



KARPAGAM ACADEMY OF HIGHER EDUCATION
(Deemed to be University, Established Under Section 3 of UGC Act 1956)
Coimbatore – 641 021.

Semester VI

16BCU602A

BIOSTATISTICS

Scope: On successful completion of this course the learner gains a clear knowledge about various aspects of Statistics, measures, hypothesis testing and application of them in their respective fields.

Objectives: To enable the students to understand the meaning, definition and functions of statistics through collection, representation, finding various measures such as mean, median, mode, correlation etc of statistics.

Unit 1

Definitions-Scope of Biostatistics- Variables in biology, collection, classification and tabulation of data- Graphical and diagrammatic representation.

Unit 2

Measures of central tendency – Arithmetic mean, median and mode. Measures of dispersion- Range, standard deviation, Coefficient of variation.

Unit 3

Correlation – Meaning and definition - Scatter diagram –Karl Pearson's correlation coefficient. Rank correlation.

Unit 4

Regression: Regression in two variables – Regression coefficient problems – uses of regression.

Unit 5

Test of significance: Tests based on Means only-Both Large sample and Small sample tests - Chi square test - goodness of fit.

TEXT BOOK

Pillai R.S.N., and Bagavathi V., 2002., Statistics, S. Chand & Company Ltd, New Delhi.

REFERENCES

Jerrold H.Z., (2003). Biostatistical Analysis, Fourth Edition, Pearson Education (Pte) .Ltd, New Delhi.

Arora, P.N., (1997). A foundation course statistics, S.Chand & Company Ltd, New Delhi.

Navnitham, P.A., (2004). Business Mathematics And Statistics, Jai Publications, Trichy,

Gupta S.P., (2001). Statistical methods, Sultan Chand & Sons, New Delhi



KARPAGAM ACADEMY OF HIGHER EDUCATION

(Deemed to be University Established Under Section 3 of UGC Act 1956)

Coimbatore – 641 021.

LECTURE PLAN DEPARTMENT OF BIOCHEMISTRY

STAFF NAME: Dr.K.Poornima

SUBJECT NAME: Biostatistics

SEMESTER:VI

SUB.CODE:16BCU602A

CLASS: III B.Sc (BC)

Sl. No	Duration of Period	Topics to be Covered	Support material
UNIT I			
1	1	Definitions-Scope of Biostatistics Variables in biology	T1:4-5 R3:1.13 R1:1-3 R2:1-2
2	1	Collection of primary data	T1:31-35 R3:3.4-3.14 R1:5-6 R2:27-31
3	1	Collection of Secondary data	T1:36-40 R3:3.2-3.4 R1:6-7 R2:31-33
4		Classification of data	T1:56-73 R3:5.2-5.9 R1:6-7 R2:44-51
5	1	Tabulation of data	T1:73-81 R3:5.18-5.35 R1:8-20 R2:52-57
6	1	Graphical representation of data	T1:101-118 R3:6.27-6.53 R2:60-78
7	1	Diagrammatic representation of data	T1:84-100 R3:6.2-6.27 R2:79-94
8	1	Revision and possible questions discussion of Unit I	
		Total no of hours planned for UNIT I = 07	

Unit II			
1	1	Measures of central tendency – Arithmetic mean problems	T1:121-146 R3:7.5-7.19 R1:21-28 R2:101-113
2	1	Median problems	T1:146-170 R3:7.19-7.32 R1:32-40 R2:113-117
3	1	Mode problem	T1:117-275 R3:7.33-7.43 R1:44-47 R2:117 -127
4	1	Measures of dispersion-Range	T1:234-239 R3:8.8 R1:64-65
5	1	Standard deviation	T1:249-304 R3:8.26-8.27 R1:75-85 R2:117-179
6	1	Practicing problems in standard deviation	T1:305-317 R3:8.27-8.32 R1:98-101
7	1	Coefficient of variation- Practicing problems	T1:270-275 R3:8.3-8.37 R1:85-88 R2:172-177
8	1	Revision and possible questions discussion of Unit II	
		Total no of hours planned for UNIT II = 08	
Unit III			
1	1	Correlation – Meaning and definition	T1:359-363 R3:10.2-10.3 R1:103-104
2	1	Scatter diagram	T1:363-365 R3:10.7-10.10
3	1	Karl Pearson's correlation coefficient	T1:365-385 R3:10.25-10.28 R1:112-122 R2:260-263
4	1	Continuation of Karl Pearson's correlation coefficient	T1:412-424 R3:10.44-10.61
5	1	Practicing problems in Karl Pearson's correlation coefficient	T1:385-391 R3:10.37.10.43
6	1	Rank correlation	T1:412-424

			R3:10.44-10.61 R1:125-130 R2:269-271
7	1	Working out of problems in Rank correlation	R3:10.44-10.61 R1:125-130 R2:269-271
8	1	Revision and possible questions discussion of Unit III	
		Total no of hours planned for UNIT III = 08	
Unit IV			
1	1	Regression - Introduction	T1:425-428 R3:11.2-11.4
2	1	Regression in two variables	T1:425-428 R3:11.2-11.4 R2:284-285
3	1	Regression equation of x on y	T1:431-445 R3:11.7 R1:149-170
4	1	Regression equation of y on x	T1:431-445 R3:11.9
5	1	Regression coefficient problems	T1:445-476 R3:11.10-11.46 R2:287-290
6	1	Uses of regression	T1:476 R3:11.47 R1:180
7	1	Revision and possible questions discussion of Unit IV	
		Total no of hours planned for UNIT IV = 07	
UNIT V			
1	1	Test of significance – An Introduction	T1:765-770 R3:3.14
2	1	Tests based on Means only – Large sample	T1:779-785 R3:3.14-3.21 R1:318-324
3	1	Small samples – students t- test	T1:785-787 R3:3.30-3.32 R1:324-327
4	1	Working out of problems in test of significance	T1:787-789 R3:3.31-3.47
5	1	Chi square test - goodness of fit.	T1:790-792 R3:4.2-4.7 R1:392-400
6	1	Practicing problems in chi square test	T1:792-805 R3:4.20-4.48
7	1	Revision and possible question discussion of unit V	
		Total no of hours planned for UNIT V = 07	

Total number of hours planned to complete this course – 40

REFERENCE**TEXT BOOK**

T1: Pillai R.S.N., and Bagavathi V., 2002., Statistics , S. Chand & Company Ltd, New Delhi.

REFERENCES

R1: Arora, P.N., (1997). A foundation course statistics, S.Chand & Company Ltd, New Delhi.

R2: Navnitham, P.A., (2004). Business Mathematics And Statistics, Jai Publications, Trichy,

R3:Gupta S.P., (2001). Statistical methods, Sultan Chand & Sons, New Delhi

UNIT-I SYLLABUS

Definitions-Scope of Biostatistics- Variables in biology, collection, classification and tabulation of data- Graphical and diagrammatic representation.

Introduction

Statistical tools are found useful in progressively increasing of disciplines. In ancient times the statistics or the data regarding the human force and wealth available in their land had been collected by the rulers. Now-a-days the fundamental concepts of statistics are considered by many to be essential part of their knowledge.

Origin and Growth

The origin of the word 'statistics' has been traced to the Latin word 'status', the Italian word 'statista', the French word 'statistique' and the German word 'statistik'. All these words mean political state.

Meaning

The word 'statistics' is used in two different meanings. As a plural word it means data or numerical statements. As a singular word it means the science of statistics and statistical methods. The word 'statistics' is also used currently as singular to mean data.

Definitions

Statistics is "the science of collection, organization, presentation, analysis and interpretation of numerical data". – Dr.S.P.Gupta.

"Statistics are numerical statement of facts in any department of enquiry, placed in relation to each other". – Dr.A.L.Bowley.

Functions

The following are the important functions of statistics.

- * Collection
- * Numerical Presentation
- * Diagrammatic Presentation
- * Condensation
- * Comparison
- * Forecasting
- * Policy Making
- * Effect Measuring
- * Estimation

- * Tests of significance.

Characteristics

- * Statistics is a Quantitative Science.
- * It never considers a single item.
- * The values should be different.
- * Inductive logic is applied.
- * Statistical results are true on the average.
- * Statistics is liable to be misused.

Collection of data

Data constitutes the base. The findings of an investigation depend on correctness and completeness of the relevant data. Sources of data are of two kinds- primary source and secondary source. The term source means origin or place from which data comes or got. A primary source is one that itself collects the data; a secondary source is one that makes available data which were collected by some other agency. Based on source, data are classified under two categories- Primary data and secondary data.

Primary data

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

Secondary Data

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

Methods of collection of primary Data

Primary Data is collected in any one of the following methods:

- * Direct personal interviews
- * Indirect oral interviews
- * Information from correspondence
- * Mailed questionnaire method.
- * Schedules sent through enumerators.

Sources of secondary data

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

Classification

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

Basis

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

i) Geographical Classification or Spatial Classification

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

ii) Chronological or Temporal or Historical Classification

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

iii) Qualitative Classification

Some data can be classified on the basis of attributes or characteristics.

iv) Quantitative Classification

Some data can be classified in terms of magnitudes.

Tabulation

Tabulation is the process of arranging data systematically in rows and columns of a table.

There are two methods or modes in which data can be presented. They are

- i) Statistical Tables
- ii) Diagrams or Graphs

Parts of a table

A good table has the following parts or components:

- * Identification number
- * Title
- * Prefatory Note or Head note
- * Stubs
- * Captions
- * Body of the table
- * Foot note
- * Source

Frequency Distribution

The easiest method of organizing data is a frequency distribution, which converts raw data into a meaningful pattern for statistical analysis.

The following are the *steps* of constructing a frequency distribution:

1. Specify the number of class intervals. A class is a group (category) of interest. No totally

accepted rule tells us how many intervals are to be used. Between 5 and 15 class intervals are generally recommended. Note that the classes must be both *mutually exclusive and all-inclusive*. Mutually exclusive means that classes must be selected such that an item can't fall into two classes, and all-inclusive classes are classes that together contain all the data.

2. When all intervals are to be the same width, the following rule may be used to find the required class interval width:

$$W = (L - S) / K$$

where:

W= class width, **L**= the largest data, **S**= the smallest data, **K**= number of classes

Example

Suppose the age of a sample of 10 students are: 20.9, 18.1, 18.5, 21.3, 19.4, 25.3, 22.0, 23.1, 23.9, and 22.5. We select $K=4$ and $W=(25.3 - 18.1)/4 = 1.8$ which is rounded-up to 2. The frequency table is as follows:

Class Interval	Class Frequency	Relative Frequency
18-20	3	30 %
20-22	2	20 %
22-24	4	40 %
24- 26	1	10 %

Cumulative Frequency Distribution

When the observations are numerical, cumulative frequency is used. It shows the total number of observations which lie above or below certain key values. Cumulative Frequency for a population = frequency of each class interval + frequencies of preceding intervals. For example, the cumulative frequency for the above problem is: 3, 5, 9, and 10.

Diagrams and Graphs

Diagrams

Diagrams are various geometrical shapes such as bars, circles etc . Diagrams are based on scales but are not confined to points or lines. They are more attractive and easier to understand than graphs and are widely used in advertisement and publicity.

Rules for construction

* Title

- * Proportion between width and height
- * Size
- * scale
- * Index
- * Suitable Diagram
- * Simplicity
- * Neatness
- * Foot-Note and source
- * Identification numbers.

Types of Diagram

The frequently used diagrams are divided into the following four heads:

1. One Dimensional diagram- Bar Diagram
2. Two Dimensional diagram – Pie Diagram, Rectangle, squares and circles
3. Three Dimensional diagram – Cubes
4. Pictograms and Cartograms.

Histograms are used to graph absolute, relative, and cumulative frequencies.

Ogive is also used to graph cumulative frequency. An ogive is constructed by placing a point corresponding to the *upper end of each class* at a height equal to the cumulative frequency of the class. These points then are connected. An ogive also shows the relative cumulative frequency distribution on the right side axis.

A **less-than ogive** shows how many items in the distribution have a value less than the upper limit of each class.

A **more-than ogive** shows how many items in the distribution have a value greater than or equal to the lower limit of each class.

A **less-than cumulative frequency polygon** is constructed by using the upper true limits and the cumulative frequencies.

A **more-than cumulative frequency polygon** is constructed by using the lower true limits and the cumulative frequencies.

Pie chart is often used in newspapers and magazines to depict budgets and other economic information. A complete circle (the pie) represents the total number of measurements. The size of a slice is proportional to the relative frequency of a particular category. For example, since a

complete circle is equal to 360 degrees, if the relative frequency for a category is 0.40, the slice assigned to that category is 40% of 360 or $(0.40)(360) = 144$ degrees.

POSSIBLE QUESTIONS

UNIT I

PART A (20 x 1 = 20 Marks)

Question number 1 – 20 online examinations

PART B (5 x 2 = 20 Marks)

1. Define Statistics.
2. Write the formula to calculate the angle in Pie Diagram.
3. Define Classification.
4. Define Mid-Value and also find the Mid-Value of 100 – 110.
5. Define Frequency Distribution.
6. What do you mean by size of the Class Interval?
7. Write a formula to calculate Percentage Bar Diagram.
8. Define Histogram.
9. What are the types of classification?
10. Write about Geographical Classification with example.
11. Define Frequency Distribution.
12. Write any two functions of Statistics.

PART C (5 X 6 = 30 Marks)

1. Explain about the Classification of data.
2. Draw a suitable Pie Diagram to represent the following submitted as a part of the budget proposal of the govt. of India for the year 1995 – 96.

Item of Expenditure	Percentage
i) Interest	25
ii) Defense	15
iii) Other non plan expenditure	20
iv) States share of taxes and duties	15
v) State and UT plan assistance	10

vi) Central plan	15
Total	100

3. You are given the average expenditure of a family for a month.

Item	Average Expenditure per month (Rs)
Food	2,400
Clothing	200
Rent	800
Medicare	150
Entertainment	450

Draw a Pie Diagram for the above data.

4. The following table shows the total sale in July 2005 of major brands of cars in India.

Represent the following data by using Simple Bar Diagram.

Brand	July 2005 Sales (in Rs. '000)
Maruti	81
Hyundai Santro	39
Honda city	14
Matiz	20
Opel Astra	10
Fait Uno	12
Ford	37
Mitsubishi Lancer	11
Mercedes	3

1	The most suitable form of questionnaire for publicity and Propaganda is	Diagram	Interval	Graph	Bar	Diagram
2	Which measure method can be adopted if the respondents are	Interval	Interval	Bar	Interval	Bar
3	Pie-chart represents the components of a factor by	Advantages	Advantages	Advantages	Advantages	Advantages
4	The number of questions of a questionnaire should be	2	2	2	2	2
5	Primary data are	Always more reliable compared to secondary data	Less reliable compared to secondary data	Depends on the case with which data have been collected	Depends on the agency collecting the data	Always more reliable compared to secondary data
6	The census data published for caste-wise population in India will be known as	Quota-wise classification	Two-way classification	Geographical classification	Quota-wise classification	Geographical classification
7	Statistics implies	Both data and science	Data only	Data, science and measures to samples	Statistics	Data, science and measures to samples
8	Data taken from the publication "Agricultural Statistics in India" will be considered as	Primary data	Secondary data	Primary and secondary data	Published data	Secondary data
9	Which classification of data of food production will be called as	Geographical classification	Chronological classification	Geographical classification	Chronological classification	Geographical classification
10	Who is the father of Biometrics?	R. A. Fisher	W. Gosset	St. Francis Galton	R. A. Fisher	St. Francis Galton
11	Number of classes of data is	2	2	2	2	2
12	Connected with primary data secondary data are	Less reliable	Less reliable	Less reliable	Less reliable	Less reliable
13	In quantitative classification data are classified on the basis of	Intervals	Intervals	Intervals	Intervals	Intervals
14	Classification according to class-intervals would yield	Discrete data	Discrete data	Discrete data	Discrete data	Discrete data
15	In qualitative classification data are classified on the basis of	Attributes	Attributes	Attributes	Attributes	Attributes
16	In non-quantitative classification data are classified on the basis of	Intervals	Intervals	Intervals	Intervals	Intervals
17	Which series is not a time series?	Time series	Time series	Time series	Time series	Time series
18	Individual observations are called	Individual observations	Individual observations	Individual observations	Individual observations	Individual observations
19	Which case is chronological classification?	Chronological	Chronological	Chronological	Chronological	Chronological
20	In discrete frequency distribution values are given as	Class intervals	Class intervals	Class intervals	Class intervals	Class intervals
21	In continuous frequency distribution values are given as	Class intervals	Class intervals	Class intervals	Class intervals	Class intervals
22	Which of the following is the one dimensional diagram?	Bar diagram	Bar diagram	Bar diagram	Bar diagram	Bar diagram
23	In bar diagram, the base line is	Vertical	Vertical	Vertical	Vertical	Vertical
24	Pictograms are drawn by	Data	Data	Data	Data	Data
25	Diagram is suitable for the data presented as	Continuous grouped frequency distribution	Continuous grouped frequency distribution	Continuous grouped frequency distribution	Continuous grouped frequency distribution	Continuous grouped frequency distribution
26	Non-quantitative data presented in diagrams form are called	Qualitative presentation	Qualitative presentation	Qualitative presentation	Qualitative presentation	Qualitative presentation
27	A simple table represents	Only one factor or variable	Only one factor or variable	Only one factor or variable	Only one factor or variable	Only one factor or variable
28	The core backbone of a table are known as	Table	Table	Table	Table	Table
29	The column of the most statistics has been named from the Latin word	Column	Column	Column	Column	Column
30	Diagram used to display	Qualitative	Qualitative	Qualitative	Qualitative	Qualitative
31	The column backbone of a table are known as	Table	Table	Table	Table	Table
32	Statistical data are collected by	Collective data without any measure	A given measure	A given measure	A given measure	A given measure
33	In grouped data, the number of classes preferred are	Minimum possible	Minimum possible	Minimum possible	Minimum possible	Minimum possible
34	Class interval is measured as	The sum of the upper and lower limit	Half of the sum of the upper and lower limit	The difference between the upper and lower limit	Upper limit - lower limit	The difference between the upper and lower limit
35	The shape of the histogram charts is that of a	Cube	Cube	Cube	Cube	Cube

UNIT-II SYLLABUS

Measures of central tendency – Arithmetic mean, median and mode. Measures of dispersion- Range, standard deviation, Coefficient of variation.

INTRODUCTION

In this chapter we are going to deal with Measures of central tendency and about the measures of dispersion. The measures of central tendency concentrate about the values in the central part of the distribution. Plainly speaking an average of a statistical series is the value of the variable which is the representative of the entire distribution. If we know the average alone we cannot form a complete idea about the distribution so for the completeness of the idea we use Measures of dispersion.

Measures of Central Tendency

According to Professor Bowley the measures of central tendency are “statistical constants which enable us to comprehend in a single effort the significance of the whole “

The following are the three measures of central tendency in this chapter we deal with

- Arithmetic Mean or simply Mean
- Median
- Mode

Arithmetic Mean or simply Mean

Arithmetic Mean or simply Mean is the total values of the item divided by their number of the items. It is usually denoted by \bar{X} .

Individual series

$$\bar{X} = \Sigma X / N$$

Example:

The expenditure of ten families are given below .Calculate arithmetic mean.

30 70 10 75 500 8 42 250 40 36

Solution

Here N=10

$$\Sigma X = 30 + 70 + 10 + 75 + 500 + 8 + 42 + 250 + 40 + 36 = 1061$$

—

$$\bar{X} = 1061 / 10 = 106.1$$

Discrete series

—

$$\bar{X} = \Sigma f X / \Sigma f$$

Example

Calculate the mean number of person per house.

No. of person : 2 3 4 5 6

No. of house : 10 25 30 25 10

Solution

X	f	f X
2	10	20
3	25	75
4	30	120
5	25	125
6	<u>10</u>	<u>60</u>

$$\Sigma f = 100 \quad \Sigma f X = 400$$

—

$$\bar{X} = 400 / 100 = 4$$

Continuous series

—

$$\bar{X} = \Sigma f m / \Sigma f \quad \text{where } m \text{ represents the mid value.}$$

$$\text{Mid-value} = (\text{upper boundary} + \text{lower boundary}) / 2.$$

Example

Calculate the mean for the following.

Marks : 20-30 30-40 40-50 50-60 60-70 70-80

No. of student : 5 8 12 15 6 4

Solution:

C.I	f	m	f m
20-30	5	25	125
30-40	8	35	280
40-50	12	45	540
50-60	15	55	825

60-70	6	65	390
70-80	<u>4</u>	75	<u>300</u>
$\Sigma f = 50$	$\Sigma f m = 2460$		

$$\bar{X} = 2460 / 50 = 49.2.$$

Median

The median is the value for the middle most items when all the items are in the order of magnitude. It is denoted by M or Me.

Individual series

For odd number of item

$$\text{Position of the median} = (N+1) / 2$$

For even number of item

$$\text{Position of the median} = [(N/2) + ((N/2)+1)] / 2$$

Example

Calculate median for the following.

22 10 6 7 12 8 5

Solution

Here N = 7

Arrange in ascending order or descending order.

5 6 7 8 10 12 22

$$\begin{aligned} (N+1) / 2 &= (7+1) / 2 \\ &= 4^{\text{th}} \text{ item} = 8 \end{aligned}$$

Discrete series

$$\text{Position of the median} = (N+1) / 2^{\text{th}} \text{ item.}$$

Example

Find the median for the following.

X : 10 15 17 18 21

F: 4 16 12 5 3

Solution

X	f	c.f
10	4	4

15	16	20
17	12	32
18	5	37
21	<u>3</u>	40

$$N = 40$$

$$\begin{aligned}(N+1)/2 &= (40+1)/2 = 20.5^{\text{th}} \text{ item} \\ &= (20^{\text{th}} \text{ item} + 21^{\text{st}} \text{ item})/2 = (15+17)/2 \\ &= 16.\end{aligned}$$

Continuous series

$$M = L + \frac{[(N/2) - c.f]}{f} \times i$$

Where L- lower boundary, f-frequency, i-size of class interval,
c.f- cumulative frequency.

Example

Calculate the median height given below.

Height	:	145-150	150-155	155-160	160-165	165-170	170-175
No. of student:		2	5	10	8	4	1

Solution :

Height	No. of student	c.f
145-150	2	2
150-155	5	7
<u>155-160</u>	<u>10</u>	17
160-165	8	25
165-170	4	29
170-175	<u>1</u>	30

$$\Sigma f = 30$$

$$\text{Position of the median} = N/2^{\text{th}} \text{ item} = 30/2 = 15.$$

$$M = L + \frac{[(N/2) - c.f]}{f} \times i$$

$$= 155 + \frac{[(15-7) \times 5]}{10} = 155 + (40/10) = 159.$$

Mode :

Mode is the value which has the greatest frequency density. Mode is usually denoted by Z.

Individual series

The value which occur more times are identified as mode.

Example

Determine the mode

32, 35, 42, 32, 42, 32.

Solution:

Unimode = 32.

Discrete series

Determine the mode

Size of dress No. of set

18	55
20	120
22	108
24	45

here mode represents highest frequency .

Mode = 20

Continuous series

$$Z = L + [i (f_1 - f_0) / (2f_1 - f_0 - f_2)]$$

Where L- lower boundary , f_1 -frequency of the modal class, f_0 – frequency of the preceding modal class, f_2 - frequency of the succeeding modal class, i-size of class interval , c.f- cumulative frequency.

Example

Determine the mode

Marks	:	0-10	10-20	20-30	30-40	40-50
No.of student	:	5	20	35	20	12

Solution

Marks	No. of student
0-10	5
10-20	20
20-30	35
30-40	20
40-50	12

$$Z = L + [i(f_1 - f_0) / (2f_1 - f_0 - f_2)]$$

$$= 20 + [10(35 - 20) / (2(35) - 20 - 20)] = 20 + 5$$

$$= 25.$$

Empirical relation

- Mode = 3 median - 2 mean.

Measures of Dispersion

Measure of dispersion deals mainly with the following three measures

- Range
- Standard deviation
- Coefficient of variation

Range

Range is the difference between the greatest and the smallest value.

- Range = L - S, where L-largest value & S-Smallest value
- Coefficient of range = (L-S) / (L+S)

Individual series

Example

Find the value of range and its coefficient of range for the following data.

8, 10, 5, 9, 12, 11

Solution

$$\text{Range} = L - S$$

$$= 12 - 5 = 7$$

$$\text{Coefficient of range} = (L - S) / (L + S)$$

$$= (12 - 5) / (12 + 5)$$

$$= 7 / 17 = 0.4118$$

Continuous series

Range = L - S, where L-Mid-value of largest boundary & S-Mid-value of smallest boundary

Calculate the range.

Marks	: 20-30	30-40	40-50	50-60	60-70	70-80
No.of student	: 5	8	12	15	6	4

Solution

C.I	f	m
20-30	5	25
30-40	8	35
40-50	12	45
50-60	15	55
60-70	6	65
70-80	4	75

Here $L=75$ & $S=25$

$$\text{Range} = L - S = 75 - 25 = 50$$

Standard deviation

The standard deviation is the root mean square deviation of the values from the arithmetic mean. It is a positive square root of variants. It is also called root mean square deviation. This is usually denoted by σ .

Individual series

$$\sigma = \sqrt{(\sum x^2 / N) - (\sum x / N)^2}$$

Example

Calculate standard deviation for the following data.

40,41,45,49,50,51,55,59,60,60.

Solution

X	X^2
40	1600
41	1681
45	2025
49	2401
50	2500
51	2601
55	3025
59	3481

$$60 \quad 3600$$

$$603600$$

$$510 \quad \Sigma x^2 = 26504$$

$$\sigma = \sqrt{(\Sigma x^2 / N) - (\Sigma x / N)^2}$$

$$= \sqrt{(26514/10) - (510/10)^2}$$

$$= 7.09$$

Discrete series

$$\sigma = \sqrt{(\Sigma fx^2 / \Sigma f) - (\Sigma fx / \Sigma f)^2}$$

Example

Calculate standard deviation for the following data.

X : 0 1 2 3 4 5

F : 1 2 4 3 0 2

Solution

X	f	fx	x ²	fx ²
0	1	0	0	0
1	2	2	1	2
2	4	8	4	16
3	3	9	9	27
4	0	0	16	0
5	<u>2</u>	25	<u>50</u>	
$\Sigma f = 12$	$\Sigma fx = 29$		$\Sigma fx^2 = 95$	

$$\sigma = \sqrt{(\Sigma fx^2 / \Sigma f) - (\Sigma fx / \Sigma f)^2}$$

$$= \sqrt{(95/12) - (29/12)^2}$$

$$= 1.44$$

Continuous series

$$\sigma = \sqrt{(\Sigma fm^2 / \Sigma f) - (\Sigma fm / \Sigma f)^2}$$

Example

C.I : 0-10 10-20 20-30 30-40 40-50

F : 2 5 9 3 1

Solution

C.I	f	m	fm	m ²	fm ²
0-10	2	5	10	25	50
10-20	5	15	75	225	1125
20-30	9	25	225	625	5625
30-40	3	35	105	1225	3675
40-50	<u>1</u>	45	<u>45</u>	2025	<u>2025</u>
	20		460		12500

$$\sigma = \sqrt{(\Sigma fm^2 / \Sigma f) - (\Sigma fm / \Sigma f)^2}$$

$$= \sqrt{(12500/20) - (460/20)^2}$$

$$= 9.79$$

Coefficient of variation

Coefficient of variation = [standard deviation / arithmetic mean] x100

Example

Calculate the coefficient of variation.

Mean= 51, standard deviation = 7.09

Solution

Coefficient of variation = [standard deviation / arithmetic mean] x100

$$= (7.09 / 51) \times 100$$

$$= 13.9$$

POSSIBLE QUESTIONS UNIT I

PART A (20 x 1 = 20 Marks) Question number 1 – 20 online examinations

PART B (5 x 2= 20Marks)

15. Calculate the Mean for the following.

X	20	30	35	15	10
f	2	3	4	3	2

16. Define Median and give Example.

17. Calculate the Range and its Coefficient for the following data.

X	:	12	14	16	18	20
f	:	1	3	5	3	1

18. What is mean by Bimodal?

19. Calculate the Median for the following data.

80	100	50	90	120	110
----	-----	----	----	-----	-----

20. Write the relation between Standard Deviation and Variance.

21. Calculate the Average number of students per class for the following data.

26	46	33	25	36	27	34	29
----	----	----	----	----	----	----	----

22. Find Median and Mode for the following data.

13	16	17	15	18	14	19	15	12
----	----	----	----	----	----	----	----	----

23. Define Range.

24. Find the Arithmetic Mean for the following data.

70	60	75	50	42	95	46
----	----	----	----	----	----	----

25. Calculate the Range and its Coefficient for the following data.

17	10	56	19	12	11	18	14
----	----	----	----	----	----	----	----

26. Find the median for 57, 58, 61, 42, 38, 65, 72, and 66

27. Write the empirical relation for Mode.

PART C (5 X 6 = 30 Marks)

1. Draw the less than Ogive and hence find the Median of the following data.

Marks	20 - 29	30 - 39	40 - 49	50 - 59	60 - 69	70 - 79	80 - 89	90 - 99
No. of students	7	11	24	32	9	14	2	1

2. Draw Percentage Bar Diagram for the following data.

Food	Rs.200
Clothing	Rs.60

Education	Rs.70
Rent	Rs.130
Miscellaneous	Rs.40

3. Draw a Multiple Bar Diagram for the following data.

Year	Sales (000 Rs)	Gross Profit (000 Rs.)	Net Profits (000 Rs)
1974	100	30	10
1975	120	40	15
1976	130	45	25
1977	150	50	25

4. Nixon Corporation manufactures computers. The following data are the numbers of computers produced at the company for sample of 25 days.

24	32	27	23	33	33	29	25	23	28	21
26	31	22	27	33	27	23	28	29	31	35
34	22	26								

Construct frequency distribution using classes 21 - 23, 24 - 26, 27 - 29, 30 - 32 and 33 - 35. And draw a Histogram to the frequency distribution.

5. The frequency distribution representing the number of days annually the employees at the Voltas Ltd. who were absent due to illness is

Number of days absent	0-2	3-5	6-8	9-11	12-14	Total
Frequency	5	12	20	10	3	50

Draw a Frequency Polygon to the above Frequency Distribution.

10. Calculate the Mode for the following Continuous Frequency Distribution.

Salary (in Rs. 1000s) :	0 - 19	20 - 39	40 - 59	60 - 79	80 - 99
No. of Employees:	5	20	35	20	12

11. Find the Mean and the Standard Deviation for the given below data set.

10	14	20	12	21	16	19	17	14	25
----	----	----	----	----	----	----	----	----	----

12. Calculate the Standard Deviation and Coefficient of Variance (CV) for the following data.

X	0 - 10	10 - 20	20 - 30	30 - 40	40 - 50
f	2	5	9	3	1

13. Calculate the Median for the following Continuous Frequency Distribution.

Wages (in Rs.) :	0 - 19	20 - 39	40 - 59	60 - 79	80 - 99
No. of Workers:	5	20	35	20	12

14. Calculate the Coefficient of Variation for the following data.

X	6	9	12	15	18
f	7	12	13	10	8

15. Calculate the Median for the following.

Hourly Wages (in Rs.)	40 - 50	50 - 60	60 - 70	70 - 80	80 - 90	90 - 100
Number of Employees	10	20	15	30	15	10

16. The following data give the details about salaries (in thousands of rupees) of seven employees randomly selected from a Pharmaceutical Company.

Serial No.	1	2	3	4	5	6	7
Salary per Annum ('000)	89	57	104	73	26	121	81

Calculate the Standard Deviation and Coefficient of variance of the given data.

17. Calculate the Arithmetic Mean for the following data.

Height (cms):	160	161	162	163	164	165	166
No. of Persons :	27	36	43	78	65	48	28

18. Calculate the Coefficient of Variance for the following data.

77 73 75 70 72 76 75 72 74 76

1	2	Mean is a measure of	central value	dispersion	correlation	regression	central value
2	2	Mean is that value in a frequency distribution which possesses maximum frequencies	maximum frequencies	maximum value	maximum frequency	max frequency	maximum frequencies
3	2	The most stable measure of central tendency is	the mean	the mode	the mode	the mean	the mean
4	2	Sum of the deviations about mean is	zero	zero	zero	zero	zero
5	2	The formula used to calculate arithmetic mean for individual series by direct method is	$\sum X/N$	$\sum fX/N$	$\sum fX/N$	$\sum fX/N$	$\sum fX/N$
6	2	The formula used to calculate arithmetic mean for individual series by short cut method is	$\sum X - Yd/N$	$\sum X - Yd/N$	$\sum X - Yd/N$	$\sum X - Yd/N$	$\sum X - Yd/N$
7	2	In continuous series the formula for A.M is	$(X_1 + X_2)/2$	$(X_1 + X_2)/2$	$(X_1 + X_2)/2$	$(X_1 + X_2)/2$	$(X_1 + X_2)/2$
8	2	A.M = 8, N = 12 then $\sum X =$	96	96	96	96	96
9	2	If 10, 20, 30, 40, 50 are the series, the mode is	40	40	40	40	40
10	2	The data series as 5, 15, 25, 35, 45 will be called as	discrete series	discrete series	discrete series	discrete series	discrete series
11	2	Which of the following divides the series into two equal parts	Mean	Mode	Median	Mean	Median
12	2	Mode of the following data is 5, 5, 5, 7, 8, 8, 9	5	7	no mode	5	no mode
13	2	Which of the following is not a measure of dispersion?	Range	Characteristic deviation	Standard deviation	Median	Median
14	2	Range of the series whose is given by	15-5	15-5	15-5	15-5	15-5
15	2	Which one of the following is relative measure of dispersion?	range	C.V	C.V	Coefficient of variation	Coefficient of variation
16	2	Range of a set of values is 40 and maximum value in the series is 85. The minimum value of the series is	45	45	45	45	45
17	2	If standard deviation is 5, then the variance is	25	25	25	25	25
18	2	If the value of mode and mean is 40 and 46 then the value of median is	44	46	44	44	44
19	2	Standard deviation is also called	Root mean square deviation	Mean square deviation	Root deviation	Root mean square deviation	Root mean square deviation
20	2	N = 10, 15 = 15-30 = 30-35 = 35-40	4	4	4	4	4
21	2	If 10, 15, 20, 25, 30 = 11.8, then	11.8	11.8	11.8	11.8	11.8
22	2	If the S.D and the C.V of a variable is 5 and 25, then the mean	200	200	200	200	200
23	2	Which one of the following is a measure of central tendency	median	median	median	median	median
24	2	Mean of the following values is = 5, 15, 20, 10, 25	15	15	15	15	15
25	2	The coefficient of the median from individual series is 2.0 and 1.0	1.0	1.0	1.0	1.0	1.0

UNIT-III SYLLABUS

Correlation: Meaning and definition - Scatter diagram –Karl Pearson's correlation coefficient. Rank correlation.

The term Correlation refers to the relationship between the variables. Simple correlation refers to the relationship between two variables. Various types of correlation are considered.

Positive or Negative when the values of two variables change in the same direction, their positive correlation between the two variables.

Example : X 50 60 70 95 100 105

Y 23 32 37 41 46 50

Example : X 34 25 18 10 7

Y 51 49 42 33 19

Simple or Partial or Multiple

When only two variables are considered as under positive or negative correlation above the correlation between them is called Simple correlation. When more than two variables as considered the correlation between two of them when all other variables are held constant, i.e., when the linear effects of all other variables on them are removed is called partial correlation. When more than two variables are considered the correlation between one of them and its estimate based on the group consisting of the other variables is called multiple correlation.

Methods

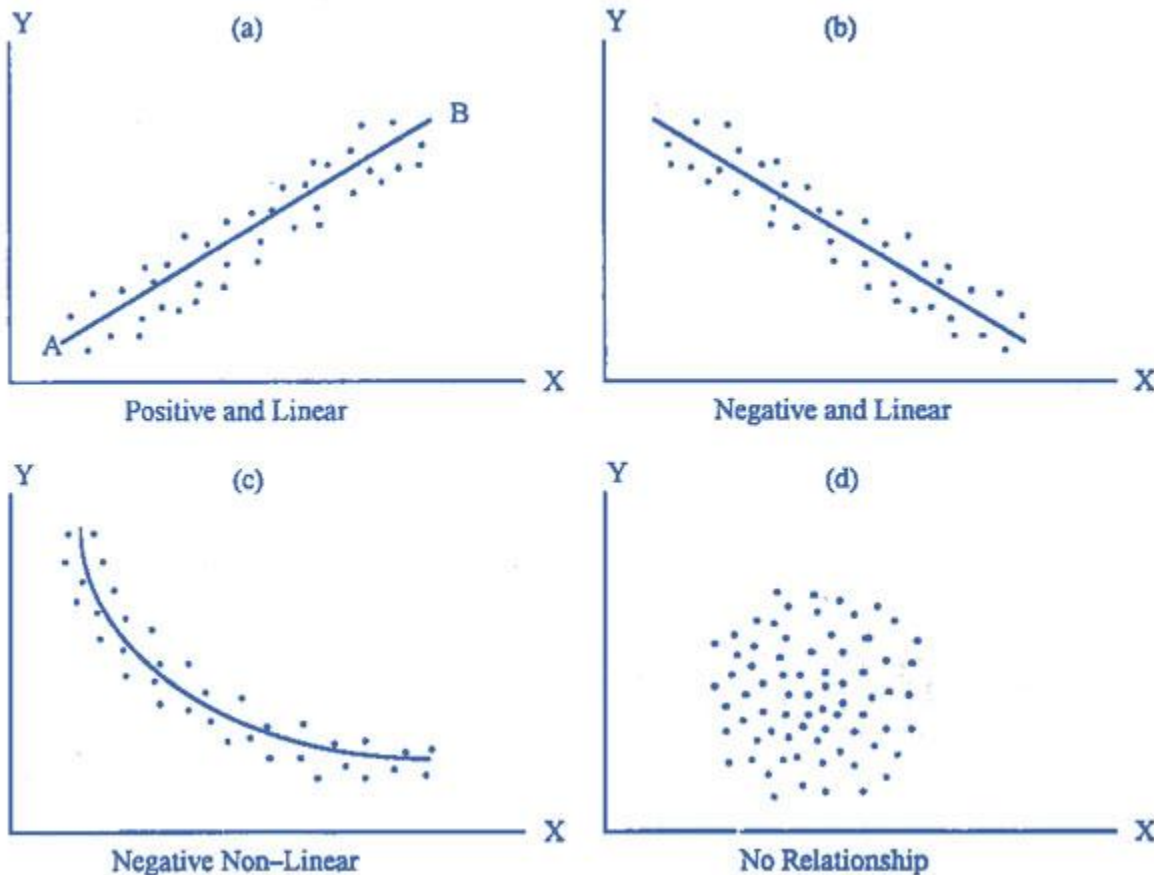
The following four methods are available under simple linear correlation and among them; product moment method is the best one.

- Scatter Diagram
- Karl Pearson's correlation coefficient or product moment correlation coefficient (r)
- Spearman's rank correlation coefficient (ρ)

- Correlation coefficient by concurrent deviation method (r_c).

Scatter Diagram

Scatter diagram is a graphic picture of the sample data. Suppose a random sample of n pairs of observations has the values $(X_1, Y_1), (X_2, Y_2), (X_3, Y_3), \dots, (X_n, Y_n)$. These points are plotted on a rectangular co-ordinate system taking independent variable on X -axis and the dependent variable on Y -axis. Whatever be the name of the independent variable, it is to be taken on X -axis. Suppose the plotted points are as shown in figure (a). Such a diagram is called scatter diagram. In this figure, we see that when X has a small value, Y is also small and when X takes a large value, Y also takes a large value. This is called direct or positive relationship between X and Y . The plotted points cluster around a straight line. It appears that if a straight line is drawn passing through the points, the line will be a good approximation for representing the original data. Suppose we draw a line AB to represent the scattered points. The line AB rises from left to the right and has positive slope. This line can be used to establish an approximate relation between the random variable Y and the independent variable X . It is nonmathematical method in the sense that different persons may draw different lines. This line is called the regression line obtained by inspection or judgment.



Making a scatter diagram and drawing a line or curve is the primary investigation to assess the type of relationship between the variables. The knowledge gained from the scatter diagram can be used for further analysis of the data. In most of the cases the diagrams are not as simple as in figure (a). There are quite complicated diagrams and it is difficult to choose a proper mathematical model for representing the original data. The scatter diagram gives an indication of the appropriate model which should be used for further analysis with the help of method of least squares. Figure (b) shows that the points in the scatter diagram are falling from the top left corner to the right. This is a relation called inverse or indirect. The points are in the neighborhood of a certain line called the regression line.

As long as the scattered points show closeness to a straight line of some direction, we draw a straight line to represent the sample data. But when the points do not lie around a straight line, we do not draw the regression line. Figure (c) shows that the plotted points have a tendency to fall from left to right in the form of a curve. This is a relation called non-linear or curvilinear. Figure (d) shows the points which apparently do not follow any pattern. If X takes a small value, Y may take a small or large value. There seems to be no sympathy between X and Y. Such a

diagram suggests that there is no relationship between the two variables.

Karl Pearson's Coefficient

Karl Pearson's Product-Moment Correlation Coefficient or simply Pearson's Correlation Coefficient for short, is one of the important methods used in Statistics to measure Correlation between two variables.

A few words about Karl Pearson. Karl Pearson was a British mathematician, statistician, lawyer and a eugenicist. He established the discipline of mathematical statistics. He founded the world's first statistics department in the University of London in the year 1911. He along with his colleagues Weldon and Galton founded the journal "Biometrika" whose object was the development of statistical theory.

The Correlation between two variables X and Y, which are measured using Pearson's Coefficient, give the values between +1 and -1. When measured in population the Pearson's Coefficient is designated the value of Greek letter rho (ρ). But, when studying a sample, it is designated the letter r. It is therefore sometimes called Pearson's r. Pearson's coefficient reflects the linear relationship between two variables. As mentioned above if the correlation coefficient is +1 then there is a perfect positive linear relationship between variables, and if it is -1 then there is a perfect negative linear relationship between the variables. And 0 denotes that there is no relationship between the two variables.

The degrees -1, +1 and 0 are theoretical results and are not generally found in normal circumstances. That means the results cannot be more than -1, +1. These are the upper and the lower limits.

Pearson's Coefficient computational formula

$$r = \frac{\sum XY - \frac{\sum X \sum Y}{N}}{\sqrt{(\sum X^2 - \frac{(\sum X)^2}{N})(\sum Y^2 - \frac{(\sum Y)^2}{N})}}$$

Sample question: compute the value of the correlation coefficient from the following table:

Subject	Age x	Weight Level y
1	43	99
2	21	65

3	25	79
4	42	75
5	57	87
6	59	81

Step 1: Make a chart. Use the given data, and add three more columns: xy , x^2 , and y^2 .

Subject	Age x	Weight Level y	xy	x^2	y^2
1	43	99			
2	21	65			
3	25	79			
4	42	75			
5	57	87			
6	59	81			

Step 2: Multiply x and y together to fill the xy column. For example, row 1 would be $43 \times 99 = 4,257$

Step 3: Take the square of the numbers in the x column, and put the result in the x^2 column.

Subject	Age x	Weight Level y	xy	x^2	y^2
1	43	99	4257	1849	
2	21	65	1365	441	
3	25	79	1975	625	
4	42	75	3150	1764	
5	57	87	4959	3249	
6	59	81	4779	3481	

Step 4: Take the square of the numbers in the y column, and put the result in the y^2 column.

Step 5: Add up all of the numbers in the columns and put the result at the bottom. The Greek letter sigma (Σ) is a short way of saying "sum of."

Subject	Age x	Weight Level y	xy	x^2	y^2
1	43	99	4257	1849	9801
2	21	65	1365	441	4225
3	25	79	1975	625	6241
4	42	75	3150	1764	5625
5	57	87	4959	3249	7569

6	59	81	4779	3481	6561
Σ	247	486	20485	11409	40022

Step 6: Use the following formula to work out the correlation coefficient. The answer is: 1.3787×10^{-4} the range of the correlation coefficient is from -1 to 1. Since our result is 1.3787×10^{-4} , a tiny positive amount, we can't draw any conclusions one way or another.

Spearman's Rank Correlation Coefficient

The Spearman correlation coefficient is often thought of as being the Pearson correlation coefficient between the ranked variables. In practice, however, a simpler procedure is normally used to calculate ρ . The n raw scores X_i, Y_i are converted to ranks x_i, y_i , and the differences $d_i = x_i - y_i$ between the ranks of each observation on the two variables are calculated.

If there are no tied ranks, then ρ is given by

$$\rho = 1 - \frac{6 \sum d^2}{N(N^2 - 1)}$$

If tied ranks exist, Pearson's correlation coefficients between ranks should be used for the calculation:

One has to assign the same rank to each of the equal values. It is an average of their positions in the ascending order of the values.

Example

X : 21 36 42 37 25

Y : 47 40 37 42 43. For the data given above, calculate the rank correlation coefficient.

Solution

		RANK			
X	Y	X	Y	d	D ²
21	47	5	1	4	16
36	40	3	4	-1	1
42	37	1	5	-4	16
37	42	2	3	-1	1
25	43	4	2	2	4

Total

$$\Sigma d = 0$$

$$\Sigma d^2 = 38$$

$$\rho = 1 - \frac{6 \Sigma d^2}{N(N^2 - 1)}$$

$$= 1 - \frac{6 \times 38}{5(5^2 - 1)}$$

$$= 1 - 1.9 = -0.9$$

Tied Ranks

When one or more values are repeated the two aspects- ranks of the repeated values and changes in the formula are to be considered.

Example

Find the rank correlation coefficient for the percentage of marks secured by a group of 8 students in Economics and Statistics.

Marks in Economics:	50	60	65	70	75	40	70	80
Marks in Statistics:	80	71	60	75	90	82	70	50

Solution

Let X be Marks in Economics and Y be Marks in Statistics

		RANK			
X	Y	X	Y	d	D ²
50	80	7	3	4	16
60	71	6	5	1	1
65	60	5	7	-2	4
70	75	3.5	4	-0.5	0.25
75	90	2	1	1	1
40	82	8	2	6	36
70	70	3.5	6	-2.5	6.25
80	50	1	8	-7	49
Total				$\Sigma d = 0$	$\Sigma d^2 = 113.5$

$$\rho = 1 - \frac{6 \{ \Sigma d^2 + m(m^2 - 1)/12 \}}{N(N^2 - 1)}$$

When $m=2$, $m(m^2-1)/12 = 0.5$

Therefore $\rho = 1 - \left[6\{113.5+0.5\}/8(8^2-1) \right]$

$$= 1 - 1.3571 = -0.3571$$

POSSIBLE QUESTIONS

UNIT III

PART A (20 x 1 = 20 Marks)

Question number 1 – 20 online examinations

PART B (5 x 2= 20Marks)

- 1) What are the types of Correlation?
- 2) Write any two properties of Correlation.
- 3) What is the range of Correlation Coefficient?
- 4) Define Positive Correlation.
- 5) What is meant by Regression?
- 6) What are the formulae for Regression co-efficients?
- 7) Distinguish between Correlation and Regression.
- 8) Write the formula for Rank Correlation, when more than one rank is repeated.
- 9) If $b_{xy} = -0.2337$ and $b_{yx} = -0.6643$ then find the Correlation Coefficient.
- 10) What is Negative Correlation? Give an example?
- 11) Write down the formula for Karl Pearson's Coefficient of Correlation.
- 12) Define Scatter Diagram.
- 13) What is Simple Correlation?

PART C (5 X 6 = 30 Marks)

- 1) Calculate the Correlation Coefficient from the following variables.

Sales in ('0000)	57	58	59	59	60	61	62	64
Advertisement Expenditure ('000)	17	16	15	18	12	14	19	11

- 2) Marks obtained by 8 students in Accountancy (X) and Statistics (Y) are given below. Compute Rank Correlation Coefficient.

X	25	20	28	22	40	60	20
---	----	----	----	----	----	----	----

Y	40	30	50	30	20	10	30
---	----	----	----	----	----	----	----

3) Calculate the two Regression Equations from the following data.

X	10	12	13	12	16	15
Y	40	38	43	45	37	43

4) Calculate Karl Pearson's Coefficient of Correlation from the following data.

Wages	100	101	102	102	100	99	97	98
Cost of Living	98	99	99	97	95	92	95	94

5) From the data given below find the two Regression Equations.

X	10	12	13	12	16	15
Y	20	28	23	25	27	30.

i) Estimate Y when X = 20.

ii) Estimate X when Y = 35.

1	2	Scatter diagram method is a	Graphic	Mathematical	Numerical	Alchemical	Graphic
2	3	If scores change in X, sometimes a corresponding decrease in Y, then X and Y are said to be	inversely related	inversely related	inversely correlated	inversely correlated	inversely correlated
3	4	Rank correlation was discovered by	R. A. Fisher	St. Francis Galton	Karl Pearson	Spearman	Spearman
4	5	Correlation is used to measure	strength of relationship between variables	one variable from another	shape of the distribution	central value	strength of relationship between variables
5	6	Coefficient of correlation lies between	-1 and +1	0 and 1	0 and 1	-1 and +1	-1 and +1
6	7	While drawing a scatter diagram if all points appear to form a straight line going downward from left to right, then it is inferred that there is	a perfect positive correlation	simple positive correlation	a perfect negative correlation	no correlation	a perfect negative correlation
7	8	The range of the rank correlation coefficient is	0 to 1	-1 to 1	0 to 1	-1 to 1	-1 to 1
8	9	The technique used to measure the strength of the relationship between two variables is called as	Ranking	Standard deviation	Ranking	Correlation analysis	Correlation analysis
9	10	Correlation is the	Relationship between two values	Relationship between two variables	Relationship between two variables	Relationship between two values	Relationship between two variables
10	11	If $r = +1$, the given two variables are	perfectly positive	perfectly negative	no correlation	negative	perfectly positive
11	12	If $r = -1$, the given two variables are	perfectly negative	perfectly positive	positive	positive	perfectly negative
12	13	If $r = 0$, the given two variables are	no correlation	perfect positive correlation	no correlation	no correlation	no correlation
13	14	Coefficient of correlation lies between	-1 and +1	0 and 1	0 and 1	-1 and +1	-1 and +1
14	15	The range of the rank correlation coefficient is	0 to 1	-1 to 1	0 to 1	-1 to 1	-1 to 1
15	16	Formula for rank correlation is	$\frac{2xy}{N(n+1)}$	$\frac{2xy}{N(n+1)}$	$\frac{2xy}{N(n+1)}$	$\frac{2xy}{N(n+1)}$	$\frac{2xy}{N(n+1)}$
16	17	The formula for computing Pearson's r is	$\frac{\sum xy}{N \sigma_x \sigma_y}$	$\frac{\sum xy}{N \sigma_x \sigma_y}$	$\frac{\sum xy}{N \sigma_x \sigma_y}$	$\frac{\sum xy}{N \sigma_x \sigma_y}$	$\frac{\sum xy}{N \sigma_x \sigma_y}$
17	18	In rank correlation the sum of the differences of ranks between two variables shall be	zero	more than 1	less than 1	0	more
18	19	Estimation of the value of one variable from the given value of another variable is done by	correlation analysis	regression analysis	regression analysis	regression analysis	regression analysis
19	20	If advertising and sales are correlated the required amount of expenditure for obtaining given amount of sales is calculated by	correlation analysis	regression analysis	regression analysis	regression analysis	regression analysis
20	21	If two ranks are equal at x^{th} place, the rank given to them is	$\frac{x}{2}$	$\frac{x}{2}$	$\frac{x}{2}$	$\frac{x}{2}$	$\frac{x}{2}$
21	22	When equal ranks are assigned the adjustments made by adding	$\frac{1}{2}(2m+1)$	$\frac{1}{2}(2m+1)$	$\frac{1}{2}(2m+1)$	$\frac{1}{2}(2m+1)$	$\frac{1}{2}(2m+1)$
22	23	This only method used with ranks not the actual value is	Karl Pearson's coefficient of correlation	rank correlation	scatter diagram method	standard deviation method	rank correlation
23	24	If the data are of a qualitative nature the homogeneity and influence on the method used is	Karl Pearson's coefficient of correlation	rank correlation	scatter diagram method	standard deviation method	rank correlation
24	25	If the data in scatter diagram is too scattered we can use that	$r = 1$	$r = 1$	$r = 0$	$r = 1$	$r = 0$

[illegible]

SYLLABUS

UNIT IV

Regression: Regression in two variables – Regression coefficient problems – uses of regression.

Simple Linear Regression

The line which gives the average relationship between the two variables is known as the regression equation. The regression equation is also called estimating equation.

Uses

1. Regression analysis is used in statistics and other disciplines.
2. Regression analysis is of practical use in determining demand curve, supply curve, consumption function, etc from market survey.
3. In Economics and Business, there are many groups of interrelated variables.
4. In social research, the relation between variables may not known; the relation may differ from place to place.
5. The value of dependent variable is estimated corresponding to any value of the independent variable using the appropriate regression equation.

Method of Least Squares

from a scatter diagram, there is virtually no limit as to the number of lines that can be drawn to make a linear relationship between the 2 variables

- the objective is to create a BEST FIT line to the data concerned
- the criterion is called the method of least squares
- i.e. the sum of squares of the *vertical deviations* from the points to the line be a minimum (based on the fact that the dependent variable is drawn on the vertical axis)
- the linear relationship between the dependent variable (Y) and the independent variable(x) can be written as $Y = a + bX$, where a and b are parameters describing the vertical intercept and the slope of the regression.
- Similarly the linear relationship between the dependent variable (XY) and the independent variable(Y) can be written as $X = a' + b'Y$, where a and b are parameters describing the vertical intercept and the slope of the regression.

○

Calculating a and b:

The values of a and b for the given pairs of values of (x_i, y_i) $i=1,2,3,\dots$ are determined,

Using the normal equations as ,

$$\sum y = Na + b\sum x$$

$$\sum xy = a\sum x + b\sum x^2$$

Similarly, the values of a' and b' for the given pairs of values of (x_i, y_i) $i=1,2,3,\dots$ are determined,

Using the normal equations as ,

$$\sum x = Na' + b'\sum y$$

$$\sum xy = a'\sum y + b'\sum y^2$$

Methods of forming the regression equations

- Regression equations on the basis of normal equations.
- Regression equations on the basis of X and Y and b_{YX} and b_{XY} .

Problem

From the following data, obtain the two regression equations.

X	6	2	10	4	8
Y	9	11	5	8	7

use normal equations.

Solution

X	Y	XY	X ²	Y ²
6	9	54	36	81
2	11	22	4	121
10	5	50	100	25
4	8	32	16	64
8	7	56	64	49
$\sum x=30$	$\sum y=40$	$\sum xy=214$	$\sum x^2=220$	$\sum y^2=340$

Let the regression equation Y on X is $Y = a + bX$

The normal equations are ,

$$\sum y = Na + b\sum x$$

$$\sum xy = a\sum x + b\sum x^2$$

By substituting the values from the table, we get

$$5a + 30b = 40 \text{ -----1}$$

$$30a + 180b = 214 \text{ -----2}$$

Solving these two equations we get,

$$a = 11.90 \text{ and } b = -0.65$$

Therefore the regression Y on X is $Y = 11.90 - 0.65X$.

Let the regression equation X on Y is $X = a' + b'Y$

The normal equations are,

$$\sum x = Na + b\sum y$$

$$\sum xy = a\sum y + b\sum y^2$$

By substituting the values from the table, we get

$$5a' + 40b' = 30 \text{ -----3}$$

$$40a' + 340b' = 214 \text{ -----4}$$

Solving these two equations we get,

$$a' = 16.40 \text{ and } b' = -1.30$$

Therefore the regression equation X on Y is $X = 16.40 - 1.30Y$

Example From the data given below, find

- the two regression equations
- The correlation coefficient between the variables X and y
- The value of Y when X = 30

X :	25	28	35	32	31	36	29	38	34	32
Y :	43	46	49	41	36	32	31	30	33	39

Solution

X	Y	$x = X - X'$	$Y = Y - Y'$	xy	x^2	y^2
25	43	-7	5	-35	49	25
28	46	-4	8	-32	16	64
35	49	3	11	33	9	121
32	41	0	3	0	0	9
31	36	-1	-2	2	1	4

36	32	4	-6	-24	16	36
29	31	-3	-7	21	9	49
38	30	6	-8	-48	36	64
34	33	2	-5	-10	4	25
32	39	0	1	0	0	1
320	380	0	0	-93	140	398

$$\bar{X} = 32, \bar{Y} = 38, b_{xy} = \sum xy / \sum y^2 = -0.2337, b_{yx} = \sum xy / \sum x^2 = -0.6643$$

iv) Regression equation of Y on X, $(Y - \bar{Y}) = b_{yx} (X - \bar{X})$

$$(Y - 38) = -0.6643(X - 32) \Rightarrow Y = 59.26 - 0.6643X$$

(ii) Regression equation of X on Y, $(X - \bar{X}) = b_{xy} (Y - \bar{Y})$

$$(X - 32) = -0.2337(Y - 38) \Rightarrow X = 40.88 - 0.2337Y$$

(iii) $r = +\sqrt{b_{yx} b_{xy}} = -0.3940$

(iv) $Y = 59.26 - 0.6643 \times 30 = 39$

Properties of Regression coefficients

1. The two regression equations are generally different and are not to be interchanged in their usage.
2. The two regression lines intersect at (\bar{X}, \bar{Y}) .
3. Correlation coefficient is the geometric mean of two regression coefficients.
4. The two regression coefficients and the correlation coefficient have the same sign.
5. Both the regression coefficients and the correlation coefficient cannot be greater than one numerically and simultaneously.
6. Regression coefficients are independent of change of origin but are affected by the change of scale.
7. Each regression coefficient is in the unit of the measurement of the dependent variable.
8. Each regression coefficient indicates the quantum of change in the dependent variable corresponding to unit increase in the independent variable.

POSSIBLE QUESTIONS

UNIT IV

PART A (20 x 1 = 20 Marks)

Question number 1 – 20 online examinations

PART B (5 x 2= 20Marks)

- 1) What is meant by Regression?
- 2) What are the formulae for Regression co-efficients?
- 3) Distinguish between Correlation and Regression.
- 4) What is Simple Correlation?
- 5) Define Regression Equation.
- 6) When $X = 40$, $Y = 60$, $\sigma_x = 10$, $\sigma_y = 15$ and $r = 0.7$ find the Regression Equation of Y on X.

PART C (5 X 6 = 30 Marks)

- 1) Calculate the two Regression Equations from the following data.

X	10	12	13	12	16	15
Y	40	38	43	45	37	43

- 2) Calculate Karl Pearson's Coefficient of Correlation from the following data.

Wages	100	101	102	102	100	99	97	98
Cost of Living	98	99	99	97	95	92	95	94

- 3) From the data given below find the two Regression Equations.

X	10	12	13	12	16	15
Y	20	28	23	25	27	30.

- i) Estimate Y when $X = 20$.
 - ii) Estimate X when $Y = 35$.
- 4) A comparison of the undergraduate Grade Point Averages of 10 corporate employees with their scores in a managerial trainee examination produced the results shown in the following table.

Exam Score	89	83	79	91	95	82	69	66	75	80
GPA	2.4	3.1	2.5	3.5	3.6	2.5	2.0	2.2	2.6	2.7

Measure the Correlation Coefficient between Exam scores and GPA by using Rank Method and also interpret the data given with the help of Scatter Diagram.

- 5) Develop the Regression Equation that best fit the data given below using annual income as an independent variable and amount of life insurance as dependent variable.

Annual Income (Rs. in 000's)	62	78	41	53	85	34
Amount of Life Insurance (Rs. in 00's)	25	30	10	15	50	7

- 6) The ranks of ten students in Economics and Statistics subjects are as follows.

Economics	3	5	8	4	7	10	2	1	6	9
Statistics	6	4	9	8	1	2	3	10	5	7

Calculate Spearman's Rank Correlation Coefficient.

- 7) You are given the following data:

	X	Y
Arithmetic Mean	36	85
Standard Deviation	11	8
Correlation coefficient between X and Y	= 0.66	

Find the two Regression Equations. And also find Correlation Coefficient.

UNIT IV					
1	The study of regression is considerably used by	Economics and Businessmen	Biological researcher	Mathematician	Engineers and Businessmen
2	The average relationship existing between X and Y variables is described by	Regression path	Regression charts	Regression line	Regression line
3	The regression equation of Y on X is expressed as	$Y = a + bX$	$Y = b_0 + b_1X$	$Y = a + bX$	$Y = a + bX$
4	The regression coefficient of X on Y is	$r = \frac{b_1}{b_2}$	$r = \frac{b_1}{b_2}$	$r = \frac{b_1}{b_2}$	$r = \frac{b_1}{b_2}$
5	The study used of the analysis of regression coefficients is used as	Mathematical	Mathematical	Mathematical	Mathematical
6	Scatter diagram method is a	Graphical	Graphical	Graphical	Graphical
7	If the regression line is horizontal, the regression coefficient is	Zero	Zero	Zero	Zero
8	Rank correlation was discovered by	R. A. Fisher	St. Francis Galton	Karl Pearson	Spearman
9	Correlation is used to measure	Strength of relationship between variables	Strength of relationship between variables	Strength of relationship between variables	Strength of relationship between variables
10	Coefficient of correlation lies between	+1 and -1	+1 and -1	+1 and -1	+1 and -1
11	While drawing a scatter diagram if all points appear to form a straight line trending downwards from left to right, then it is	A perfect positive correlation	A perfect positive correlation	A perfect positive correlation	A perfect positive correlation
12	The value of the rank correlation coefficient is	0 to 1	0 to 1	0 to 1	0 to 1
13	The technique used to measure the strength of the relationship between the variables is referred to as	Rank correlation	Rank correlation	Rank correlation	Rank correlation
14	If all the variables are varying in the same direction it is known as	Positive correlation	Positive correlation	Positive correlation	Positive correlation
15	X = 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90, 95, 100. Y = 10, 20, 30, 40, 50, 60, 70, 80, 90, 100. These two variables are	Positively correlated	Positively correlated	Positively correlated	Positively correlated
16	X = 10, 20, 30, 40, 50, 60, 70, 80, 90, 100. Y = 10, 20, 30, 40, 50, 60, 70, 80, 90, 100. These two variables are	Positively correlated	Positively correlated	Positively correlated	Positively correlated
17	If the two variables are varying in opposite direction the correlation is said to be	Negative	Negative	Negative	Negative
18	When we study the relationship between the yield of rice per acre and both the amount of rainfall and the amount of fertilizer used, it is the problem of	Multiple correlation	Multiple correlation	Multiple correlation	Multiple correlation
19	Study of three or more variables simultaneously is known as	Multiple correlation	Multiple correlation	Multiple correlation	Multiple correlation
20	If the amount of change in one variable tends to bear constant ratio to the amount of change in the other variable then the correlation is said to be	Linear	Linear	Linear	Linear
21	X = 10, 20, 30, 40, 50, 60, 70, 80, 90, 100. Y = 10, 20, 30, 40, 50, 60, 70, 80, 90, 100. These two variables are	Linear	Linear	Linear	Linear
22	If the amount of change in one variable does not bear a constant ratio to the amount of change in the other variable it is known as	Non linear	Non linear	Non linear	Non linear
23	In scatter diagram if the points are lying in a straight line from lower left hand corner to the upper right hand corner the correlation is said to be	Perfectly positive	Perfectly positive	Perfectly positive	Perfectly positive
24	In scatter diagram if the points are lying in a straight line from upper left hand corner to the lower right hand corner the correlation is said to be	Perfectly negative	Perfectly negative	Perfectly negative	Perfectly negative
25	A simple and non mathematical method of studying correlation between the variables is	Karl Pearson's coefficient of correlation	Karl Pearson's coefficient of correlation	Karl Pearson's coefficient of correlation	Karl Pearson's coefficient of correlation
26	Correlation is the	Relationship between two values	Relationship between two values	Relationship between two values	Relationship between two values
27	If $r = 1$, the given two variables are	Perfectly positive	Perfectly positive	Perfectly positive	Perfectly positive
28	If $r = -1$, the given two variables are	Perfectly negative	Perfectly negative	Perfectly negative	Perfectly negative
29	If $r = 0$, the given two variables are having	Zero correlation	Zero correlation	Zero correlation	Zero correlation
30	Coefficient of correlation lies between	+1 and -1	+1 and -1	+1 and -1	+1 and -1
31	The value of the rank correlation coefficient is	0 to 1	0 to 1	0 to 1	0 to 1
32	Formula for rank correlation is	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$
33	The formula for computing Pearson's r is	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$	$r = \frac{C - \frac{N(N+1)}{6}}{\frac{N(N^2-1)}{12}}$
34	In rank correlation the sum of the differences of ranks between two variables shall be	Zero	Zero	Zero	Zero
35	Translation of the value of one variable from the given value of another variable is done by	Regression analysis	Regression analysis	Regression analysis	Regression analysis
36	If advertising and sales are correlated the expected amount of sales for a given advertising expenditure is calculated by	Regression analysis	Regression analysis	Regression analysis	Regression analysis
37	If advertising and sales are correlated the required amount of expenditure for attaining given amount of sales is calculated by	Regression analysis	Regression analysis	Regression analysis	Regression analysis
38	If yield of rice and rainfall are correlated the amount of rain required to achieve a certain production figure is calculated by	Regression analysis	Regression analysis	Regression analysis	Regression analysis
39	The statistical measure of the linear relationship is	Correlation	Correlation	Correlation	Correlation
40	The term 'Regression' was first used by	R. A. Fisher	St. Francis Galton	Karl Pearson	Spearman
41	In 1877, the relationship between the height of the fathers and sons was studied by	R. A. Fisher	St. Francis Galton	Karl Pearson	Spearman
42	The regression equation of X on Y is expressed as	$X = a + bY$	$X = a + bY$	$X = a + bY$	$X = a + bY$
43	The measure of the average relationship between two or more variables in terms of their original units of the data is referred to as	Regression	Regression	Regression	Regression
44	One of the most famous and useful techniques in Economics and Business research is	Regression	Regression	Regression	Regression
45	The variable which is used to predict the variable of interest is called the	Independent variable	Independent variable	Independent variable	Independent variable
46	The variable which is used to predict the variable of interest is called the	Dependent variable	Dependent variable	Dependent variable	Dependent variable
47	In correlation analysis, b_1 and b_2 are	Intercept	Intercept	Intercept	Intercept
48	In regression analysis the regression coefficients b_1 and b_2 are	Intercept	Intercept	Intercept	Intercept
49	In the regression equation $Y = a + bX$, Y is a	Dependent variable	Dependent variable	Dependent variable	Dependent variable
50	In the regression equation $Y = a + bX$, X is a	Independent variable	Independent variable	Independent variable	Independent variable
51	The regression coefficient of Y on X is	b_1	b_1	b_1	b_1
52	In the regression equation $Y = a + bX$, b is the	Intercept	Intercept	Intercept	Intercept
53	In the regression equation $Y = a + bX$, a is the	Intercept	Intercept	Intercept	Intercept
54	In the regression equation $Y = a + bX$, a is the	Intercept	Intercept	Intercept	Intercept

UNIT-V SYLLABUS

Test of significance: Tests based on Means only-Both Large sample and Small sample tests – Student's t test, F-test, Chi square test - goodness of fit

Hypothesis:

A statistical hypothesis is an assumption that we make about a population parameter, which may or may not be true concerning one or more variables.

According to Prof. Morris Hamburg "A hypothesis in statistics is simply a quantitative statement about a population".

Hypothesis testing:

Hypothesis testing is to test some hypothesis about parent population from which the sample is drawn.

Example:

A coin may be tossed 200 times and we may get heads 80 times and tails 120 times, we may now be interested in testing the hypothesis that the coin is unbiased.

To take another example we may study the average weight of the 100 students of a particular college and may get the result as 110lb. We may now be interested in testing the hypothesis that the sample has been drawn from a population with average weight 115lb.

Hypotheses are two types

1. Null Hypothesis
2. Alternative hypothesis

Null hypothesis:

The hypothesis under verification is known as *null hypothesis* and is denoted by H_0 and is always set up for possible rejection under the assumption that it is true.

For example, if we want to find out whether extra coaching has benefited the students or not, we shall set up a null hypothesis that "*extra coaching has not benefited the students*". Similarly, if we want to find out whether a particular drug is effective in curing malaria we will take the null hypothesis that "*the drug is not effective in curing malaria*".

Alternative hypothesis:

The rival hypothesis or hypothesis which is likely to be accepted in the event of rejection of the null hypothesis H_0 is called alternative hypothesis and is denoted by H_1 or H_a .

For example, if a psychologist who wishes to test whether or not a certain class of people have a mean I.Q. 100, then the following null and alternative hypothesis can be established.

The null hypothesis would be

$$H_0 : \mu = 100$$

Then the alternative hypothesis could be any one of the statements.

$$H_1 : \mu \neq 100$$

$$(or) H_1 : \mu > 100$$

$$(or) H_1 : \mu < 100$$

Errors in testing of hypothesis:

After applying a test, a decision is taken about the acceptance or rejection of null hypothesis against an alternative hypothesis. The decisions may be four types.

- 1) The hypothesis is true but our test rejects it.(type-I error)
- 2) The hypothesis is false but our test accepts it. .(type-II error)
- 3) The hypothesis is true and our test accepts it.(correct)
- 4) The hypothesis is false and our test rejects it.(correct)

The first two decisions are called errors in testing of hypothesis.

i.e.1) Type-I error

2) Type-II error

1)Type-I error: The type-I error is said to be committed if the null hypothesis (H_0) is true but our test rejects it.

2)Type-II error: The type-II error is said to be committed if the null hypothesis (H_0) is false but our test accepts it.

Level of significance:

The maximum probability of committing type-I error is called level of significance and is denoted by α .

$$\alpha = P (\text{Committing Type-I error})$$

$$= P (H_0 \text{ is rejected when it is true})$$

This can be measured in terms of percentage i.e. 5%, 1%, 10% etc.....

Power of the test:

The probability of rejecting a false hypothesis is called power of the test and is denoted by $1 - \beta$.

$$\text{Power of the test} = P (H_0 \text{ is rejected when it is false})$$

$$= 1 - P (H_0 \text{ is accepted when it is false})$$

$$= 1 - P (\text{Committing Type-II error})$$

$$= 1 - \beta$$

- A test for which both α and β are small and kept at minimum level is considered desirable.
- The only way to reduce both α and β simultaneously is by increasing sample size.
- The type-II error is more dangerous than type-I error.

Critical region:

A statistic is used to test the hypothesis H_0 . The test statistic follows a known distribution. In a test, the area under the probability density curve is divided into two regions i.e. the region of acceptance and the region of rejection. The region of rejection is the region in which H_0 is rejected. It indicates that if the value of test statistic lies in this region, H_0 will be rejected. This region is called critical region. The area of the critical region is equal to the level of significance α . The critical region is always on the tail of the distribution curve. It may be on both sides or on one side depending upon the alternative hypothesis.

One tailed and two tailed tests:

A test with the null hypothesis $H_0: \theta = \theta_0$ against the alternative hypothesis $H_1: \theta \neq \theta_0$, it is called a two tailed test. In this case the critical region is located on both the tails of the distribution.

A test with the null hypothesis $H_0: \theta = \theta_0$ against the alternative hypothesis $H_1: \theta > \theta_0$ (right tailed alternative) or $H_1: \theta < \theta_0$ (left tailed alternative) is called one tailed test. In this case the critical region is located on one tail of the distribution.

$H_0: \theta = \theta_0$ against $H_1: \theta > \theta_0$ ----- right tailed test

$H_0: \theta = \theta_0$ against $H_1: \theta < \theta_0$ ----- left tailed test

Sampling distribution:

Suppose we have a population of size 'N' and we are interested to draw a sample of size 'n' from the population. In different time if we draw the sample of size n, we get different samples of different observations i.e. we can get ${}^N C_n$ possible samples. If we calculate some particular statistic from each of the ${}^N C_n$ samples, the distribution of sample statistic is called sampling distribution of the statistic. For example if we consider the mean as the statistic, then the distribution of all possible means of the samples is a distribution of the sample mean and it is called sampling distribution of the mean.

Standard error:

Standard deviation of the sampling distribution of the statistic t is called standard error of t.

$$\text{i.e. } S.E(t) = \sqrt{\text{Var}(t)}$$

Utility of standard error:

1. It is a useful instrument in the testing of hypothesis. If we are testing a hypothesis at 5% l.o.s and if the test statistic i.e. $|Z| = \left| \frac{t - E(t)}{S.E(t)} \right| > 1.96$ then the null hypothesis is rejected at 5% l.o.s otherwise it is accepted.
2. With the help of the S.E we can determine the limits within which the parameter value expected to lie.
3. S.E provides an idea about the precision of the sample. If S.E increases the precision decreases and vice-versa. The reciprocal of the S.E i.e. $\frac{1}{S.E}$ is a measure of precision of a sample.
4. It is used to determine the size of the sample.

Test statistic:

The test statistic is defined as the difference between the sample statistic value and the hypothetical value, divided by the standard error of the statistic.

$$\text{i.e. test statistic } Z = \frac{t - E(t)}{S.E(t)}$$

Procedure for testing of hypothesis:

1. Set up a null hypothesis i.e. $H_0 : \theta = \theta_0$.
2. Set up an alternative hypothesis i.e. $H_1 : \theta \neq \theta_0$ or $H_1 : \theta > \theta_0$ or $H_1 : \theta < \theta_0$
3. Choose the level of significance i.e. α .
4. Select appropriate test statistic Z .
5. Select a random sample and compute the test statistic.
6. Calculate the tabulated value of Z at $\alpha\%$ l.o.s i.e. Z_α .
7. Compare the test statistic value with the tabulated value at $\alpha\%$ l.o.s. and make a decision whether to accept or to reject the null hypothesis.

Large sample tests:

The sample size which is greater than or equal to 30 is called as large sample and the test depending on large sample is called large sample test.

The assumption made while dealing with the problems relating to large samples are

Assumption-1: The random sampling distribution of the statistic is approximately normal.

Assumption-2: Values given by the sample are sufficiently closed to the population value and can be used on its place for calculating the standard error of the statistic.

Large sample test for single mean (or) test for significance of single mean:

For this test

The null hypothesis is $H_0 : \mu = \mu_0$
against the two sided alternative $H_1 : \mu \neq \mu_0$
where μ is population mean
 μ_0 is the value of μ

Let $x_1, x_2, x_3, \dots, x_n$ be a random sample from a normal population with mean μ and variance σ^2

i.e. if $X \sim N(\mu, \sigma^2)$ then $\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$, Where \bar{x} be the sample mean

Now the test statistic $Z = \frac{t - E(t)}{S.E(t)} \sim N(0,1)$

$$= \frac{\bar{x} - E(\bar{x})}{S.E(\bar{x})} \sim N(0,1)$$

$$\Rightarrow Z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1)$$

Now calculate $|Z|$

Find out the tabulated value of Z at $\alpha\%$ l.o.s i.e. Z_α

If $|Z| > Z_\alpha$, reject the null hypothesis H_0

If $|Z| < Z_\alpha$, accept the null hypothesis H_0

Note: if the population standard deviation is unknown then we can use its estimate s, which will be calculated from the sample. $s = \sqrt{\frac{1}{n-1} \sum (x - \bar{x})^2}$.

Large sample test for difference between two means:

If two random samples of size n_1 and n_2 are drawn from two normal populations with means μ_1 and μ_2 , variances σ_1^2 and σ_2^2 respectively

Let \bar{x}_1 and \bar{x}_2 be the sample means for the first and second populations respectively

Then $\bar{x}_1 \sim N\left(\mu_1, \frac{\sigma_1^2}{n_1}\right)$ and $\bar{x}_2 \sim N\left(\mu_2, \frac{\sigma_2^2}{n_2}\right)$

Therefore $\bar{x}_1 - \bar{x}_2 \sim N\left(\mu_1 - \mu_2, \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)$

For this test

The null hypothesis is $H_0 : \mu_1 = \mu_2 \Rightarrow \mu_1 - \mu_2 = 0$

against the two sided alternative $H_1 : \mu_1 \neq \mu_2$

$$\begin{aligned} \text{Now the test statistic } Z &= \frac{t - E(t)}{S.E(t)} \sim N(0,1) \\ &= \frac{(\bar{x}_1 - \bar{x}_2) - E(\bar{x}_1 - \bar{x}_2)}{S.E(\bar{x}_1 - \bar{x}_2)} \sim N(0,1) \\ \Rightarrow Z &= \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{S.E(\bar{x}_1 - \bar{x}_2)} \sim N(0,1) \\ \Rightarrow Z &= \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}} \sim N(0,1) \quad [\text{since } \mu_1 - \mu_2 = 0 \text{ from } H_0] \end{aligned}$$

Now calculate $|Z|$

Find out the tabulated value of Z at $\alpha\%$ l.o.s i.e. Z_α

If $|Z| > Z_\alpha$, reject the null hypothesis H_0

If $|Z| < Z_\alpha$, accept the null hypothesis H_0

Note: If σ_1^2 and σ_2^2 are unknown then we can consider S_1^2 and S_2^2 as the estimate value of σ_1^2 and σ_2^2 respectively..

Large sample test for single standard deviation (or) test for significance of standard deviation:

Let $x_1, x_2, x_3, \dots, x_n$ be a random sample of size n drawn from a normal population with mean μ and variance σ^2 ,

for large sample, sample standard deviation s follows a normal distribution with mean σ and variance $\sigma^2/2n$ i.e. $s \sim N\left(\sigma, \sigma^2/2n\right)$

For this test

The null hypothesis is $H_0 : \sigma = \sigma_0$

against the two sided alternative $H_1 : \sigma \neq \sigma_0$

Now the test statistic $Z = \frac{t - E(t)}{S.E(t)} \sim N(0,1)$

$$= \frac{s - E(s)}{S.E(s)} \sim N(0,1)$$

$$\Rightarrow Z = \frac{s - \sigma}{\sigma / \sqrt{2n}} \sim N(0,1)$$

Now calculate $|Z|$

Find out the tabulated value of Z at $\alpha\%$ l.o.s i.e. Z_α

If $|Z| > Z_\alpha$, reject the null hypothesis H_0

If $|Z| < Z_\alpha$, accept the null hypothesis H_0

Large sample test for difference between two standard deviations:

If two random samples of size n_1 and n_2 are drawn from two normal populations with means μ_1 and μ_2 , variances σ_1^2 and σ_2^2 respectively

Let s_1 and s_2 be the sample standard deviations for the first and second populations respectively

Then $s_1 \sim N\left(\sigma_1, \frac{\sigma_1^2}{2n_1}\right)$ and $\bar{x}_2 \sim N\left(\sigma_2, \frac{\sigma_2^2}{2n_2}\right)$

Therefore $s_1 - s_2 \sim N\left(\sigma_1 - \sigma_2, \frac{\sigma_1^2}{2n_1} + \frac{\sigma_2^2}{2n_2}\right)$

For this test

The null hypothesis is $H_0 : \sigma_1 = \sigma_2 \Rightarrow \sigma_1 - \sigma_2 = 0$

against the two sided alternative $H_1 : \sigma_1 \neq \sigma_2$

Now the test statistic $Z = \frac{t - E(t)}{S.E(t)} \sim N(0,1)$

$$= \frac{(s_1 - s_2) - E(s_1 - s_2)}{S.E(s_1 - s_2)} \sim N(0,1)$$

$$\Rightarrow Z = \frac{(s_1 - s_2) - (\sigma_1 - \sigma_2)}{S.E(s_1 - s_2)} \sim N(0,1)$$

$$\Rightarrow Z = \frac{(s_1 - s_2)}{\sqrt{\frac{\sigma_1^2}{2n_1} + \frac{\sigma_2^2}{2n_2}}} \sim N(0,1) \quad [\text{since } \sigma_1 - \sigma_2 = 0 \text{ from } H_0]$$

Now calculate $|Z|$

Find out the tabulated value of Z at $\alpha\%$ l.o.s i.e. Z_α

If $|Z| > Z_\alpha$, reject the null hypothesis H_0

If $|Z| < Z_\alpha$, accept the null hypothesis H_0

Large sample test for single proportion (or) test for significance of proportion:

Let x is number of success in n independent trials with constant probability p , then x follows a binomial distribution with mean np and variance npq .

In a sample of size n let x be the number of persons possessing a given attribute

then the sample proportion is given by $\hat{p} = \frac{x}{n}$

$$\text{Then } E(\hat{p}) = E\left(\frac{x}{n}\right) = \frac{1}{n}E(x) = \frac{1}{n}np = p$$

$$\text{And } V(\hat{p}) = V\left(\frac{x}{n}\right) = \frac{1}{n^2}V(x) = \frac{1}{n^2}npq = \frac{pq}{n}$$

$$S.E(\hat{p}) = \sqrt{\frac{pq}{n}}$$

For this test

The null hypothesis is $H_0 : p = p_0$
against the two sided alternative $H_1 : p \neq p_0$

$$\text{Now the test statistic } Z = \frac{t - E(t)}{S.E(t)} \sim N(0,1)$$

$$= \frac{\hat{p} - E(\hat{p})}{S.E(\hat{p})} \sim N(0,1)$$

$$\Rightarrow Z = \frac{\hat{p} - p}{\sqrt{\frac{pq}{n}}} \sim N(0,1)$$

Now calculate $|Z|$

Find out the tabulated value of Z at $\alpha\%$ l.o.s i.e. Z_α

If $|Z| > Z_\alpha$, reject the null hypothesis H_0

If $|Z| < Z_\alpha$, accept the null hypothesis H_0

Large sample test for single proportion (or) test for significance of proportion:

let x_1 and x_2 be the number of persons processing a given attribute in a random sample of size n_1 and n_2 then the sample proportions are given by $\hat{p}_1 = \frac{x_1}{n_1}$ and

$$\hat{p}_2 = \frac{x_2}{n_2}$$

Then $E(\hat{p}_1) = p_1$ and $E(\hat{p}_2) = p_2 \Rightarrow E(\hat{p}_1 - \hat{p}_2) = p_1 - p_2$

And $V(\hat{p}_1) = \frac{p_1 q_1}{n_1}$ and $V(\hat{p}_2) = \frac{p_2 q_2}{n_2} \Rightarrow V(\hat{p}_1 - \hat{p}_2) = \frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}$

$$S.E(\hat{p}_1) = \sqrt{\frac{p_1 q_1}{n_1}} \text{ and } S.E(\hat{p}_2) = \sqrt{\frac{p_2 q_2}{n_2}} \Rightarrow S.E(\hat{p}_1 - \hat{p}_2) = \sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}$$

For this test

The null hypothesis is $H_0 : p_1 = p_2$
against the two sided alternative $H_1 : p_1 \neq p_2$

Now the test statistic $Z = \frac{t - E(t)}{S.E(t)} \sim N(0,1)$

$$= \frac{\hat{p}_1 - \hat{p}_2 - E(\hat{p}_1 - \hat{p}_2)}{S.E(\hat{p}_1 - \hat{p}_2)} \sim N(0,1)$$

$$\Rightarrow Z = \frac{\hat{p}_1 - \hat{p}_2 - (p_1 - p_2)}{S.E(\hat{p}_1 - \hat{p}_2)} \sim N(0,1)$$

$$\Rightarrow Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{p_1 q_1}{n_1} + \frac{p_2 q_2}{n_2}}} \sim N(0,1)$$

$$\Rightarrow Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{pq}{n_1} + \frac{pq}{n_2}}} \sim N(0,1) \quad \text{Since } p_1 = p_2 \text{ from } H_0$$

$$\Rightarrow Z = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{pq \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}} \sim N(0,1)$$

When p is not known p can be calculated by $p = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}$ and $q = 1 - p$

Now calculate $|Z|$

Find out the tabulated value of Z at $\alpha\%$ l.o.s i.e. Z_α

If $|Z| > Z_\alpha$, reject the null hypothesis H_0

If $|Z| < Z_\alpha$, accept the null hypothesis H_0

- As σ is unknown,

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}} = \left[\bar{x} - z_{\alpha/2} \frac{s}{\sqrt{n}}, \bar{x} + z_{\alpha/2} \frac{s}{\sqrt{n}} \right]$$

Step 2: If μ_0 falls into the above confidence intervals, then

do *not* reject H_0 . Otherwise, reject H_0 .

Example 1:

The average starting salary of a college graduate is \$19000 according to government's report. The average salary of a random sample of 100 graduates is \$18800. The standard error is 800.

- Is the government's report reliable as the level of significance is 0.05.
- Find the p-value and test the hypothesis in (a) with the level of significance $\alpha = 0.01$.
- The other report by some institute indicates that the average salary is \$18900. Construct a 95% confidence interval and test if this report is reliable.

[solutions:]

(a)

$$H_0 : \mu = \mu_0 = 19000 \text{ vs. } H_a : \mu \neq \mu_0 = 19000,$$

$$n = 100, \bar{x} = 18800, s = 800, \alpha = 0.05$$

Then,

$$|z| = \left| \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \right| = \left| \frac{18800 - 19000}{800/\sqrt{100}} \right| = |-2.5| = 2.5 > z_{\alpha/2} = z_{0.025} = 1.96.$$

Therefore, reject H_0 .

(b)

$$\text{p-value} = P(|Z| > |z|) = P(|Z| > 2.5) = 2 \cdot P(Z > 2.5) = 0.0124 > 0.01$$

Therefore, *not* reject H_0 .

(c)

$$H_0 : \mu = \mu_0 = 18900 \text{ vs } H_a : \mu \neq \mu_0 = 18900,$$

A 95% confidence interval is

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}} = 18800 \pm z_{0.025} \frac{800}{\sqrt{100}} = 18800 \pm 1.96 \cdot 80 = [18643.2, 18956.8].$$

Since $\mu_0 = 18900 \in [18643.2, 18956.8]$, Therefore, *not* reject H_0 .

Example 2:

A sample of 49 provides a sample mean of 38 and a sample standard deviation of 7. Let

$\alpha = 0.05$. Please test the hypothesis

$$H_0 : u = 40 \text{ vs. } H_a : u \neq 40 .$$

based on

- (a) classical hypothesis test
- (b) p-value
- (c) confidence interval.

[solution:]

$$\bar{x} = 38, s = 7, u_0 = 40, n = 49, z = \frac{\bar{x} - u_0}{s/\sqrt{n}} = \frac{38 - 40}{7/\sqrt{49}} = -2 .$$

(a)

$$|z| = 2 > 1.96 = z_{0.025}$$

we reject H_0 .

(b)

$$p\text{-value} = P(|Z| > |z|) = P(|Z| > 2) = 2 * (1 - 0.9772) = 0.0456 < 0.05 = \alpha$$

we reject H_0 .

(c)

$100 \times (1 - \alpha)\% = 95\%$ confidence interval is

$$\bar{x} \pm z_{\alpha/2} \frac{s}{\sqrt{n}} = 38 \pm z_{0.025} \frac{7}{\sqrt{49}} = 38 \pm 1.96 = [36.04, 39.96] .$$

Since $40 \notin [36.04, 39.96]$, we reject H_0 .

Hypothesis Testing for the Mean (Small Samples)

For samples of size less than 30 and when σ is unknown, if the population has a normal, or nearly normal, distribution, the t -distribution is used to test for the mean μ .

Using the t-Test for a Mean μ when the sample is small		
Procedure	Equations	Example 4
State the claim mathematically and verbally. Identify the null and alternative hypotheses	State H_0 and H_a	$H_0 : \mu \geq 16500$ $H_a : \mu < 16500$ $n = 14, \bar{x} = 15700, s = 1250$
Specify the level of	Specify α	$\alpha = 0.05$

significance		
Identify the degrees of freedom and sketch the sampling distribution	$d.f. = n - 1$	$d.f. = 13$
Determine any critical values. If test is left tailed, use One tail, α column with a negative sign. If test is right tailed, use One tail, α column with a positive sign. If test is two tailed, use Two tails, α column with a negative and positive sign.	Table 5 (<i>t</i> -distribution) in appendix B	The test is left-tailed. Since test is left tailed and $d.f. = 13$, the critical value is $t_0 = -1.771$
Determine the rejection regions.	The rejection region is $t < t_0$	The rejection region is $t < -1.771$
Find the standardized test statistic	$t = \frac{\bar{x} - \mu}{\sigma_{\bar{x}}} \approx \frac{\bar{x} - \mu}{s/\sqrt{n}}$	$t = \frac{15700 - 16500}{1250/\sqrt{14}} \approx -2.39$
Make a decision to reject or fail to reject the null hypothesis	If t is in the rejection region, reject H_0 , Otherwise do not reject H_0	Since $-2.39 < -1.771$, reject H_0
Interpret the decision in the context of the original claim.		Reject claim that mean is at least 16500.

Chi-Square Tests and the F-Distribution

Goodness of Fit

DEFINITION A **chi-square goodness-of-fit test** is used to test whether a frequency distribution fits an expected distribution.

The test is used in a multinomial experiment to determine whether the number of results in each category fits the null hypothesis:

H_0 : The distribution fits the proposed proportions

H_1 : The distribution differs from the claimed distribution.

To calculate the test statistic for the chi-square goodness-of-fit test, you can use observed frequencies and expected frequencies.

DEFINITION The **observed frequency O** of a category is the frequency for the category observed in the sample data.

The **expected frequency E** of a category is the calculated frequency for the category. Expected frequencies are obtained assuming the specified (or hypothesized) distribution. The expected frequency for the i th category is

$$E_i = np_i$$

where n is the number of trials (the sample size) and p_i is the assumed probability of the i th category.

The Chi-square Goodness of Fit Test: The sampling distribution for the goodness-of-fit test is a chi-square distribution with $k - 1$ degrees of freedom where k is the number of categories. The test statistic is

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

where O represents the observed frequency of each category and E represents the expected frequency of each category. To use the chi-square goodness of fit test, *the following must be true*.

1. The observed frequencies must be obtained using a random sample.
2. The expected frequencies must be ≥ 5 .

Performing the Chi-Square Goodness-of-Fit Test (p 496)		
Procedure	Equations	Example (p 497)
Identify the claim. State the null and alternative hypothesis.	State H_0 and H_1	H_0 : Classical 4% Country 36% Gospel 11% Oldies 2% Pop 18% Rock 29%
Specify the significance level	Specify α	$\alpha = 0.01$
Determine the degrees of freedom	d.f. = #categories - 1	$d.f. = 6 - 1 = 5$
Find the critical value	χ^2_α : Obtain from Table 6 Appendix B	$\phi^2_{0.01}(d.f = 5) = 15.086$

Identify the rejection region	$\chi^2 \geq \chi^2_{\alpha}$	$\chi^2 \geq 15.086$
Calculate the test statistic	$\chi^2 = \sum \frac{(O - E)^2}{E}$	<p>Survey results, n = 500</p> <p>Classical O = 8 E = .04*500 = 20</p> <p>Country O = 210 E = .36*500 = 180</p> <p>Gospel O = 7 E = .11*500 = 55</p> <p>Oldies O = 10 E = .02*500 = 10</p> <p>Pop O = 75 E = .18*500 = 90</p> <p>Rock O = 125 E = .29*500 = 145</p> <p>Substituting $\chi^2 = 22.713$</p>
Make the decision to reject or fail to reject the null hypothesis	<p>Reject if χ^2 is in the rejection region</p> <p>Equivalently, we reject if the P-value (the probability of getting as extreme a value or more extreme) is $\leq \alpha$</p>	<p>Since $22.713 > 15.086$ we reject the null hypothesis</p> <p>Equivalently</p> <p>$P(X \geq 22.713) < 0.01$ so reject the null hypothesis. (Note Table 6 of Appendix B doesn't have a value less than 0.005.)</p>
Interpret the decision in the context of the original claim		Music preferences differ from the radio station's claim.

Using Minitab to perform the Chi-Square Goodness-of-Fit Test (Manual p 237)

The data from the example above (Example 2 p 497) will be used.

Enter Three columns: Music Type: Classical, etc, Observed: 8 etc, Distribution 0.04, etc. (Note the names of the columns 'Music Types', 'Observed' and 'Distribution' are entered in the gray row at the top.)

Select **Calc->Calculator**, **Store the results in C4**, and calculate the **Expression** C3*500, click **OK**, Name C4 'Expected' since it now contains the expected frequencies

Music Type	Observed	Distribution	Expected
Classical	8	0.04	20
Country	210	0.36	180
Gospel	72	0.11	55
Oldies	10	0.02	10

Pop	75	0.18	90
Rock	125	0.29	145

Next calculate the chi-square statistic, $(O-E)^2/E$ as follows: Click **Calc->Calculator**. Store the results in C5 and calculate the Expression $(C2-C4)**2/C4$. Click on **OK** and C5 should contain the calculated values.

7.2000
5.0000
41.8909
0.0000
2.5000
2.7586

Next add up the values in C5 and the sum is the test statistic as follows: Click on **Calc->Column Statistics**. Select **Sum** and enter C5 for the **Input Variable**. Click **OK**. The chi-square statistic is displayed in the session window as follows:

Sum of C5
Sum of C5 = 22.7132

Next calculate the P-value: Click on **Calc->Probability Distributions->Chi-square**. Select **Cumulative Probability** and enter 5 **Degrees of Freedom**. Enter the value of the test statistic 22.7132 for the **Input Constant**. Click **OK**.

The following is displayed on the Session Window.

Cumulative Distribution Function
Chi-Square with 5 DF
x P(X <= x)
22.7132 0.999617

$P(X \leq 22.7132) = 0.999617$ So the P-value = $1 - 0.999617 = 0.000383$. This is less than $\alpha = 0.01$ so we reject the null hypothesis.

Instead of calculating the P-value, we could have found the critical value from the Chi-Square

table (Table 6 Appendix B) for 5 degrees of freedom as we did above. The value is 15.086, and since our test statistic is 22.7132, we reject the null hypothesis.

Chi-Square with M&M's

H_0 : Brown: 13%, Yellow: 14%, Red: 13%, Orange: 20%, Green 16%, Blue 24%
Significance level: $\alpha = 0.05$
Degrees of freedom: number of categories – 1 = 5
Critical Value: $\chi^2_{0.05}(d.f. = 5) = 11.071$
Rejection Region: $\chi^2 \geq 11.071$
Test Statistic: $\chi^2 = \sum \frac{(O - E)^2}{E}$, where O is the actual number of M&M's of each color in the bag and E is the proportions specified under H_0 times the total number.
Reject H_0 if the test statistic is greater than the critical value (1.145)

Section 10.2 Independence

This section describes the chi-square test for independence which tests whether two random variables are independent of each other.

DEFINITION An $r \times c$ **contingency table** shows the observed frequencies for the two variables. The observed frequencies are arranged in r rows and c columns. The intersection of a row and a column is called a **cell**.

The following is a contingency table for two variables A and B where f_{ij} is the frequency that A equals A_i and B equals B_j .

	A_1	A_2	A_3	A_4	A
B_1	f_{11}	f_{12}	f_{13}	f_{14}	$f_{1.}$
B_2	f_{21}	f_{22}	f_{23}	f_{24}	$f_{2.}$
B_3	f_{31}	f_{32}	f_{33}	f_{34}	$f_{3.}$
B	$f_{.1}$	$f_{.2}$	$f_{.3}$	$f_{.4}$	f

If A and B are independent, we'd expect

$$f_{ij} = \text{prob}(A = A_i) * \text{prob}(B = B_j) * f = \left(\frac{f_{i.}}{f} \right) \left(\frac{f_{.j}}{f} \right) f = \frac{(f_{i.})(f_{.j})}{f}$$

$$\frac{(\text{sum of row } i) * (\text{sum of column } j)}{\text{sample size}}$$

Example 1 Determining the expected frequencies of CEO's ages as a function of company size under the assumption that age is independent of company size.

	<= 39	40 - 49	50 - 59	60 - 69	>= 70	Total
Small/midsize	42	69	108	60	21	300
Large	5	18	85	120	22	250
Total	47	87	193	180	43	550

	<= 39	40 - 49	50 - 59	60 - 69	>= 70	Total
Small/midsize	$\frac{300 * 47}{550}$ ≈ 25.64	$\frac{300 * 87}{550}$ ≈ 47.45	$\frac{300 * 193}{550}$ ≈ 105.27	$\frac{300 * 180}{550}$ ≈ 98.18	$\frac{300 * 43}{550}$ ≈ 23.45	300
Large	$\frac{250 * 47}{550}$ ≈ 21.36	$\frac{250 * 87}{550}$ ≈ 39.55	$\frac{250 * 193}{550}$ ≈ 87.73	$\frac{250 * 180}{550}$ ≈ 81.82	$\frac{250 * 43}{550}$ ≈ 19.55	250
Total	47	87	193	180	43	550

After finding the expected frequencies under the assumption that the variables are independent, you can test whether they are independent using the chi-square independence test.

DEFINITION A chi-square independence test is used to test the independence of two random variables. Using a chi-square test, you can determine whether the occurrence of one variable affects the probability of occurrence of the other variable.

To use the test,

1. The observed frequencies must be obtained from a random sample
2. Each expected frequency must be ≥ 5

The sampling distribution for the test is a chi-square distribution with

$$(r-1)(c-1)$$

degrees of freedom, where r and c are the number of rows and columns, respectively, of the contingency table. The test statistic for the chi-square independence test is

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

where O represents the observed frequencies and E represents the expected frequencies.

To begin the test we state the null hypothesis that the variables are independent and the alternative hypothesis that they are dependent.

Performing a Chi-Square Test for Independence (p 507)		
Procedure	Equations	Example2 (p 507)
Identify the claim. State the null and alternative hypotheses.	State H_0 and H_1	H_0 : CEO's ages are independent of company size H_1 : CEO's ages are dependent on company size.
Specify the level of significance	Specify α	$\alpha = 0.01$
Determine the degrees of freedom	$d.f. = (r-1)(c-1)$	$d.f. = (2-1)(5-1) = 4$
Find the critical value.	χ^2_α : Obtain from Table 6, Appendix B	$\chi^2_\alpha \geq 13.277$
Identify the rejection region	$\chi^2 \geq \chi^2_\alpha$	$\chi^2 \geq 13.277$
Calculate the test statistic	$\chi^2 = \sum \frac{(O-E)^2}{E}$	$\sum \frac{(O-E)^2}{E} \approx 77.9$ Note that O is in the table of actual CEO's ages above, and E is in the table of Expected CEO's ages (if independent of size) above

Make a decision to reject or fail to reject the null hypothesis	Reject if χ^2 is in the rejection region. Equivalently, we reject if the P-value (the probability of getting as extreme a value or more extreme) is $\leq \alpha$	Since $77.9 > 13.277$ we reject the null hypothesis Equivalently $P(X \geq 77.0) < \alpha$ so reject the null hypothesis. (Note Table 6 of Appendix B doesn't have a value less than 0.005.)
Interpret the decision in the context of the original claim		CEO's ages and company size are dependent.

The test statistic (77.887) is greater than the critical value obtained from Table 6, Appendix B (13.277) so the null hypothesis is rejected. (Alternatively the P-Value (0.000) is less than the level of significance, α (0.01) so the null hypothesis is rejected.)

- An urban geographer randomly samples 20 new residents of a neighborhood to determine their ratings of local bus service. The scale used is as follows: 0—very dissatisfied, 1—dissatisfied, 2—neutral, 3—satisfied, 4—very satisfied. The 20 responses are 0,4,3,2,2,1,1,2,1,0,0,1,2,1,3,4,2,0,4,1. Use the sign test to see whether the population median is 2.

Solution:

There are 5 observations above the hypothesized median. Because the sample size is larger than 10, we test using the sample proportion $p = 5/20 = 0.25$. Using the PROB-VALUE method the steps in this test are:

- $H_0: \pi = 0.5$ and $H_A: \pi \neq 0.5$
- We will use the Z-distribution
- We will use the 5%-level, thus $\alpha = 0.05$
- The test statistic is $z = (0.25 - 0.5) / \sqrt{0.25 / 20} = -2.24$
- Table A-4 shows that $P(|Z| > 2.24) \gg 0.025$.
- Because $\text{PROB-VALUE} < \alpha$, we reject H_0 . We conclude π is different than 0.5, and thus the median is different than 2.

4. A course in statistical methods was team-taught by two instructors, Professor Jovita Fontanez and Professor Clarence Old. Professor Fontanez used many active learning techniques, whereas Old employed traditional methods. As part of the course evaluation, students were asked to indicate their instructor preference. There was reason to think students would prefer Fontanez, and the sample obtained was consistent with that idea: of the 12 students surveyed, 8 preferred Professor Fontanez and 2 preferred Professor Old. The remaining students were unable to express a preference. Test the hypothesis that the students prefer Fontanez. (*Hint: Use the sign test.*)

Solution:

Although the sample is large enough for a normal approximation, we will use the binomial distribution to illustrate its application. Of the 12 observations, 8 preferred Prof. Fontanez, thus we need the probability of observing 8 or more successes in 12 trials of a Bernoulli process with the probability of success equal to 0.5. From Table A-1, we get

$$P(X \geq 8) = P(X = 8) + P(X = 9) + P(X = 10) + P(X = 11) + P(X = 12)$$

$$P(X \geq 8) = 0.1208 + 0.0537 + 0.0161 + 0.0029 + 0.0002 = 0.1937$$

Adopting the 5% uncertainty level, we see that $\text{PROB-VALUE} > \alpha$. Thus we fail to reject H_0 . We cannot conclude students prefer Fontanez.

5. Use the data in Table 10-8 to perform two Mann–Whitney tests: (a) compare uncontrolled intersections and intersections with yield signs, and (b) compare uncontrolled intersections and intersections with stop signs.

Solution:

- (a) The rank sums are 119.5 and 90.5 for the yield-signed and uncontrolled intersections respectively. Given the small sample size, we use an exact test rather than the normal approximation. The associated PROB-VALUE is 0.272. Adopting a 5% level of uncertainty, we fail to reject the hypothesis of no difference. We cannot conclude the samples were drawn from different populations.

(b) The rank sums are 130.5 and 59.5 for the stop-signed and uncontrolled intersections respectively. Given the small sample size, we use an exact test rather than the normal approximation. The associated PROB-VALUE is 0.013. Adopting a 5% level of uncertainty, we reject the hypothesis of no difference. We conclude the samples were drawn from different populations.

6. Solid-waste generation rates measured in metric tons per household per year are collected randomly in selected areas of a township. The areas are classified as high-density, low density, or sparsely settled. It is thought that generation rates probably differ because of differences in waste collection and opportunities for on-site storage. Do the following data support this hypothesis?

High Density	Low Density	Sparsely Settled
1.84	2.04	1.07
3.06	2.28	2.31
3.62	4.01	0.91
4.91	1.86	3.28
3.49	1.42	1.31

Solution:

We will use the multi-sample Kruskal-Wallis test with an uncertainty level $\alpha = 0.1$. The null hypothesis is that all samples have been drawn from the same population. The rank sums are 55, 39 and 26 for the high density, low density, and sparsely settled samples respectively. The Kruskal-Wallis statistic is

$$H = \frac{12}{15(15+1)} \left(\frac{55^2}{5} + \frac{39^2}{5} + \frac{26^2}{5} \right) - 3(15+1) = 4.22$$

Using the χ^2 distribution with $3 - 1 = 2$ degrees of freedom, the associated PROB-VALUE is 0.121. We fail to reject the null hypothesis. The sample does not support the hypothesis of differing waste generation rates.

7. The distances travelled to work by a random sample of 12 people to their places of work in 1996 and again in 2006 are shown in the following table.

Person	Distance (km)		Person	Distance (km)	
	1996	2006		1996	2006
1	8.6	8.8	7	7.7	6.5
2	7.7	7.1	8	9.1	9
3	7.7	7.6	9	8	7.1
4	6.8	6.4	10	8.1	8.8
5	9.6	9.1	11	8.7	7.2
6	7.2	7.2	12	7.3	6.4

Has the length of the journey to work changed over the decade?

Solution:

The sample can be considered as twelve paired observations. By taking differences between paired values, we get measures of the change for each individual. If the median change for the population is zero, we expect a sample to have a median difference near zero. Thus we will do a sign test for the median difference with a hypothesized value of zero. In other words, the hypotheses are $H_0: \eta = 0$ and $H_A: \eta \neq 0$. We denote samples values whose distance decreased with a minus sign. Sample values with a positive difference get a plus sign. The sample becomes

$$S = \{-, +, +, +, +, 0, +, -, +, -, +, +\}$$

Ignoring the tie, this is a sample of size 11 with 8 values above the hypothesized median. We are using Format (C) of Table 10-1, thus the PROB-VALUE is $2P(X \geq 8)$ where X is a binomial variable with $\pi = 0.5$. From the equation for the binomial, the PROB-VALUE is found to be 0.113. At the $\alpha = 10\%$ level, we fail to reject the null hypothesis. We cannot conclude there has been a change in distance.

8. One hundred randomly sampled residents of a city subject to periodic flooding are classified according to whether they are on the floodplain of the major river bisecting the

city or off the floodplain. These households are then surveyed to determine whether they currently have flood insurance of any kind. The survey results are as follows:

	On the Floodplain	Off the Floodplain
Insured	50	10
No Insurance	15	25

Test a relevant hypothesis.

Solution:

We will do a χ^2 test for a relationship between insurance and house location. The null hypothesis is no relationship (independence). Augmenting the data with expected frequencies, we have:

	On the Floodplain	Off the Floodplain
Insured	50 (39)	10 (21)
No Insurance	15 (26)	25 (14)

The corresponding χ^2 value is 22.16. Table A-8 shows that with 1 degree of freedom, $P(\chi^2 > 20)$ is zero to 3 decimal places. Thus for any reasonable level of uncertainty (any $\alpha < 0.0005$), we can reject the null hypothesis.

9. The occurrence of sunshine over a 30-day period was calculated as the percentage of time the sun was visible (i.e., not obscured by clouds). The daily percentages were:

Day	Percentage of sunshine	Day	Percentage of sunshine	Day	Percentage of sunshine
1	75	11	21	21	77
2	95	12	96	22	100
3	89	13	90	23	90

4	80	14	10	24	98
5	7	15	100	25	60
6	84	16	90	26	90
7	90	17	6	27	100
8	18	18	0	28	90
9	90	19	22	29	58
10	100	20	44	30	0

If we define a sunny day as one with over 50% sunshine, determine whether the pattern of occurrence of sunny days is random.

Solution:

For this we can use the number-of-runs test. Rather than calculate runs across two samples, here we will simply note if a day has 50% or more sunshine. The sample becomes

$$S = \{+, +, +, +, -, +, +, -, +, +, -, +, +, +, -, -, -, +, +, +, +, +, +, -, -\}$$

We see that the sample consists of 12 runs. There are $n_x = 21$ sunny days, and $n_y = 9$ cloudy days. Because $n_x < 20$, we cannot use the normal approximation given in Table 10-5. Instead the probability is computed using combinatorial rules, and is approximately 0.4. This is far too large for rejection of the randomness hypothesis. We cannot conclude the pattern is non-random.

10. Test the normality of the DO data (a) using the Kolmogorov–Smirnov test with the ungrouped data of Table 2-4 and (b) using the χ^2 test with $k = 6$ classes of Table 2-6.

Solution:

- (a) We will take the mean and standard deviation as known rather than estimated from the sample. Doing so results in calculated PROB-VALUES that are smaller than the true values (i.e., we are more likely to reject the null hypothesis). For the DO data the mean and standard deviation are 5.58 and 0.39 respectively. We sort the data, and then find the differences between the observed and expected

cumulative distributions. The table below shows the results for a few of the 50 observations:

x_i	$S(x_i)$	$F(x_i)$	$ S(x_i)-F(x_i) $
4.2	0.020	0.015	0.005
4.3	0.040	0.023	0.017
4.4	0.060	0.032	0.028
...
5.9	0.780	0.692	0.088
...
6.7	0.960	0.960	0.000
6.8	0.980	0.972	0.008
6.9	1.000	0.981	0.019

The maximum difference is 0.088. Table A-9 shows that with 50 degrees of freedom, the corresponding PROB-VALUE is about 0.6. We obviously cannot reject the hypothesis of normality.

- (b) Here we will take the mean and standard deviation as unknown, to be estimated from the sample. In other words, we estimate two parameters from the sample. In building the χ^2 table, we combine the first two and the last two categories in Table 2-6 to ensure at least 2 expected frequencies per cell. This reduces the number of categories 4, as seen in the table below:

Group	Minimum	Maximum	O_j	E_j	$(O_j-E_j)^2/E_j$
1	4.000	4.990	9	3.3	10.13
2	5.000	5.490	10	17.0	2.89
3	5.500	5.990	20	21.7	0.14
4	6.000	6.990	11	7.0	2.24

The observed Chi-square value is 15.4. With $k - p - 1 = 4 - 2 - 1 = 1$ degrees of freedom, Table A-8 shows that the PROB-VALUE is less than 0.0005. We therefore reject the null hypothesis.

Note that with only 4 classes, we can obtain only a rough idea of the distribution of DO. The 4 classes given in Table 2-6 do not yield a distribution that is at all similar to the normal distribution. In practice one would need many classes (and observations) for the χ^2 test to be reliable.

UNIT V

1	5	The hypothesis under test is	Null hypothesis	Alternative hypothesis	Null hypothesis	Alternative hypothesis	Null hypothesis
2	6	A test based on a test statistic is classified as	Randomized test	Non-randomized test	Randomized test	Non-randomized test	Randomized test
3	7	Student's t test is applicable in case of	Small samples	The sample of size between 5 and 30	Small samples	The sample size more than 100	Small samples
4	8	Reject H_0 when it is false is known as	Type I error	Type II error	Correct decision	Wrong decision	Correct decision
5	9	Accept H_0 when it is true is known as	Type I error	Type II error	Correct decision	Wrong decision	Correct decision
6	10	Accept H_0 when it is false is known as	Type I error	Type II error	Correct decision	Wrong decision	Type II error
7	11	The sample of 10 from the degree of freedom for student's t test is	df	df	df	df	df
8	12	If a sample mean is less than 20 then those samples may be considered as	One sample	Two sample	One sample	Two sample	One sample
9	13	The range of statistic is	0 to 1	0 to 1	0 to 1	0 to 1	0 to 1
10	14	The distribution used to test goodness of fit is	χ^2 distribution	χ^2 distribution	χ^2 distribution	χ^2 distribution	χ^2 distribution
11	15	For exactly known population value	$\mu = \mu_0$	$\mu \neq \mu_0$	$\mu = \mu_0$	$\mu \neq \mu_0$	$\mu \neq \mu_0$
12	16	95% of observed items of population mean are	$\Delta.M \pm 1.96 S.E$	$\Delta.M \pm 1.96 S.E$	$\Delta.M \pm 1.96 S.E$	$\Delta.M \pm 1.96 S.E$	$\Delta.M \pm 1.96 S.E$
13	17	99% of observed items of population mean are	$\Delta.M \pm 2.58 S.E$	$\Delta.M \pm 2.58 S.E$	$\Delta.M \pm 2.58 S.E$	$\Delta.M \pm 2.58 S.E$	$\Delta.M \pm 2.58 S.E$
14	18	A sort of the population selected for study is called as	Sample	Sample	Sample	Sample	Sample
15	19	Test statistic Z is	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$
16	20	In statistical test, standard error of sample mean is	S / \sqrt{n}	S / \sqrt{n}	S / \sqrt{n}	S / \sqrt{n}	S / \sqrt{n}
17	21	Null hypothesis is denoted by	H_0	H_0	H_0	H_0	H_0
18	22	Alternative hypothesis is denoted by	H_1	H_1	H_1	H_1	H_1
19	23	The value of the level of significance is	α	α	α	α	α
20	24	The value of the level of significance is	α	α	α	α	α
21	25	If the sample size is n , the standard deviation reduces to	Normal distribution	χ^2 distribution	Cauchy distribution	χ^2 distribution	Cauchy distribution
22	26	If n the sample size is larger than 30, the median distribution reduces to	Normal distribution	χ^2 distribution	Cauchy distribution	χ^2 distribution	Normal distribution
23	27	The $d.f$ for student's t based on a random sample of size n is	$n-1$	n	$n-1$	$n-1$	$n-1$
24	28	Degrees of freedom for chi-square test is	$n-1$	n	$n-1$	$n-1$	$n-1$
25	29	Student's t test is defined by the ratio	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$	$(\bar{X} - \mu_0) / (S / \sqrt{n})$
26	30	Chi square value for degree of freedom	χ^2	χ^2	χ^2	χ^2	χ^2
27	31	Degree of freedom is denoted by	df	df	df	df	df
28	32	Student's distribution the ratio was given by	Arthur H. Student	William S. Gosset	Carl Pearson	William S. Gosset	William S. Gosset
29	33	Chi Square test was first used in testing statistical hypothesis by	Karl Pearson	Robert Fisher	A. Fisher	Karl Pearson	Karl Pearson
30	34	While testing significance of the difference of two sample means in case of small samples, the degree of freedom is calculated by	$V = n_1 + n_2 - 1$	$V = n_1 + n_2 - 1$	$V = n_1 + n_2 - 2$	$V = n_1 + n_2 - 2$	$V = n_1 + n_2 - 2$
31	35	Analysis of variance utilizes	F test	Chi square test	Z test	t test	F test
32	36	Analysis of variance technique originated in	Archibald research	Industrial research	Business research	Archibald research	Archibald research
33	37	Analysis of variance technique was developed by	R. A. Fisher	Carl Pearson	Carl Pearson	R. A. Fisher	R. A. Fisher
34	38	Change of experiment was introduced by	R. A. Fisher	Sturgeson	Karl Pearson	Sturgeson	R. A. Fisher
35	39	The variance within samples is known as	Block	Block	Block	Block	Block



Reg.No. _____
[16BCU602 A]

KARPAGAM ACADEMY OF HIGHER EDUCATION

COIMBATORE - 21

DEPARTMENT OF BIOCHEMISTRY

III B.SC., BIOCHEMISTRY - SIX SEMESTER

16BCU602 A - BIOSTATISTICS

FIRST INTERNAL EXAMINATION - DECEMBER 2018

DATE: 17.12.18 AN

TIME: 2 HOURS

MAXIMUM: 50 MARKS

PART- A (20 x 1 = 20 marks)

Answer ALL the following:

1. Data taken from the publication, 'Agricultural Situation in India' will be considered as
(a) primary data (b) secondary data
(c) both primary and secondary data (d) tertiary data
2. Statistics deals with
(a) qualitative information (b) quantitative information
(c) both qualitative and quantitative information (d) only numerals
3. Year wise recording of data of food production will be called as:
(a) Geographical classification (b) Chronological classification
(c) Quantitative classification (d) qualitative classification
4. Who is the father of Biostatistics?
(a) R.A. Fisher (b) W. Gosset
(c) Sir Francis Galton (d) S.C. Gupta
5. Statistics can be considered as
a. an art b. a science c. both an art and a science
d. neither an art nor a science
6. The most suitable form of presentation for publicity and Propaganda is
(a) diagram (b) graph
(c) map (d) numerals

7. Histogram is suitable for the data presented as
(a) continuous grouped frequency distribution
(b) discrete grouped frequency distribution
(c) individual series
(d) discontinuous series
8. Numerical data presented in descriptive form are called
(a) classified presentation (b) tabular presentation
(c) graphical presentation (d) textual presentation
9. Mean is a measure of
(a) central value (b) dispersion
(c) correlation (d) significance
10. Mode is that value in a frequency distribution which possesses
(a) minimum frequency (b) frequency one
(c) maximum frequency (d) only two frequency
11. The most stable measure of central tendency is
(a) the mean (b) the median
(c) the mode (d) percentile
12. Sum of the deviations about mean is:
(a) minimum (b) zero (c) maximum (d) two
13. Mode of the following data 3, 6, 5, 7, 8, 4, 9
(a) 3 (b) 7 (c) no mode (d) 5
14. Which of the following is a measure of central tendency?
(a) Range (b) Quartile deviation
(c) Standard deviation (d) Median
15. Median is also called as
(a) first quartile (b) second quartile
(c) third quartile (d) fourth quartile
16. The census data published for state wise population in India will be known as
(a) Quantitative classification (b) Two-way classification
(c) Geographical classification (d) one way classification
17. Classification according to class-intervals would yield
(a) raw data (b) discrete data
(c) qualitative data (d) grouped data
18. In qualitative classification data are classified on the basis of
(a) attributes (b) time (c) location (d) class intervals

19 In geographical classification data are classified on the basis of
(a) area (b) attributes (c) time (d) location

20. Which source is one that itself collects the data?
(a) primary data (b) secondary data
(c) published data (d) tertiary data

PART B (3x2=6 marks)
Answer ALL the questions

21. List out the parts of table.
22. What are the various methods of collecting primary data?
23. Define Median and give Example

PART C (3x8=24 marks)
Answer ALL the questions

24.a. Explain about the Classification of data.

OR

b. Draw a suitable Pie Diagram to represent the following submitted as a part of the budget proposal of the govt. of India for the year 1995 - 96.

Item of Expenditure	Percentage
i) Interest	25
ii) Defense	15
iii) Other non plan expenditure	20
iv) States share of taxes and duties	15
v) State and UT plan assistance	10
vi) Central plan	15
Total	100

25.a. Calculate the Arithmetic Mean for the following data.

Height (cms):	160	161	162	163	164	165	166
No. of Persons:	27	36	43	78	65	48	28

OR

b. Explain in detail the various methods of collecting primary data.

26. a. Calculate the Median for the following Continuous Frequency Distribution.

Wages (in Rs.):	0 - 19	20 - 39	40 - 59	60 - 79	80 - 99
No. of Workers:	5	20	35	20	12

OR

b. Differentiate diagrams and graph. In what way the graphic presentation is advantageous than any other method?

class: M Bsc Botany

sub: Biostatistics

subcode: 16 Bw 60 2 A

No of copies: (30)

[Signature]