_2^2$

$$S_1^2 = \frac{\sum x_1^2 - \frac{(\sum x_1)^2}{n1}}{n1 - 1} = 16.9318$$

$$S_2^2 = \frac{\sum x_2^2 - \frac{(\sum x_2)^2}{n2}}{n2 - 1} = 67.0606$$

$$\therefore F = \frac{S_2^2}{S_1^2} \sim F_{(n_2 - 1, n_1 - 1)} \text{ d.f if } S_2^2 > S_1^2$$

$$F_{cal} = \frac{67.0606}{16.9318} = 3.961$$

F_{tab}(11,11) d.f=2.82

$$\Rightarrow F_{cal} > F_{tab}$$

We reject the null hypothesis H₀.(i.e) the variances are unequal.

Here the variances are unequal with equal sample size then the test statistic is

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{S \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \sim t_{\left(\frac{n_1+n_2}{2}-1\right)} \text{ d.f}$$

Where

$$S^2 = \frac{\left[\sum x_1^2 - \frac{(\sum x_1)^2}{n1} \right] + \left[\sum x_2^2 - \frac{(\sum x_2)^2}{n2} \right]}{n_1 + n_2 - 2}$$

$$S^2 = \frac{SS1 + SS2}{n_1 + n_2 - 2} = \frac{186.25 + 737.6667}{12 + 12 - 2} = 41.9962$$

$$t = \frac{25.25 - 28.83}{\sqrt{41.9962\left(\frac{1}{12} + \frac{1}{12}\right)}} = \frac{3.58}{\sqrt{6.9994}} = 1.353$$

$t_{cal}=1.353$

table value

$$t_{\left(\frac{n1+n2-1}{2}\right)} = t_{\left(\frac{12+12-1}{2}\right)} = t_{11} \text{ d.f at } 5\% \text{ l.os} = 2.201$$

Inference:

$$t_{cal} < t_{tab}$$

We accept the null hypothesis H_0

We conclude that the two propagation methods do not differ with regard to onion yield.

Example 4

The following data related the rubber percentage of two types of rubber plants, where the sample have been drawn independently. Test for their mean difference.

Type I	6.21	5.70	6.04	4.47	5.22	4.45	4.84	5.84	5.88	5.82	6.09	5.59
	6.06	5.59	6.74	5.55								

Type II	4.28	7.71	6.48	7.71	7.37	7.20	7.06	6.40	8.93	5.91	5.51	6.36
---------	------	------	------	------	------	------	------	------	------	------	------	------

Solution

$H_0: \mu_1 = \mu_2$ (i.e) there is no significance difference between the two samples.

$H_1: \mu_1 \neq \mu_2$ (i.e) there is a significance difference between the two samples.

Level of significance = 5%

Here

n1=16	n2=12
$\sum x_1 = 90.09$	$\sum x_2 = 80.92$
$\bar{x}_1 = 5.63$	$\bar{x}_2 = 6.7431$
$\sum x_1^2 = 513.085$	$\sum x_2^2 = 561.64$

Before we go to test the means first we have to test their variability using F-test.

F-test

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$H_1: \sigma_1^2 \neq \sigma_2^2$$

$$S_1^2 = \frac{\sum x_1^2 - \frac{(\sum x_1)^2}{n_1}}{n_1 - 1} = 0.388$$

$$S_2^2 = \frac{\sum x_2^2 - \frac{(\sum x_2)^2}{n_2}}{n_2 - 1} = 1.452$$

$$\therefore F = \frac{S_2^2}{S_1^2} \sim F_{(n_2 - 1, n_1 - 1)} \text{ d.f if } S_2^2 > S_1^2$$

$$F_{cal} = \frac{1.452}{0.388} = 3.742$$

$$F_{tab}(11, 15) \text{ d.f} = 2.51$$

$$\Rightarrow F_{cal} > F_{tab}$$

We reject the null hypothesis H_0 (i.e) the variances are unequal.

Here the variances are unequal with unequal sample size then the test statistic is

$$t = \frac{|\bar{x}_1 - \bar{x}_2|}{\sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}} \sim t_w$$

$$\text{Where } t_w = \frac{t_1 \left(\frac{S_1^2}{n_1} \right) + t_2 \left(\frac{S_2^2}{n_2} \right)}{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

$$t_1 = t_{(16-1)} \text{ d.f} = 2.131$$

$$t_2 = t_{(12-1)} \text{ d.f} = 2.201$$

$$t_w = \frac{2.131 \left(\frac{0.388}{16} \right) + 2.201 \left(\frac{1.452}{12} \right)}{\frac{0.388}{16} + \frac{1.425}{12}} = 2.187$$

$$t_{cal} = \frac{(5.63 - 6.7431_2)}{\sqrt{\frac{0.388}{16} + \frac{1.452}{12}}} = 2.912$$

Inference

$$t_{cal} > t_{tab}$$

We reject the null hypothesis H_0

(i.e) there is significant difference between the two rubber plants.

Equality of two means (Dependant samples)

Paired t test

Null hypothesis

$H_0 : \mu_1 = \mu_2$ ie the increments are just by chance

Alternative Hypothesis

$H_1 : \mu_1 \neq \mu_2$ ($\mu_1 > \mu_2$ (or) $\mu_1 < \mu_2$)

test statistic

$$t = \frac{|\bar{d}|}{S/\sqrt{n}} \sim t(n-1) d.f$$

Where $\bar{d} = \frac{\sum di}{n}$, $S = \sqrt{\frac{\sum di^2 - \frac{(\sum di)^2}{n}}{n-1}}$

$d_i = X_i - Y_i$ ($i = 1, 2, \dots, n$)

Example 5

In certain food experiment to compare two types of baby foods A and B, the following results of increase in weight (lbs) we observed in 8 children as follows.

Food A(x)	49	53	51	52	47	50	52	53
Food B(y)	52	55	52	53	50	54	54	53

Examine the significance of increase in weight of children due to food B.

Solution

$H_0 : \mu_1 = \mu_2$, there is no significant difference between the two foods.

$H_1 : \mu_1 \neq \mu_2$, there is significant difference between the two foods.

Level of significance = 5%

test statistic

$$t = \frac{|\bar{d}|}{S/\sqrt{n}} \sim t(n-1) d.f$$

x	y	d=x-y	d ²
49	52	-3	9
53	55	-2	4
51	52	-1	1
51	52	-1	1
47	50	-3	16
50	54	-4	16
52	54	-2	4
53	53	0	0
Total		-16	44

$$\bar{d} = \frac{\sum di}{n} = \frac{-16}{8} = -2,$$

$$S = \sqrt{\frac{\sum di^2 - \frac{(\sum di)^2}{n}}{n-1}} = 1.3093$$

$$t_{cal} = \frac{|-2|}{1.3093/\sqrt{8}} = 4.32$$

Table value:

$t(8-1)$ d.f at 5% l.o.s = 2.365

Inference:

$$t_{cal} > t_{tab}$$

We reject the null hypothesis H_0

(i.e) there is significant difference between the two foods A and B.

Learning Exercise

- 10 samples of leaves of the plant are chosen at random from a large population and their weight in grams are found to be as follows

63 63 64 65 66 69 69 70 70 71

From this data mean weight in universe is 65 gm. Can we assume this mean weight?

- A health status survey in a few villages revealed that the normal serum protein value of children in that locality is 7.0 g/100ml. A group of 16 children, who received high protein food for a period of 6 months had serum protein values shown below. Can we consider that the mean serum protein level of these who were fed on high protein diet is different from that of the general population.

Children	1	2	3	4	5	6	7	8	9	10
Protein level g %	7.1	7.70	8.2	7.56	7.05	7.08	7.21	7.25	7.36	6.59
Children	11	12	13	14	15	16				
Protein level g %	6.85	7.9	7.27	6.56	7.93	8.5				

- The following data related to the rate of diffusion of CO_2 through two series of different porosity, find out whether the diffusion rate same for both sides.

Diffusion through fine soil (x_1)	20	31	31	23	28	23	26	27	26	17	17	25
Diffusion through coarse soil (x_2)	19	30	32	28	15	26	35	18	25	27	35	34

- A new variety of cotton was evolved by a breed. In order to compare its yielding ability with that of a ruling variety, an experiment was conducted in Completely Randomised Design. The yield (kg/plot) was observed. The summary of the results are given below. Test whether the new variety of cotton gives higher yield than the ruling variety.

New Variety	$n_1 = 9$	$\bar{x}_1 = 28.2$	$S_1^2 = 5.4430$
Ruling Variety	$n_2 = 11$	$\bar{x}_2 = 25.9$	$S_2^2 = 1.2822$

5. The iron contents of fruits before and after applying farm yard manure were observed as follows.

Fruit No:	1	2	3	4	5	6	7	8	9	10
Before Applying	7.7	8.5	7.2	6.3	8.1	5.2	6.5	9.4	8.3	7.5
After Applying	8.1	8.9	7.0	6.1	8.2	8.0	5.8	8.9	8.7	8.0

Is there any Significant difference between the mean iron contents in the fruits before & after the farm yarn manure?

1. 10 samples of leaves of the plant are chosen at random from a large population and their weight in grams are found to be as follows

63	63	64	65	66	69	69	70	70	71
----	----	----	----	----	----	----	----	----	----

From this data mean weight in universe is 65 gm. Can we assume this mean weight?

2. A health status survey in a few villages revealed that the normal serum protein value of children in that locality is 7.0 g/100ml. A group of 16 children, who received high protein food for a period of 6 months had serum protein values shown below. Can we consider that the mean serum protein level of these who were fed on high protein diet is different from that of the general population.

Children	1	2	3	4	5	6	7	8	9	10
Protein level g %	7.1	7.70	8.2	7.56	7.05	7.08	7.21	7.25	7.36	6.59
Children	11	12	13	14	15	16				
Protein level g %	6.85	7.9	7.27	6.56	7.93	8.5				

3., The following data related to the rate of diffusion of CO₂ through two series of different porosity, find out whether the diffusion rate same for both sides.

Diffusion through fine soil (x ₁)	20	31	31	23	28	23	26	27	26	17	17	25
Diffusion through coarse soil (x ₂)	19	30	32	28	15	26	35	18	25	27	35	34

4. A new variety of cotton was evolved by a breed. In order to compare its yielding ability with that of a ruling variety, an experiment was conducted in Completely Randomised Design. The yield (kg/plot) was observed. The summary of the results are given below. Test whether the new variety of cotton gives higher yield than the ruling variety.

New Variety	n ₁ = 9	$\bar{x}_1 = 28.2$	S ₁ ² = 5.4430
Ruling Variety	n ₂ = 11	$\bar{x}_2 = 25.9$	S ₂ ² = 1.2822

5. The iron contents of fruits before and after applying farm yard manure were observed as follows.

Is there any Significant difference between the mean iron contents in the fruits before & after the farm yarn manure?

Exercise.8

Chi square test – test for association and goodness of fit

χ^2 – test for goodness of fit

If O_i , ($i=1,2,\dots,n$) is a set of observed (experimental frequencies) and E_i ($i=1,2,\dots,n$) is the corresponding set of expected (theoretical or hypothetical) frequencies, then,

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{(n-1)} \text{ df}$$

Example1:

The number of yeast cells counted in a haemocytometer is compared to the theoretical value is given below. Does the experimental result support the theory?

No. of Yeast cells in the square	Observed Frequency	Expected Frequency
0	103	106
1	143	141
2	98	93
3	42	41
4	8	14
5	6	5

Solution:

H_0 : the experimental results support the theory

H_1 : the experimental results does not support the theory.

Level of significance=5%

Test Statistic:

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{(n-1) \text{ df}}$$

O_i	E_i	$O_i \cdot E_i$	$(O_i - E_i)^2$	$(O_i - E_i)^2 / E_i$
103	106	-3	9	0.0849
143	141	2	4	0.0284
98	93	5	25	0.2688
42	41	1	1	0.0244
8	14	-6	36	2.5714
6	5	1	1	0.2000
400	400			3.1779

$$\therefore \chi^2 = 3.1779$$

Table value:

$$\chi^2_{(6-1=5 \text{ at } 5\% \text{ l.os})} = 11.070$$

Inference

$$\chi^2_{\text{cal}} < \chi^2_{\text{tab}}$$

We accept the null hypothesis.

(i.e) there is a good correspondence between theory and experiment.

χ^2 test for independence of attributes

Formula:

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{(O_{ij} - E_{ij})^2}{E_{ij}} \sim \chi^2_{(m-1)(n-1) \text{ df}}$$

O_{ij} – observed frequencies

E_{ij} – Expected frequencies

m= number of rows

n= number of columns

$$\sum O_{ij} = \sum E_{ij}$$

Exapmle 2:

The severity of a disease and blood group were studied in a research project. The findings sre given in the following table, knowmn as the m xn contingency table. Can this severity of the condition and blood group are associated.

Severity of a disease classified by blood group in 1500 patients.

Condition	Blood Groups				Total
	O	A	B	AB	
Severe	51	40	10	9	110
Moderate	105	103	25	17	250
Mild	384	527	125	104	1140
Total	540	670	160	130	1500

Solution:

H₀: The two attributes severity of the condition and blood groups are not associated.

H₁: The two attributes severity of the condition and blood groups are associated.

Calculation of Expected frequencies

Condition	Blood Groups				Total
	O	A	B	AB	
Severe	39.6	49.1	11.7	9.5	110
Moderate	90.0	111.7	26.7	21.7	250
Mild	410.4	509.2	121.6	98.8	1140

Total	540	670	160	130	1500
-------	-----	-----	-----	-----	------

Test statistic:

$$\chi^2 = \sum_{i=1}^m \sum_{j=1}^n \frac{(o_{ij} - E_{ij})^2}{E_{ij}} \sim \chi^2_{(m-1)(n-1)} \text{ df}$$

Here m= 3.,n =4

Calculations:

O _i	E _i	O _i ·E _i	(O _i ·E _i) ²	(O _i ·E _i) ² /E _i
51	39.6	11.4	129.96	3.2818
40	49.1	-9.1	82.81	1.6866
10	11.7	-1.7	2.89	0.2470
9	9.5	-0.5	0.25	0.0263
105	90.0	15	225.00	2.5000
103	111.7	-8.7	75.69	0.6776
25	26.7	-1.7	2.89	0.1082
17	21.7	-4.7	22.09	1.0180
384	410.4	-26.4	696.96	1.6982
527	509.2	17.8	316.84	0.6222
125	121.6	3.4	11.56	0.0951
104	98.8	5.2	27.04	0.2737
				12.2347

$$\therefore \chi^2 = 12.2347$$

Table value:

$$\chi^2_{(3-1)(4-1)} = \chi^2_{(6)} \text{ at } 5\% \text{ l.os} = 12.59$$

Inference

$$\chi^2_{\text{cal}} < \chi^2_{\text{tab}}$$

We accept the null hypothesis.

(i.e) the two attributes severity of the condition and blood group are independent.

2x2 – contingency table

$$\chi^2 = \frac{N(ad - bc)^2}{(a + b)(c + d)(a + c)(b + d)} \sim \chi^2_{(2-1)(2-1) \text{ df}} = \chi^2_{(1) \text{ df}}$$

Example 3:

In order to determine the possible effect of a chemical treatment on the rate of germination of cotton seeds a pot culture experiment was conducted. The results are given below

Chemical treatment and germination of cotton seeds

	Germinated	Not germinated	Total
Chemically Treated	118	22	140
Untreated	120	40	160
Total	238	62	300

Does the chemical treatment improve the germination rate of cotton seeds at 1 % level?

Solution:

H₀: The chemical treatment does not improve the germination rate of cotton seeds.

H₁: The chemical treatment improves the germination rate of cotton seeds.

L.O.S = 1 %

Test statistic

$$\chi^2 = \frac{N(ad - bc)^2}{(a + b)(c + d)(a + c)(b + d)} \sim \chi^2_{(2-1)(2-1) \text{ df}} = \chi^2_{(1) \text{ df}}$$

$$\chi^2 = \frac{300(118 \times 40 - 22 \times 120)^2}{140 \times 160 \times 62 \times 238} = 3.927$$

Table value:

$$\chi^2 (1) \text{ df at } 1\% \text{ L.O.S} = 6.635$$

Inference

$$\chi^2_{\text{cal}} < \chi^2_{\text{tab}}$$

We accept the null hypothesis.

(i.e)The chemical treatment will not improve the germination rate of cotton seeds significantly.

Yates correction for continuity

$$\text{Then use, } \chi^2 = \frac{N(ad - bc)^2}{(a + b)(c + d)(a + c)(b + d)} \sim \chi^2 (1) \text{ df}$$

(or)

Directly use the χ^2 - statistic as

$$\chi^2 = \frac{N \left(|ad - bc| - \frac{N}{2} \right)^2}{(a + b)(c + d)(a + c)(b + d)} \sim \chi^2 (1) \text{ df}$$

Example 4:

In an experiment on the effect of a growth regulator on fruit setting in muskmelon the following results were obtained. Test whether the fruit setting in muskmelon and the application of growth regulator are independent at 1% level.

	Fruit set	Fruit not set	Total
Treated	16	9	25

Control	4	21	25
Total	20	30	50

Solution:

H_0 :Fruit setting in muskmelon does not depend on the application of growth regulator.

H_1 : Fruit setting in muskmelon depend on the application of growth regulator.

L.O.S = 1 %

Tet statistic

$$\chi^2 = \frac{N \left(|ad - bc| - \frac{N}{2} \right)^2}{(a+b)(c+d)(a+c)(b+d)} \sim \chi^2 (1) \text{ df}$$

$$\chi^2 = \frac{50 \left[|16 \times 21 - 9 \times 4| - \frac{50}{2} \right]^2}{25 \times 25 \times 20 \times 30} = 10.08$$

Table value:

$$\chi^2 (1) \text{ df at } 1 \% \text{ L.O.S} = 6.635$$

Inference

$$\chi^2_{\text{cal}} > \chi^2_{\text{tab}}$$

We reject the null hypothesis.

(i.e) Fruit setting in muskmelon is influenced by the growth regulator.

Learning Exercise

1. The theory predicts the proportion of beans in the 4 groups A, B, C, D should be 9:3:3:1. In an experiment among 1600 beans, the number in the four groups were 882, 313, 287 and 118. Does the experimental result support the theory.

2. A study was conducted, among 100 professors from 3 different divisions for the preference on beverages of 3 categories test if there is any relationship between the field of teaching and preference of beverage.

Field of teaching				
Beverage	Business	Social Sciences	Agri	Total
Tea	20	10	10	40
Coffee	10	10	15	35
Cold drinks	10	8	7	25
Total	40	28	32	100

3. A random sample of 600 students from Delhi University are selected and asked their opinion about autonomous Status of Colleges. The results were given below. Test the hypothesis at 5% level that opinions are independent of class groupings.

Class grouping	Favour of	Against	
Commerce	120	80	200
Science	130	70	200
Arts	70	30	100
Total	400	200	600

4. In a survey of preference of new coverage 100 persons are collected and taste preference of average was surveyed according to sex of the person. We conclude that the taste preference and sex of the person are associated.

	Male	Female	

Favour	35	25	60
Against	25	15	40
Total	60	40	100

5. Two new food stuffs were introduced and public opinion was sought based on the taste of the food stuff. The results are given below. Examine whether there is an association between the category of the food stuffs and the taste of the food stuffs.

Category			
	A	B	
Tasty	620	380	1000
Not tasty	550	450	1000
Total	1170	830	2000

Exercise.9

Calculation of Karl Pearson’s correlation coefficient

Pearsons Correlation coefficient

. The correlation coefficient r is known as Pearson’s correlation coefficient as it was discovered by Karl Pearson.

$$r = \frac{\frac{1}{n-1} (\sum (x - \bar{x})(y - \bar{y}))}{\sqrt{\frac{1}{n-1} \sum (x - \bar{x})^2} \sqrt{\frac{1}{n-1} \sum (y - \bar{y})^2}}$$

Which can be simplified as

$$r = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sqrt{\sum x^2 - \frac{(\sum x)^2}{n}} \sqrt{\sum y^2 - \frac{(\sum y)^2}{n}}}$$

Testing the significance of r

The significance of r can be tested by Student’s t test. The test statistics is given by

$$t = \frac{|r|}{\sqrt{\frac{1-r^2}{n-2}}}$$

Example.1

Compute Pearsons coefficient of correlation between advertisement cost and sales as per the data given below:

Advertisement Cost in 1000’s	39	65	62	90	82	75	25	98	36	78
Sales in lakhs	47	53	58	86	62	68	60	91	51	84

Solution

H_0 : The correlation coefficient r is not significant

H_1 : The correlation coefficient r is significant.

Level of significance 5%

From the data

$n = 10$

$$\sum x = 650 \quad \sum y = 660 \quad \sum xy = 45604 \quad \sum x^2 = 47648 \quad \sum y^2 = 45784$$

$$r = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sqrt{\sum x^2 - \frac{(\sum x)^2}{n}} \sqrt{\sum y^2 - \frac{(\sum y)^2}{n}}}$$

$$= \frac{45604 - \frac{(650)(660)}{10}}{\sqrt{47648 - \frac{(650)^2}{10}} \sqrt{45784 - \frac{(660)^2}{10}}}$$

$$= \frac{45604 - 42900}{(73.47)(47.1)} = 0.7804$$

Correlation coefficient is positively correlated.

Test Statistic

$$t = \frac{|r|}{\sqrt{\frac{1-r^2}{n-2}}} \sim (n-2) \text{ d.f}$$

$$t = \frac{0.7804}{\sqrt{\frac{1 - (0.7804)^2}{10 - 2}}} = 3.530$$

$$t_{\text{tab}} = t_{(10-2, 5\% \text{los})} = 2.306$$

Inference

$t_{\text{cal}} > t_{\text{tab}}$, we reject null hypothesis.

∴ The correlation coefficient r is significant. (i.e) There is a relation between advertisement company and the sales.

Learning Exercise

1. Calculate the simple correlation coefficient between wing length & tail length of the following 12 birds of a particular species. Also test its significant.

Wing length (cm)x	1	2	3	4	5	6	7	8	9	10	11	12
	10.4	10.8	11.1	10.2	10.3	10.2	10.7	10.5	10.8	11.2	10.6	11.4
Tail length (cm)y	7.4	7.6	7.9	7.2	7.4	7.1	7.4	7.2	7.8	7.7	7.8	8.3

2. The data refer to the yield of grain in gms|plant(y) and the number of productive tillers (x) and 15 paddy plants

Y	37	20	42	36	20	30	26	21	43	44	22	31	26	37	26
X	15	12	17	14	12	13	12	9	24	20	14	18	13	15	7

Find the correlation

3. The following data relates to the yield in grams(y) and the matured pods (x) of 10 groundnut plants. Work out the correlation coefficient and test its significance.

X:	14	34	20	16	11	11	20	17	22	17
Y:	16	40	21	18	14	13	20	35	17	27

4. Find the persons coefficient of correlation between price and demand from the following data.

Price	11	13	15	17	18	19	20
Demand	30	29	24	24	21	18	15

Practical.10
Fitting of simple linear regression of y on x

Testing the significance of regression co-efficient

To test the significance of the regression coefficient we can apply either a t test or analysis of variance (F test). The ANOVA table for testing the regression coefficient will be as follows:

Sources of variation	DF	SS	MS	F ratio
Due to regression	1	SS(b)	S_b^2	S_b^2 / S_e^2
Deviation from regression	n-2	SS(Y)-SS(b)	S_e^2	
Total	n-1	SS(Y)		

In case of t test the test statistic is given by

$$t = b / SE(b) \text{ where } SE(b) = s_e^2 / SS(X)$$

Example 1

Form a paddy field, 36 plants were selected at random. The length of panicles(x) and the number of grains per panicle (y) of the selected plants were recorded. The results are given below. Fit a regression line y on x. Also test the significance (or) regression coefficient.

The length of panicles in cm (x) and the number of grains per panicle (y) of paddy plants.

S.No.	Y	X	S.No.	Y	X	S.No.	Y	X
1	95	22.4	13	143	24.5	25	112	22.9
2	109	23.3	14	127	23.6	26	131	23.9
3	133	24.1	15	92	21.1	27	147	24.8
4	132	24.3	16	88	21.4	28	90	21.2
5	136	23.5	17	99	23.4	29	110	22.2
6	116	22.3	18	129	23.4	30	106	22.7
7	126	23.9	19	91	21.6	31	127	23.0
8	124	24.0	20	103	21.4	32	145	24.0
9	137	24.9	21	114	23.3	33	85	20.6
10	90	20.0	22	124	24.4	34	94	21.0

11	107	19.8	23	143	24.4	35	142	24.0
12	108	22.0	24	108	22.5	36	111	23.1

Null Hypothesis H_0 : regression coefficient is not significant.

Alternative Hypothesis H_1 : regression coefficient is significant.

$$\sum y = 4174 \quad \sum y^2 = 496258 \quad \bar{y} = \frac{\sum y}{n} = 115.94$$

$$\sum x = 822.9 \quad \sum x^2 = 18876.83 \quad \bar{x} = \frac{\sum x}{n} = 22.86$$

$$\sum xy = 96183.4$$

$$SS(Y) = \sum y^2 - \frac{(\sum y)^2}{n} = 496258 - \frac{(4174)^2}{36} = 12305.8889$$

$$SS(X) = \sum x^2 - \frac{(\sum x)^2}{n} = 18876.83 - \frac{(822.9)^2}{36} = 66.7075$$

The regression line y on x is $\bar{y} = a_1 + b_1 \bar{x}$

$$b_1 = \frac{\sum xy - \frac{\sum x \sum y}{n}}{\sum x^2 - \frac{(\sum x)^2}{n}} = \frac{96183.4 - \frac{(822.9)(4174)}{36}}{66.7075} = 11.5837$$

$$\bar{y} = a_1 + b_1 \bar{x}$$

$$115.94 = a + (11.5837)(22.86)$$

$$A = 115.94 - 264.8034$$

$$A = -148.8633$$

The fitted regression line is $y = -148.8633 + 11.5837x$

$$SS(b) = \frac{\left(\sum xy - \frac{\sum x \sum y}{n} \right)^2}{\sum x^2 - \frac{(\sum x)^2}{n}} = \frac{(722.7167)^2}{66.7075} = 8950.8841$$

ANOVA Table:

Sources of Variation	d.f	SS	MSS	F-value
Replication	1	8950.8841	8950.8841	90.7093
Error	36-2=34	3355.0048	98.6766	
Total	35	12305.8889		

For t-test

$$t = \frac{b}{SE(b)} \sim t_{(n-2)} d.f$$

$$SE(b) = \sqrt{\frac{Se^2}{SS(X)}} = \sqrt{\frac{98.6776}{66.7075}} = 1.2162$$

$$t = \frac{11.5837}{1.2162} = 9.5245$$

Table Value:

$t_{(n-2)}$ d.f= t_{34} d.f at 5% level=2.032

$t_{cal} > t_{tab}$. we reject H_0 .

Hence t is significant.

Learning Exercise

1. The following data are using length of 13 sparrows of various ages.

Age (days) (x)	3	4	5	6	8	9	10	11	12	14	15	16	17
Wing length (cm)	1.4	1.5	2.2	2.4	3.1	3.2	3.2	3.9	4.1	4.7	4.5	5.2	5.0

Y

Fit the regression line of y on x. Also test the significance of the regression coefficient.

2. Obtain the regression line of the form $y = a + bx$ between the average number of tillers (x) and the yield in kgs (y) of turmeric crop from the following data

Average number of tillers(x)	3.5	3.2	3.5	3.8	3.6	3.74	2.8	4.2	4.0	4.5
Yield in (Kgs) (y)	2.0	1.8	1.9	2.1	2.0	2.3	1.7	2.5	2.6	3.0

3. Find out the regression equations y on x between the number of root fibers(x) and yields in kgs(y) of ginger crop from the following data. Also test the significance of regression coefficient.

Number of root fibers (x)	6	5	3	7	4	9	10	11	15	8
Yield in kgs(y)	1.0	0.8	0.5	1.1	0.6	1.2	1.5	1.6	1.9	0.9

4. Find out the regression equation of y on x and test the significance of regression coefficient.

X	2.1	2.7	3.5	4.6	4.9	5.7	6.0	7.4	8.3	4.8
Y	11.2	12.8	16.4	17.8	19.5	21.1	22.8	26.3	31.0	21.1

Practical.11

Formation of ANOVA table for completely Randomised design (CRD) with equal replication and comparison of means using critical difference values

Completely Randomized Design (CRD)

CRD is the basic single factor design. In this design the treatments are assigned completely at random so that each experimental unit has the same chance of receiving any one treatment. But CRD is appropriate only when the experimental material is homogeneous. As there is generally large variation among experimental plots due to many factors CRD is not preferred in field experiments.

In laboratory experiments and greenhouse studies it is easy to achieve homogeneity of experimental materials and therefore CRD is most useful in such experiments.

Layout of a CRD

Completely randomized Design is the one in which all the experimental units are taken in a single group which are homogeneous as far as possible.

The randomization procedure for allotting the treatments to various units will be as follows.

Step 1: Determine the total number of experimental units.

Step 2: Assign a plot number to each of the experimental units starting from left to right for all rows.

Step 3: Assign the treatments to the experimental units by using random numbers.

The statistical model for CRD with one observation per unit

$$Y_{ij} = \mu + t_i + e_{ij}$$

μ = overall mean effect

t_i = true effect of the i th treatment

e_{ij} = true effect of the j th unit receiving i th treatment

The arrangements of data in CRD is as follows:

	Treatments				
	T ₁	T ₂	T _i	T _K	
	Y ₁₁	Y ₂	Y _{i1}	Y _{K1}	
	Y ₁₂	Y ₂₂	Y _{i2}	Y _{K2}	
	Y _{1r1}	Y _{2r2}	Y _{iri}	Y _{k rk}	
Total	T ₁	T ₂	T _i	T _k	GT

The null hypothesis will be

H₀: μ₁ = μ₂=.....=μ_k or There is no significant difference between the treatments

And the alternative hypothesis is

H₁: μ₁ = μ₂=.....=μ_k There is significant difference between the treatments

The different steps in forming the analysis of variance table for a CRD are:

$$1. \quad C.F = \frac{Y_i^2}{n} = \frac{(GT)^2}{n}$$

n= Total number of observations

$$2. \quad \text{Total SS} = \text{TSS} = \sum_{i=1}^k \sum_{j=1}^{v_i} Y_{ij}^2 - C.F$$

$$3. \quad \text{Treatment SS} = \frac{Y_1^2}{r_1} + \frac{Y_2^2}{r_2} + \dots + \frac{Y_k^2}{r_k} - \frac{Y^2}{n}$$

$$= \sum_{i=1}^k \frac{y_i^2}{r_i} - C.F$$

$$4. \quad \text{Error SS} = \sum_{i=1}^k \sum_{j=1}^{r_i} y_{ij}^2 - \sum_{i=1}^k \frac{y_i^2}{r_i}$$

$$\text{Error SS} = \text{TSS} - \text{TrSS}$$

5. Form the following ANOVA table and calculate F value.

Source of variation	df	SS	Ms	F
Treatments	t-1	TrSS	$TrMS = \frac{TSS}{t-1}$	$\frac{TrMS}{EMS}$
Error	n-t	ESS	$EMS = \frac{ESS}{n-t}$	
Total	n-1	TSS		

6. Compare the calculate F with the critical value of F so that acceptance or rejection of the null hypothesis can be determined.

7. If null hypothesis is rejected that indicates there is significant differences between the different treatments.

8. Calculate C D value.

$$C.D. = t \times SE (d)$$

$$\text{where S.E (d)} = \sqrt{EMS \left(\frac{1}{r_i} + \frac{1}{r_j} \right)}$$

and t is the critical t value for error degrees of freedom at 5% level.

Problem

The following table gives the yield in kgs per plot of five varieties of wheat after being applied to each of four plots in a completely randomized design.

Varieties	Yield in kgs				Totals	Treatment means
A	8	8	6	10	32 (T1)	8 (T1)
B	10	12	13	9	44 (T2)	11 (T2)
C	18	17	13	16	64 (T3)	16 (T3)
D	12	10	15	11	48 (T4)	12 (T4)
E	8	11	9	8	36 (T5)	9 (T5)
Grand Total					224	

Analysis

1. Correction Factor (C.F) = $\frac{(\text{Grand total})^2}{r \times t}$ where r is the number of replications per treatment and t is the number of treatments.

$$= \frac{(224)^2}{4 \times 5} = 2508.8$$

2. Total sum of squares TSS = $\sum_{i=1}^k \sum_{j=1}^{v_i} Y_{ij}^2 - C.F$

$$= 8^2 + 8^2 + \dots + 9^2 + 8^2 - 2508.8 = 207.2$$

3. varieties Sum of squares VSS = $\frac{T_1^2 + T_2^2 + T_3^2 + T_4^2 + T_5^2}{r} - C.F$

$$= \frac{32^2 + 44^2 + 64^2 + 48^2 + 36^2}{4} - 2508.8$$

4. Error sum of Squares = Total sum of squares – Variety sum of squares = 207.2 - 155.2 = 52.0

Table for Analysis of Variance

Source of variation	D.F.	S.S.	MS (variance)	F (variance ratio)	F at 5%
Between varieties	4	155.2	38.80	11.80*	3.06
Within varieties (error)	15	52.0	3.47		
Total	19	207.2			

*Significant at 5% level of significances

Here, F test indicates that there are significant difference between the variety means since the observed value of the variance ratio is significant at 5% level of significance. Now we wish to know as to which variety is the best and also which varieties show the significant difference among themselves. This can be done with the help of critical difference (c.d.)

Now, standard error of the difference between two treatment means is

S.E.d = $\sqrt{\frac{2XEMS}{r}}$ where E Ms is the error mean square and 'r' is the no. of replications.

$$= \sqrt{\frac{2X3.47}{4}} = 1.32$$

∴ Critical difference (CD) = SEd x 't' 5% value for error d.f.

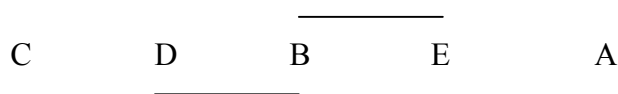
$$= 1.32 \times 2.131 = 2.81$$

Since the 'F' test in the analysis of variance indicated significant differences between the varieties, we are justified in comparing the varieties with the help of the above C.D. value.

Varieties	A	B	C	D	E	CD
Many yields	8	11	16	12	9	2.81

Conclusions represented symbolically (bar chart notation):

The varieties can be compared by setting them in the descending (or ascending) order of their mean yields in the following manner.



The varieties which do not differ significantly have been underlined by a bar. This method of underlining the treatments which do not differ significantly is the way of indicating the significance and non-significance of individual comparisons.

Learning Exercise

1) The following table gives the yields of five varieties of paddy with four replications each by using completely randomized Design.

Varieties	Yield in kg.			
A	8	8	6	10
B	10	12	13	9
C	18	17	13	16
D	12	10	15	11
E	8	11	9	8

Analyse the data and draw your conclusions.

2) Below are given the plan and yield in kg per plot of a completely randomized design for testing the effect of five different fertilizers A, B, C, D & E.

D	E	B	E	D
20	17	21	16	15
A	C	A	D	B
8	17	9	13	17
B	D	E	C	A
12	19	18	18	15
C	A	C	A	E
16	8	18	10	15
E	B	D	B	C

Statistics

13	16	23	14	19
----	----	----	----	----

Analyse the data and state your conclusions.

Exercise.12

Formation of ANOVA table for Randomised blocks design (RBD) and comparison of means using critical difference values

Randomized blocks design (RBD)

When the experimental material is heterogeneous, the experimental material is grouped into homogenous sub-groups called blocks. As each block consists of the entire set of treatments a block is equivalent to a replication.

If the fertility gradient runs in one direction say from north to south or east to west then the blocks are formed in the opposite direction. Such an arrangement of grouping the heterogeneous units into homogenous blocks is known as randomized blocks design.

Layout of RBD

Each block consists of as many experimental units as the number of treatments. The treatments are allocated randomly to the experimental units within each block independently such that each treatment occurs once. Each treatment is replicated as many times as the number of blocks or the number of blocks are chosen to be equal to the number of replications for the treatments.

The analysis of variance model for RBD is

$$Y_{ij} = \mu + t_i + r_j + e_{ij}$$

where

μ = the overall mean

t_i = the i th treatment effect

r_j = the j th replication effect

e_{ij} = the error term

Analysis of RBD

The results of RBD can be arranged in a two way table according to the replications (blocks) and treatments.

There will be $r \times t$ observations in total where r stands for number of replications and t for number of treatments. .

The data are arranged in a two way table form by representing treatments in rows and replications in columns.

Treatment	Replication					Total
	1	2	3	r	
1	Y11	Y12	Y13	Y1r	T1
2	Y21	Y22	Y23	Y2r	T2
3	Y31	Y32	Y33	Y3r	T3
t	Yt1	Yt2	Yt3	Ytr	Tt
Total	R1	R2	R3		Rr	G.T

In this design the total variance is divided into three sources of variation viz., between replications, between treatments and error

$$C.F = \frac{(GT)^2}{rt}$$

$$\text{Total SS} = \sum \sum Y_{ij}^2 - C-F$$

$$\text{Replication SS} = \frac{1}{t} \sum R_j^2 - C-F$$

$$\text{Treatments SS} = \frac{1}{r} \sum T_i^2 - C-F$$

$$\text{Error SS} = \text{Total SS} - \text{Replication SS} - \text{Treatment SS}$$

The skeleton ANOVA table for RBD with t treatments and r replications

Sources of variation	d-f	SS	ms	F Value
Replication	r-1	R S S	R M S	R M S/ E M S
Treatment	t-1	Tr S S	Tr M S	Tr M S/E M S
Error	(r-1) (t-1)	E S S	E M S	
Total	Rt -1	T S S		

$$CD = SE (d) \times t$$

$$\text{Where S.E (d)} = \sqrt{\frac{2EMS}{r}}$$

t = critical value of t for a specified level of significance and error degrees of freedom

Based on the CD value the bar chart can be drawn.

From the bar chart conclusion can be written.

Problem

The yields of six nitrogen treatments on a crop in kgs along with the plan of the experiment are given below. The number of blocks is five and the nitrogen treatments have been represented by A, B, C, D, E and F.

Block I	Block II	Block III	Block IV	Block V
D 17 C 12	B 12 C 15	E 23 A 30	A 28 F 64	F 75 C 14
F 70 B 6	E 26 A 26	C 16 D 20	B 9 D 23	D 20 B 7
A 20 E 28	D 10 F 62	F 56 B 10	E 33 C 14	E 30 A 23

It is required to analyse the data.

Analysis

i) Tabulation of the data

The first step in the analysis of data is to tabulate yield figures according to block and treatments in the follow manner.

Varieties	Blocks					Treatment totals	Treatment means
	I	II	III	IV	V		
A	20	26	30	28	23	1277 (T ₁)	25.4
B	9	12	10	9	7	47 (T ₂)	9.4
C	12	15	16	14	14	71 (T ₃)	14.2
D	17	10	20	23	20	90 (T ₄)	18.0
E	28	26	23	35	30	142 (T ₅)	28.4
F	70	62	56	64	75	327 (T ₆)	65.4
Totals	156 (B1)	151 (B2)	155 (B3)	173 (B4)	169 (B5)	804 (GT)	

ii) Sums of squares for different sources

a. Correction factor (CF) = $\frac{(GT)^2}{bxt}$ Where GT is the grand total; 'b' blocks; 't'

No of treat ments = $\frac{(804)^2}{5x6} = 21547.2$

b. Total S.S = S.S of all the observation – CF
 $= 20^2 + 9^2 + \dots + 75^2 - 21547.2 = 10466.8$

c. S.S. due to blocks (Bss) = $\frac{B_1^2 + B_2^2 + \dots + B_5^2}{t} - CF$

$$= \frac{156^2 + 151^2 + \dots + 169^2}{6} - 21547.2$$

$$= 61.4$$

$$\text{d. S.S due to treatments (tss)} = \frac{T_1^2 + T_2^2 + \dots + T_5^2}{r} - CF$$

$$= \frac{127^2 + 47^2 + \dots + 327^2}{6} - 21547.2$$

$$= 10167.2$$

$$\text{e. S.S. due to error} = \text{Total SS} - \text{Bss} - \text{tss} = 10646.8 - 61.4 - 10167.2$$

$$= 418.2$$

iii) Table of analysis of variance

Now these values will be set down in a table of analysis of variance as given below:

Analysis of variance table

Source of variation	D.F	S.S	M.S	Variance ratio 'F'	F 15%
Blocks	4	61.4	15.35		
Treatments	5	10167.2	2033.44	97.24*	2.71
Error	20	418.2	20.19		
Total	29	10646.8			

* Significant at 5% level of significance

It is clear from the table that this observed value of 'F' is significant at 5% level of significance which proves that there are significant differences between the treatment means. Now, we have to test the significance of the difference between the individual treatments, and this will be done with the help of C, D as usual.

iv) Critical difference

S.E. of the difference between any two treatment means is

$$SE_d = \sqrt{\frac{2XEMS}{r}} = \sqrt{\frac{2 \times 20.91}{5}} = 2.89$$

∴ Critical difference = $SE_d \times t_{5\%} = 2.89 \times 2.086 = 6.03$

v) Conclusions represented symbolically

The treatments have been compared by setting them in the descending order of their mean yields in the following manner

Varieties	F	E	A	D	C	B
Mean yield	65.4	<u>28.4</u>	<u>25.4</u>	18.0	<u>14.2</u>	<u>9.4</u>

The treatments which do not differ significantly have been underlined by a bar. The treatment ‘F’ has been found to be the best.

Learning Exercise

- 1) The yield of rice (in kg) with five fertilizers tested in four blocks using RBD is given the following layout. Analyse the data & interpret your conclusion.

Block 1	Block 2	Block 3	Block 4
B 10	C 13	A 19	D 20
C 16	A 21	D 24	E 36
A 20	D 21	E 32	B 9
D 23	E 31	B 10	C 13
E 33	B 11	C 14	A 24

- 2) An experiment was conducted in RBD to study to comparative performance of yield of six varieties of oranges (kg/plot) are given below. Analyse the data and give your conclusion.

Treatments	Blocks				
	B1	B2	B3	B4	B5
V1	5.5	5.9	6.3	6.5	6.7
V2	7.4	7.7	7.9	7.5	8.1
V3	4.6	5.1	5.3	4.9	4.7
V4	5.0	5.8	5.6	6.1	5.3
V5	6.7	6.2	6.9	6.8	6.0
V6	8.2	7.9	7.5	7.2	6.9

Exercise.13
Formation of ANOVA table for Latin square design (LSD) and comparison of means using critical difference values

Latin Square Design

When the experimental material is divided into rows and columns and the treatments are allocated such that each treatment occurs only once in each row and each column, the design is known as L S D.

In LSD the treatments are usually denoted by A B C D etc.

For a 5 x 5 LSD the arrangements may be

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>B</i>	<i>A</i>	<i>E</i>	<i>C</i>	<i>D</i>
<i>C</i>	<i>D</i>	<i>A</i>	<i>E</i>	<i>B</i>
<i>D</i>	<i>E</i>	<i>B</i>	<i>A</i>	<i>C</i>
<i>E</i>	<i>C</i>	<i>D</i>	<i>B</i>	<i>A</i>

Square 1

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>B</i>	<i>A</i>	<i>D</i>	<i>E</i>	<i>C</i>
<i>C</i>	<i>E</i>	<i>A</i>	<i>B</i>	<i>D</i>
<i>D</i>	<i>C</i>	<i>E</i>	<i>A</i>	<i>B</i>
<i>E</i>	<i>D</i>	<i>B</i>	<i>C</i>	<i>A</i>

Square 2

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>A</i>
<i>C</i>	<i>D</i>	<i>E</i>	<i>A</i>	<i>B</i>
<i>D</i>	<i>E</i>	<i>A</i>	<i>B</i>	<i>C</i>
<i>E</i>	<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>

Square 3

Analysis

The ANOVA model for LSD is

$$Y_{ijk} = \mu + r_i + c_j + t_k + e_{ijk}$$

r_i is the i th row effect

c_j is the j th col effect

t_k is the k th treatment effect

The analysis of variance table for LSD is as follows:

Sources of Variation	Degrees of Freedom	Sum of Squares	Mean Squares	F-ratio
Rows	t-1	RSS	RMS	RMS/EMS
Columns	t-1	CSS	CMS	CMS/EMS
Treatments	t-1	TrSS	TrMS	TrMS/EMS
Error	(t-1)(t-2)	ESS	EMS	
Total	t ² -1	TSS		

F table value

F_[t-1,(t-1)(t-2)] degrees of freedom at 5% or 1% level of significance

Steps to calculate the above Sum of Squares are as follows:

$$\text{Correction Factor (CF)} = \frac{(\text{GrandTotal})^2}{(\text{treatment})^2}$$

$$\text{Total Sum of Squares (TSS)} = \sum (y_{ijk})^2 - CF$$

$$\text{Row sum of squares (RSS)} = \frac{1}{t} \sum_{i=1}^t (R_i)^2 - CF$$

$$\text{Column sum of squares (CSS)} = \frac{1}{t} \sum_{j=1}^t (C_j)^2 - CF$$

$$\text{Treatment sum of squares (TrSS)} = \frac{1}{t} \sum_{k=1}^t (T_k)^2 - CF$$

$$\text{Error Sum of Squares} = \text{TSS} - \text{RSS} - \text{CSS} - \text{TrSS}$$

These results can be summarized in the form of analysis of variance table.

Calculation of SE, SE(d) and CD values

$$SE = \sqrt{\frac{EMS}{r}}$$

where r is the number of rows

$$SE(d) = \sqrt{2} \times SE$$

$$CD = t \times SE(d)$$

where t =table value of t for a specified level of significance and error degrees of freedom

Using CD value the bar chart can be drawn and the conclusion may be written.

Problem

Below are given the plan and yield in kgs/plot of a 5x5 Latin square experiment on the wheat crop carried out for testing the effects of five, manorial treatments A, B, C, D, and E. ‘A’ denotes control.

B	15	A	8	E	17	D	20	C	17	R1	=	77
A	9	D	21	C	19	E	16	B	13	R2	=	78
C	18	B	12	D	23	A	8	E	17	R3	=	78
E	18	C	16	A	10	B	15	D	23	R4	=	82
D	22	E	15	B	13	C	18	A	10	R5	=	78

$$C_1 = 82, C_2 = 72, C_3 = 82, C_4 = 77, C_5 = 80 ; GT = 393$$

Analyze the data and state your conclusions.

Analysis

1. Correction factor = $\frac{(GT)^2}{rxc}$ where GT is the grand total ‘r’ is number of rows, and ‘c’ is number of columns

$$= \frac{(393)^2}{5 \times 5} = 6177.96$$

$$2. \text{ Total SS} = 152 + 82 = \dots + 102 - CF$$

$$= 666_1 - 6177.96 = 483.04$$

$$3. \text{ SS due to rows (SSR)} = \frac{R_1^2 + R_2^2 + \dots + R_5^2}{t} - CF$$

$$= \frac{77^2 + \dots + 78^2}{5} - 6177.96 = 3.04$$

$$4. \text{ SS due to columns (SSC)} = \frac{C_1^2 + C_2^2 + \dots + C_5^2}{t} - CF$$

$$= \frac{82^2 + \dots + 80^2}{5} - 6177.96$$

$$= 14.24$$

5. To get SS due to treatments, first find the totals for each treatment using the given data as follows:

Treatment (A)	B	C	D	E
8	15	17	20	17
9	13	19	21	16
8	12	18	23	17
10	15	16	23	18
10	13	18	22	15
T ₁ =45	T ₂ = 68	T ₃ = 88	T ₄ = 109	T ₅ = 83

$$\therefore \text{ SS due to treatments} = \frac{T_1^2 + T_2^2 + \dots + T_5^2}{r} - CF$$

$$= \frac{45^2 + \dots + 83^2}{5} - 6177.96$$

$$= 454.64$$

$$6. \text{ SS due to error} = \text{TSS} - \text{SSR} - \text{SSC} - \text{SST}$$

$$= 483.04 - 3.04 - 14.24 - 454.6 = 11.12$$

7. Table for analysis of variance

Source of variation	Df	SS	MS	Variance ratio F	F value at 5% level & 1% level
Rows	4	3.04	0.76	123.34**	3.26 5.41
Columns	4	14.24	3.56		
Treatments	4	454.24	113.66		
Error	12	11.12	0.92		
Total	24	483.04			

** Highly significant

The observed highly significant value of the variance ratio indicates that there are significant differences between the treatment means.

S.E. of the difference between the treatment means (SED)

$$= \sqrt{\frac{2XEMS}{r}}$$

where EMS indicates the error mean square and 'r'

indicates the number of replications.

$$\text{i.e. SEd} = \sqrt{\frac{2 \times 0.92}{5}} = 0.61$$

∴ Critical difference = SEd x t 5% at df = 12 = 0.61 x 2.179

$$= 1.33$$

Summary of results

Treatment means will be calculated from the original table on treatment totals.

Treatments	A	B	C	D	E	CD 5%
Mean yield in kgs / plot	9.0	13.6	17.6	21.8	16.6	1.33

Conclusion represented symbolically

The treatment have been compared by setting them in the descending order of their yields.

Treatments : D C E B A

Mean yields 21.8 17.6 16.6 13.6 9.0
In kgs/plot

The treatment 'D' is the best of all. The treatments 'C' and 'E' do not differ significantly each other.

The yield obtained by applying every one of the manorial treatment is significantly higher that obtained without applying any manure.

Learning Exercise

1. An oil company tested four different blends of gasoline for fuel efficiency according to a Latin square design in order to control for the variability of four different drivers and four different models of cars. Fuel efficiency was measured in miles per gallon (mpg) after driving cars over a standard course.

Fuel Efficiencies (mpg) For 4 Blends of Gasoline
(Latin Square Design: Blends Indicated by Letters A-D)

		Car Model			
Driver		I	II	III	IV
1		D 15.5	B 33.9	C 13.2	A 29.1
2		B 16.3	C 26.6	A 19.4	D 22.8
3		C 10.8	A 31.1	D 17.1	B 30.3
4		A 14.7	D 34.0	B 19.7	C 21.6

These data are from Ott: *Statistical Methods and Data Analysis*, 4th ed., Duxbury, 1993, page 866.
(Similar data are given in the 5th edition by Ott/Longnecker, in problem 15.10, page 889.)

Analyse the data and draw your conclusion.

2. The numbers of wireworms counted in the plots of Latin square following soil fumigations (L,M,N,O,P) in the previous year were

		Columns				
Rows		P(4)	O(2)	N(5)	L(1)	M(3)
		M(5)	L(1)	O(6)	N(5)	P(3)
		O(4)	M(8)	L(1)	P(5)	N(4)
		N(12)	P(7)	M(7)	O(10)	L(5)
		L(5)	N(4)	P(3)	M(6)	O(9)

Analyse the data and draw your conclusions.

3. The following layout presents the observations made on 5 treatments A,B, C,D and E in an experiment of paddy crop by adopting LSD. The figures indicate the grain yield of paddy in kg/plot. Analyse the data and draw your conclusion.

		Columns				
Rows		B	D	E	A	C
		5	6	3	10	12
		C	A	B	E	D
	9	4	6	5	5	
	D	C	A	B	E	

Statistics

			8	15	7	6	5	
			E	B	C	D	A	
			5	8	13	9	5	
			A	E	D	C	B	
			9	6	12	16	8	



**This Book Download From e-course of ICAR
Visit for Other Agriculture books, News,
Recruitment, Information, and Events at
WWW.AGRIMOON.COM**

Give FeedBack & Suggestion at info@agrimoon.com

Send a Massage for daily Update of Agriculture on WhatsApp

+91-8148663744

DISCLAIMER:

The information on this website does not warrant or assume any legal liability or responsibility for the accuracy, completeness or usefulness of the courseware contents.

The contents are provided free for noncommercial purpose such as teaching, training, research, extension and self learning.



Connect With Us:

