

UNIT 1: STATISTICAL DESCRIPTION OF DATA



LEARNING OBJECTIVES

After reading this chapter, students will be able to understand:

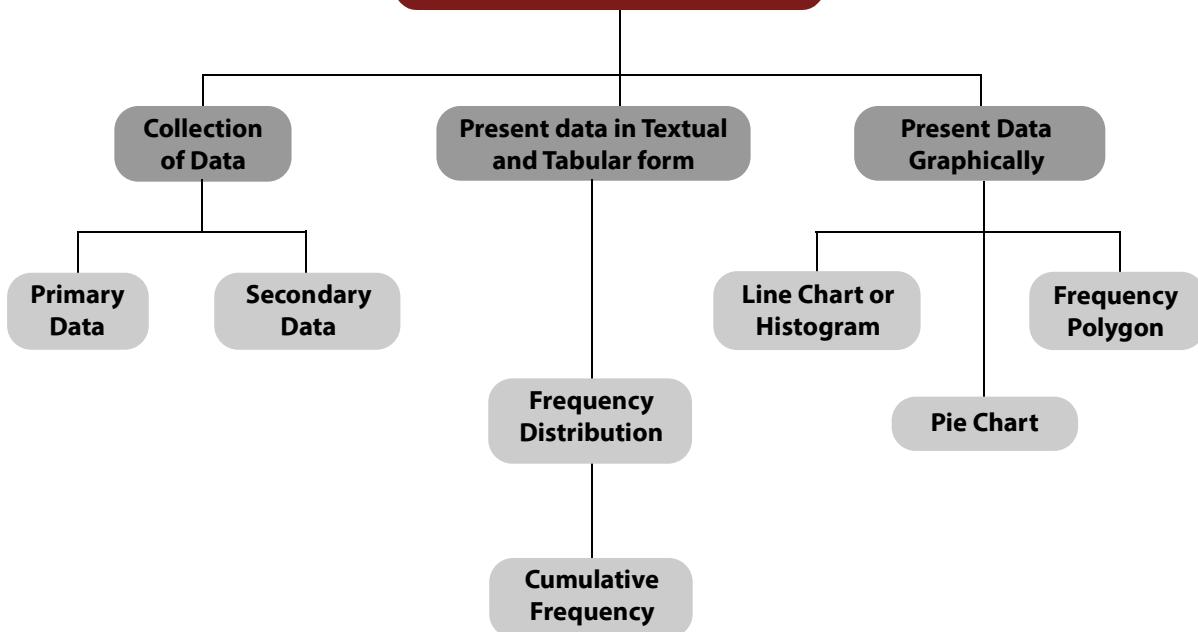
- ◆ Have a broad overview of the subject of statistics and application thereof.
- ◆ Know about data collection technique including the distinction of primary and secondary data.
- ◆ Know how to present data in textual and tabular format including the technique of creating frequency distribution and working out cumulative frequency.
- ◆ Know how to present data graphically using histogram, frequency polygon and pie chart.

UNIT OVERVIEW



Broad overview of the subject of statistics

Applications of Statistics





13.1.1 INTRODUCTION OF STATISTICS

The modern development in the field of not only Management, Commerce, Economics, Social Sciences, Mathematics and so on but also in our life like public services, defence, banking, insurance sector, tourism and hospitality, police and military etc. are dependent on a particular subject known as statistics. Statistics does play a vital role in enriching a specific domain by collecting data in that field, analysing the data by applying various statistical techniques and finally making statistical inferences about the domain. In the present world, statistics has almost a universal application. Our government applies statistics to make the economic planning in an effective and a pragmatic way. The businessman plan and expand their horizons of business on the basis of the analysis of the feedback data. The political parties try to impress the general public by presenting the statistics of their performances and accomplishments. Most of the research scholars of today also apply statistics to present their research papers in an authoritative manner. Thus the list of people using statistics goes on and on and on. Due to these factors, it is necessary to study the subject of statistics in an objective manner.

History of Statistics

Going through the history of ancient period and also that of medieval period, we do find the mention of statistics in many countries. However, there remains a question mark about the origin of the word 'statistics'. One view is that statistics originated from the Latin word 'status'. According to another school of thought, it had its origin in the Italian word 'statista'. Some scholars believe that the German word 'statistik' was later changed to statistics and another suggestion is that the French word 'statistique' was made as statistics with the passage of time.

In those days, statistics was analogous to state or, to be more precise, the data that are collected and maintained for the welfare of the people belonging to the state. We are thankful to Kautilya who had kept a record of births and deaths as well as some other precious records in his famous book 'Arthashastra' during Chandragupta's reign in the fourth century B.C. During the reign of Akbar in the sixteenth century A.D. We find statistical records on agriculture in Ain-i-Akbari written by Abu Fazl. Referring to Egypt, the first census was conducted by the Pharaoh during 300 B.C. to 2000 B.C.

Definition of Statistics

We may define statistics either in a singular sense or in a plural sense. Statistics, when used as a plural noun, may be defined as data qualitative as well as quantitative, that are collected, usually with a view of having statistical analysis.

However, statistics, when used as a singular noun, may be defined, as the scientific method that is employed for collecting, analysing and presenting data, leading finally to drawing statistical inferences about some important characteristics it means it is 'science of counting' or 'science of averages'.

Application of statistics

Among various applications of statistics, let us confine our discussions to the fields of Economics, Business Management and Commerce and Industry.

Economics

Modern developments in Economics have the roots in statistics. In fact, Economics and Statistics are closely associated. Time Series Analysis, Index Numbers, Demand Analysis etc. are some

overlapping areas of Economics and Statistics. In this connection, we may also mention Econometrics – a branch of Economics that interact with statistics in a very positive way. Conducting socio-economic surveys and analysing the data derived from it are made with the help of different statistical methods. Regression analysis, one of the numerous applications of statistics, plays a key role in Economics for making future projection of demand of goods, sales, prices, quantities etc. which are all ingredients of Economic planning.

Business Management

Gone are the days when the managers used to make decisions on the basis of hunches, intuition or trials and errors. Now a days, because of the never-ending complexity in the business and industry environment, most of the decision making processes rely upon different quantitative techniques which could be described as a combination of statistical methods and operations research techniques. So far as statistics is concerned, inferences about the universe from the knowledge of a part of it, known as sample, plays an important role in the development of certain criteria. Statistical decision theory is another component of statistics that focuses on the analysis of complicated business strategies with a list of alternatives – their merits as well as demerits.

Statistics in Commerce and Industry

In this age of cut-throat competition, like the modern managers, the industrialists and the businessmen are expanding their horizons of industries and businesses with the help of statistical procedures. Data on previous sales, raw materials, wages and salaries, products of identical nature of other factories etc are collected, analysed and experts are consulted in order to maximise profits. Measures of central tendency and dispersion, correlation and regression analysis, time series analysis, index numbers, sampling, statistical quality control are some of the statistical methods employed in commerce and industry.

Limitations of Statistics

Before applying statistical methods, we must be aware of the following limitations:

- I Statistics deals with the aggregates. An individual, to a statistician has no significance except the fact that it is a part of the aggregate.
- II Statistics is concerned with quantitative data. However, qualitative data also can be converted to quantitative data by providing a numerical description to the corresponding qualitative data.
- III Future projections of sales, production, price and quantity etc. are possible under a specific set of conditions. If any of these conditions is violated, projections are likely to be inaccurate.
- IV The theory of statistical inferences is built upon random sampling. If the rules for random sampling is not strictly adhered to, the conclusion drawn on the basis of these unrepresentative samples would be erroneous. In other words, the experts should be consulted before deciding the sampling scheme.



13.1.2 COLLECTION OF DATA

We may define 'data' as quantitative information about some particular characteristic(s) under consideration. Although a distinction can be made between a qualitative characteristic and a quantitative characteristic but so far as the statistical analysis of the characteristic is concerned,

we need to convert qualitative information to quantitative information by providing a numeric descriptions to the given characteristic. In this connection, we may note that a quantitative characteristic is known as a variable or in other words, a variable is a measurable quantity. Again, a variable may be either discrete or continuous. When a variable assumes a finite or a countably infinite number of isolated values, it is known as a discrete variable. Examples of discrete variables may be found in the number of petals in a flower, the number of misprints a book contains, the number of road accidents in a particular locality and so on. A variable, on the other hand, is known to be continuous if it can assume any value from a given interval. Examples of continuous variables may be provided by height, weight, sale, profit and so on. Finally, a qualitative characteristic is known as an attribute. The gender of a baby, the nationality of a person, the colour of a flower etc. are examples of attributes.

We can broadly classify data as

- (a) Primary;
- (b) Secondary.

Collection of data plays the very important role for any statistical analysis. The data which are collected for the first time by an investigator or agency are known as primary data whereas the data are known to be secondary if the data, as being already collected, are used by a different person or agency. Thus, if Prof. Das collects the data on the height of every student in his class, then these would be primary data for him. If, however, another person, say, Professor Bhargava uses the data, as collected by Prof. Das, for finding the average height of the students belonging to that class, then the data would be secondary for Prof. Bhargava.

Collection of Primary Data

The following methods are employed for the collection of primary data:

- (i) Interview method;
- (ii) Mailed questionnaire method;
- (iii) Observation method;
- (iv) Questionnaires filled and sent by enumerators.

Interview method again could be divided into (a) Personal Interview method, (b) Indirect Interview method and (c) Telephone Interview method.

In personal interview method, the investigator meets the respondents directly and collects the required information then and there from them. In case of a natural calamity like a super cyclone or an earthquake or an epidemic like plague, we may collect the necessary data much more quickly and accurately by applying this method.

If there are some practical problems in reaching the respondents directly, as in the case of a rail accident, then we may take recourse for conducting Indirect Interview where the investigator collects the necessary information from the persons associated with the problems.

Telephone interview method is a quick and rather non-expensive way to collect the primary data where the relevant information can be gathered by the researcher himself by contacting the interviewee over the phone. The first two methods, though more accurate, are inapplicable for covering a large area whereas the telephone interview, though less consistent, has a wide coverage.

The number of non-responses is maximum for this third method of data collection.

Mailed questionnaire method comprises of framing a well-drafted and soundly-sequenced questionnaire covering all the important aspects of the problem under consideration and sending them to the respondents with pre-paid stamp after providing all the necessary guidelines for filling up the questionnaire. Although a wide area can be covered using the mailed questionnaire method, the amount of non-responses is likely to be maximum in this method.

In observation nuclear, data are collected, as in the case of obtaining the data on the height and weight of a group of students, by direct observation or using instrument. Although this is likely to be the best method for data collection, it is time consuming, laborious and covers only a small area. Questionnaire form of data collection is used for larger enquiries from the persons who are surveyed. Enumerators collects information directly by interviewing the persons having information : Questions are explained and hence data is collected.

Sources of Secondary Data

There are many sources of getting secondary data. Some important sources are listed below:

- (a) International sources like WHO, ILO, IMF, World Bank etc.
- (b) Government sources like Statistical Abstract by CSO, Indian Agricultural Statistics by the Ministry of Food and Agriculture and so on.
- (c) Private and quasi-government sources like ISI, ICAR, NCERT etc.
- (d) Unpublished sources of various research institutes, researchers etc.

Scrutiny of Data

Since the statistical analyses are made only on the basis of data, it is necessary to check whether the data under consideration are accurate as well as consistent. No hard and fast rules can be recommended for the scrutiny of data. One must apply his intelligence, patience and experience while scrutinising the given information.

Errors in data may creep in while writing or copying the answer on the part of the enumerator. A keen observer can easily detect that type of error. Again, there may be two or more series of figures which are in some way or other related to each other. If the data for all the series are provided, they may be checked for internal consistency. As an example, if the data for population, area and density for some places are given, then we may verify whether they are internally consistent by examining whether the relation

$$\text{Density} = \frac{\text{Population}}{\text{Area}} \quad \text{holds.}$$

A good statistician can also detect whether the returns submitted by some enumerators are exactly of the same type thereby implying the lack of seriousness on the part of the enumerators. The bias of the enumerator also may be reflected by the returns submitted by him. This type of error can be rectified by asking the enumerator(s) to collect the data for the disputed cases once again.



13.1.3 PRESENTATION OF DATA

Once the data are collected and verified for their homogeneity and consistency, we need to present them in a neat and condensed form highlighting the essential features of the data. Any statistical analysis is dependent on a proper presentation of the data under consideration.

Classification or Organisation of Data

It may be defined as the process of arranging data on the basis of the characteristic under consideration into a number of groups or classes according to the similarities of the observations. Following are the objectives of classification of data:

- It puts the data in a neat, precise and condensed form so that it is easily understood and interpreted.
- It makes comparison possible between various characteristics, if necessary, and thereby finding the association or the lack of it between them.
- Statistical analysis is possible only for the classified data.
- It eliminates unnecessary details and makes data more readily understandable.

Data may be classified as -

- Chronological or Temporal or Time Series Data;
- Geographical or Spatial Series Data;
- Qualitative or Ordinal Data;
- Quantitative or Cardinal Data.

When the data are classified in respect of successive time points or intervals, they are known as time series data. The number of students appeared for CA final for the last twenty years, the production of a factory per month from 2000 to 2015 etc. are examples of time series data.

Data arranged region wise are known as geographical data. If we arrange the students appeared for CA final in the year 2015 in accordance with different states, then we come across Geographical Data.

Data classified in respect of an attribute are referred to as qualitative data. Data on nationality, gender, smoking habit of a group of individuals are examples of qualitative data. Lastly, when the data are classified in respect of a variable, say height, weight, profits, salaries etc., they are known as quantitative data.

Data may be further classified as *frequency data* and *non-frequency data*. The qualitative as well as quantitative data belong to the frequency group whereas time series data and geographical data belong to the non-frequency group.

Mode of Presentation of Data

Next, we consider the following mode of presentation of data:

- Textual presentation;
- Tabular presentation or Tabulation;
- Diagrammatic representation.

(a) Textual presentation

This method comprises presenting data with the help of a paragraph or a number of paragraphs. The official report of an enquiry commission is usually made by textual presentation. Following is an example of textual presentation.

'In 2009, out of a total of five thousand workers of Roy Enamel Factory, four thousand and two hundred were members of a Trade Union. The number of female workers was twenty per cent of the total workers out of which thirty per cent were members of the Trade Union.

In 2010, the number of workers belonging to the trade union was increased by twenty per cent as compared to 2009 of which four thousand and two hundred were male. The number of workers not belonging to trade union was nine hundred and fifty of which four hundred and fifty were females.'

The merit of this mode of presentation lies in its simplicity and even a layman can present data by this method. The observations with exact magnitude can be presented with the help of textual presentation. Furthermore, this type of presentation can be taken as the first step towards the other methods of presentation.

Textual presentation, however, is not preferred by a statistician simply because, it is dull, monotonous and comparison between different observations is not possible in this method. For manifold classification, this method cannot be recommended.

(b) Tabular presentation or Tabulation

Tabulation may be defined as systematic presentation of data with the help of a statistical table having a number of rows and columns and complete with reference number, title, description of rows as well as columns and foot notes, if any.

We may consider the following guidelines for tabulation :

- I A statistical table should be allotted a serial number along with a self-explanatory title.
- II The table under consideration should be divided into caption, Box-head, Stub and Body. Caption is the upper part of the table, describing the columns and sub-columns, if any. The Box-head is the entire upper part of the table which includes columns and sub-column numbers, unit(s) of measurement along with caption. Stub is the left part of the table providing the description of the rows. The body is the main part of the table that contains the numerical figures.
- III The table should be well-balanced in length and breadth.
- IV The data must be arranged in a table in such a way that comparison(s) between different figures are made possible without much labour and time. Also, the row totals, column totals, the units of measurement must be shown.
- V The data should be arranged intelligently in a well-balanced sequence and the presentation of data in the table should be appealing to the eyes as far as practicable.
- VI Notes describing the source of the data and bringing clarity and, if necessary, about any rows or columns known as footnotes, should be shown at the bottom part of the table.

The textual presentation of data, relating to the workers of Roy Enamel Factory is shown in the following table.

Table 13.1.1

Status of the workers of Roy Enamel factory on the basis of their trade union membership for 2009 and 2010.

Status	Member of TU			Non-member			Total		
	Year	M (1)	F (2)	T (3)=(1)+(2)	M (4)	F (5)	T (6)=(4)+(5)	M (7)	F (8)
2009	3900	300	4200	300	500	800	4200	800	5000
2010	4200	840	5040	500	450	950	4700	1290	5990

Source:

Footnote: TU, M, F and T stand for trade union, male, female and total respectively.

The tabulation method is usually preferred to textual presentation as

- (i) It facilitates comparison between rows and columns.
- (ii) Complicated data can also be represented using tabulation.
- (iii) It is a must for diagrammatic representation.
- (iv) Without tabulation, statistical analysis of data is not possible.

(c) Diagrammatic representation of data

Another alternative and attractive representation of statistical data is provided by charts, diagrams and pictures. Unlike the first two methods of representation of data, diagrammatic representation can be used for both the educated section and uneducated section of the society. Furthermore, any hidden trend present in the given data can be noticed only in this mode of representation. However, compared to tabulation, this is less accurate. So, if there is a priority for accuracy, we have to recommend tabulation.

We are going to consider the following types of diagrams :

- I Line diagram or Historiagram;
- II Bar diagram;
- III Pie chart.

I Line diagram or Historiagram

When the data vary over time, we take recourse to line diagram. In a simple line diagram, we plot each pair of values of (t, y_t) , y_t representing the time series at the time point t in the $t-y_t$ plane. The plotted points are then joined successively by line segments and the resulting chart is known as line-diagram.

When the time series exhibit a wide range of fluctuations, we may think of logarithmic or ratio chart where $\log y_t$ and not y_t is plotted against t . We use Multiple line chart for representing two or more related time series data expressed in the same unit and multiple-axis chart in somewhat similar situations if the variables are expressed in different units.

II Bar diagram

There are two types of bar diagrams namely, Horizontal Bar diagram and Vertical Bar diagram. While horizontal bar diagram is used for qualitative data or data varying over space, the vertical bar diagram is associated with quantitative data or time series data. Bars i.e. rectangles of equal width and usually of varying lengths are drawn either horizontally or vertically. We consider Multiple or Grouped Bar diagrams to compare related series. Component or sub-divided Bar diagrams are applied for representing data divided into a number of components. Finally, we use Divided Bar charts or Percentage Bar diagrams for comparing different components of a variable and also the relating of the components to the whole. For this situation, we may also use Pie chart or Pie diagram or circle diagram.



ILLUSTRATIONS:

Example 13.1.1: The profits in lakhs of Rupees of an industrial house for 2009, 2010, 2011, 2012, 2013, 2014, and 2015 are 5, 8, 9, 6, 12, 15 and 24 respectively. Represent these data using a suitable diagram.



SOLUTION:

We can represent the profits for 7 consecutive years by drawing either a line chart or a vertical bar chart. Fig. 13.1.1 shows a line chart and figure 13.1.2 shows the corresponding vertical bar chart.

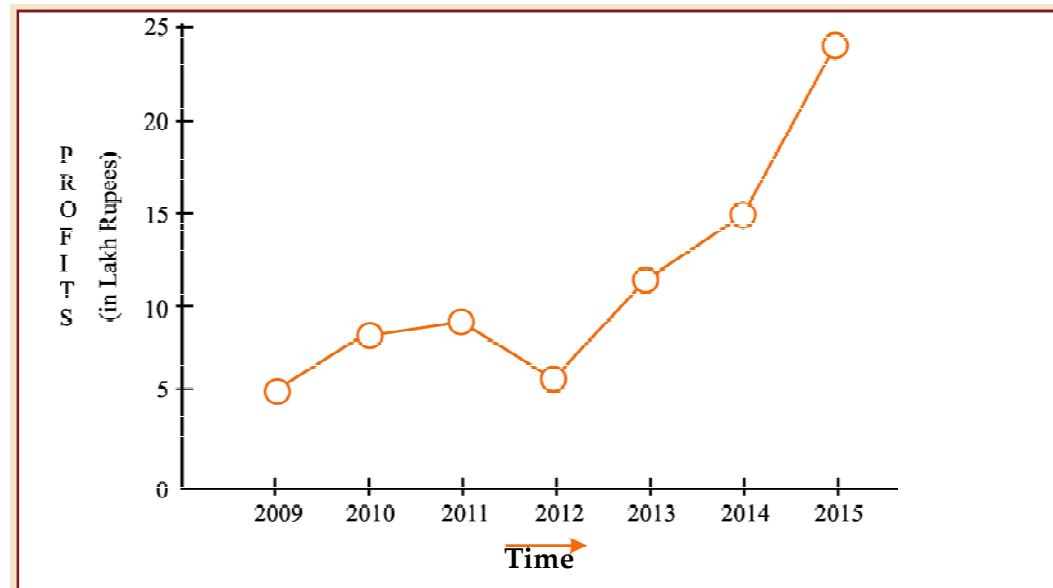
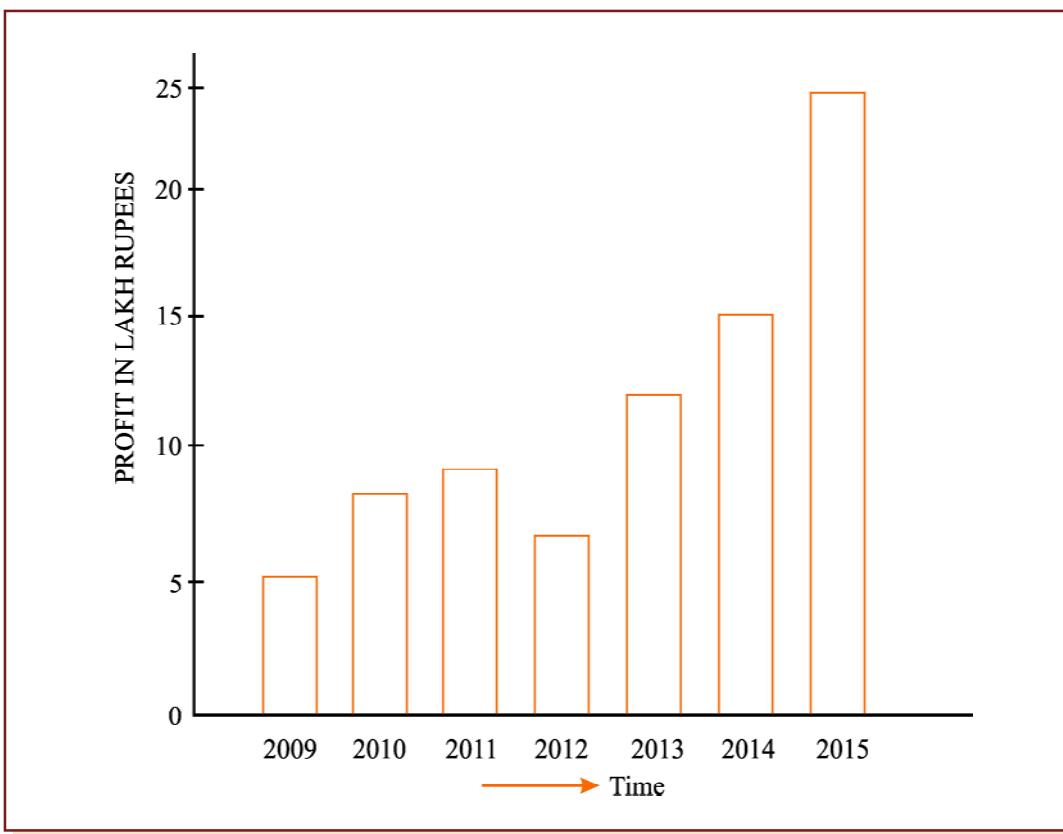


Figure 13.1.1

Showing line chart for the Profit of an Industrial House during 2009 to 2015.

**Figure 13.1.2**

Showing vertical bar diagram for the Profit of an Industrial house from 2009 to 2015.

Example 13.1.2: The production of wheat and rice of a region are given below :

Year	Production in metric tones	
	Wheat	Rice
2012	12	25
2013	15	30
2014	18	32
2015	19	36

Represent this information using a suitable diagram.

Solution:

We can represent this information by drawing a multiple line chart. Alternately, a multiple bar diagram may be considered. These are depicted in figure 13.1.3 and 13.1.4 respectively.

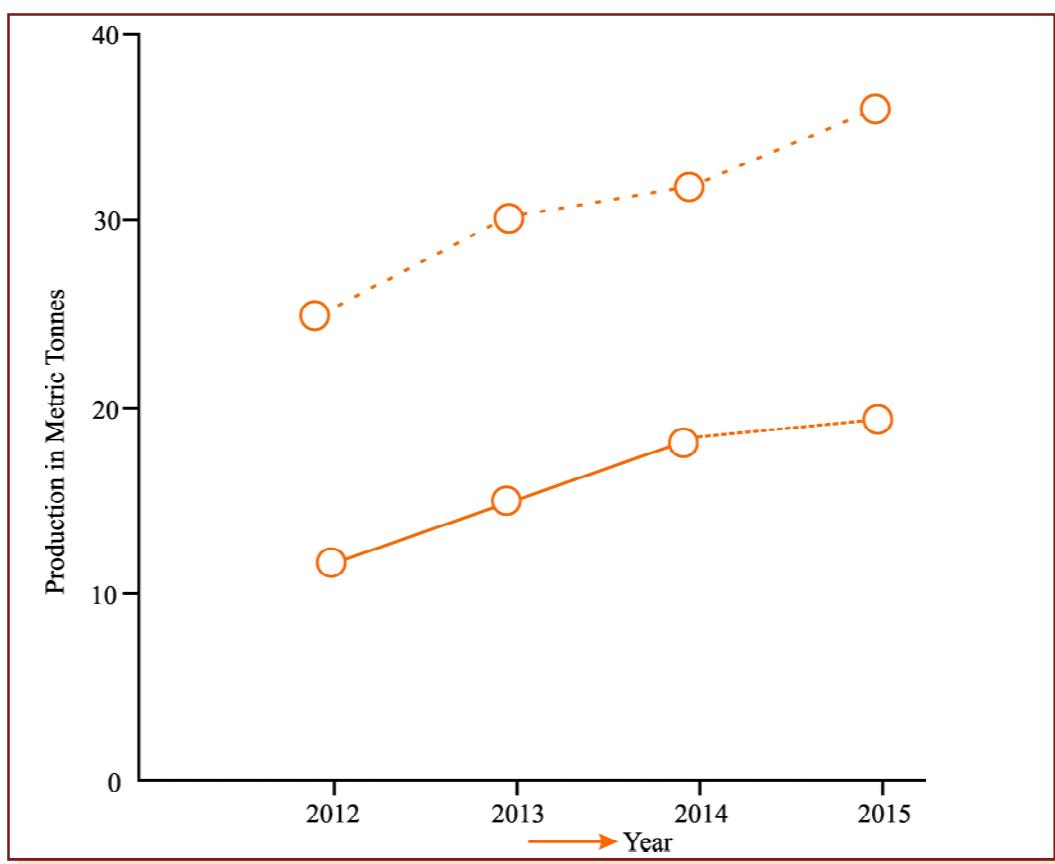


Figure 13.1.3

Multiple line chart showing production of wheat and rice of a region during 2012–2015.
(Dotted line represent production of rice and continuous line that of wheat).

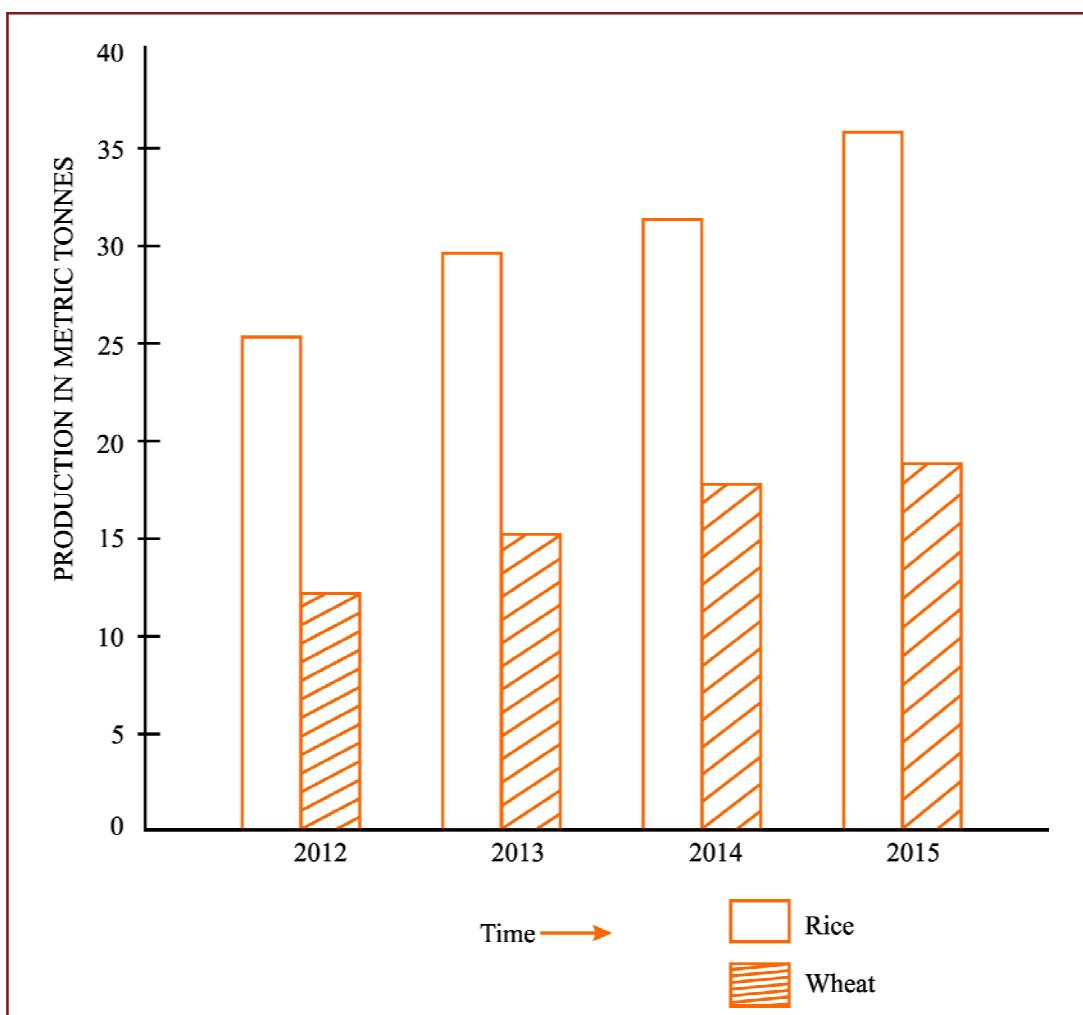


Figure 13.1.4

Multiple bar chart representing production of rice and wheat from 2012 to 2015.

Example 13.1.3: Draw an appropriate diagram with a view to represent the following data :

Source	Revenue in millions of (₹)
Customs	80
Excise	190
Income Tax	160
Corporate Tax	75
Miscellaneous	35

Solution:

Pie chart or divided bar chart would be the ideal diagram to represent this data. We consider Pie chart.

Table 13.1.2

Computation for drawing Pie chart

Source (1)	Revenue in Million rupees (2)	Central angle $= \frac{(2)}{\text{Total of (2)}} \times 360^\circ$
Customs	80	$\frac{80}{540} \times 360^\circ = 53^\circ$ (approx.)
Excise	190	$\frac{190}{540} \times 360^\circ = 127^\circ$
Income Tax	160	$\frac{160}{540} \times 360^\circ = 107^\circ$
Corporate Tax	75	$\frac{75}{540} \times 360^\circ = 50^\circ$
Miscellaneous	35	$\frac{35}{540} \times 360^\circ = 23^\circ$
Total	540	360°

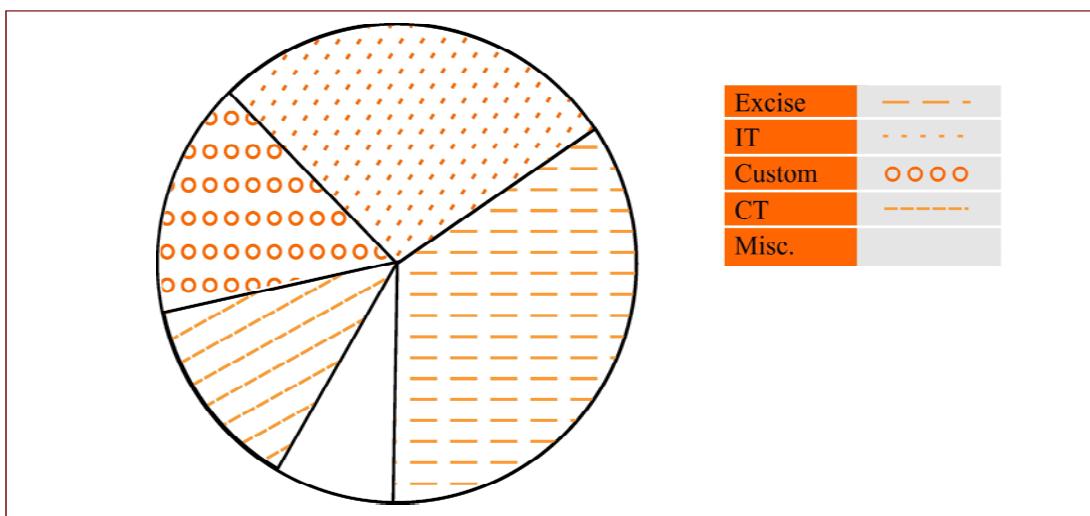


Figure 13.1.5
Pie chart showing the distribution of Revenue



13.1.4 FREQUENCY DISTRIBUTION

As discussed in the previous section, frequency data occur when we classify statistical data in respect of either a variable or an attribute. A frequency distribution may be defined as a tabular representation of statistical data, usually in an ascending order, relating to a measurable characteristic according to individual value or a group of values of the characteristic under study.

In case, the characteristic under consideration is an attribute, say nationality, then the tabulation is made by allotting numerical figures to the different classes the attribute may belong like, in this illustration, counting the number of Indian, British, French, German and so on. The qualitative characteristic is divided into a number of categories or classes which are mutually exclusive and exhaustive and the figures against all these classes are recorded. The figure corresponding to a particular class, signifying the number of times or how frequently a particular class occurs is known as the frequency of that class. Thus, the number of Indians, as found from the given data, signifies the frequency of the Indians. So frequency distribution is a statistical table that distributes the total frequency to a number of classes.

When tabulation is done in respect of a discrete random variable, it is known as Discrete or Ungrouped or simple Frequency Distribution and in case the characteristic under consideration is a continuous variable, such a classification is termed as Grouped Frequency Distribution. In case of a grouped frequency distribution, tabulation is done not against a single value as in the case of an attribute or a discrete random variable but against a group of values. The distribution of the number of car accidents in Delhi during 12 months of the year 2005 is an example of a ungrouped frequency distribution and the distribution of heights of the students of St. Xavier's College for the year 2004 is an example of a grouped frequency distribution.

Example 13.1.4: Following are the records of babies born in a nursing home in Bangalore during a week (B denoting Boy and G for Girl) :

B	G	G	B	G	G	B	B	G	G
G	G	B	B	B	G	B	B	G	B
B	B	G	B	B	B	G	G	B	G

Construct a frequency distribution according to gender.

Solution:

In order to construct a frequency distribution of babies in accordance with their gender, we count the number of male births and that of female births and present this information in the following table.

Table 13.1.3

Frequency distribution of babies according to Gender

Category	Number of births
Boy (B)	16
Girl (G)	14
Total	30

Frequency Distribution of a Variable

For the construction of a frequency distribution of a variable, we need to go through the following steps :

- I Find the largest and smallest observations and obtain the difference between them, known as Range, in case of a continuous variable.
- II Form a number of classes depending on the number of isolated values assumed by a discrete variable. In case of a continuous variable, find the number of class intervals using the relation, No. of class Interval \times class length \geq Range.
- III Present the class or class interval in a table known as frequency distribution table.
- IV Apply 'tally mark' i.e. a stroke against the occurrence of a particular value in a class or class interval.
- V Count the tally marks and present these numbers in the next column, known as frequency column, and finally check whether the total of all these class frequencies tally with the total number of observations.

Example 13.1.5: A review of the first 30 pages of a statistics book reveals the following printing mistakes:

0	1	3	3	2	5	6	0	1	0
4	1	1	0	2	3	2	5	0	4
2	3	2	2	3	3	4	6	1	4

Make a frequency distribution of printing mistakes.

Solution:

Since x , the printing mistakes, is a discrete variable, x can assume seven values 0, 1, 2, 3, 4, 5 and 6. Thus we have 7 classes, each class comprising a single value.

Table 13.1.4

Frequency Distribution of the number of printing mistakes of the first 30 pages of a book

Printing Mistake	Tally marks	Frequency (No. of Pages)
0	III	5
1	III	5
2	III I	6
3	III I	6
4	III	4
5	II	2
6	II	2
Total	-	30

Example 13.1.6: Following are the weights in kgs. of 36 BBA students of St. Xavier's College.

70 73 49 61 61 47 57 50 59
 59 68 45 55 65 68 56 68 55
 70 70 57 44 69 73 64 49 63
 65 70 65 62 64 73 67 60 50

Construct a frequency distribution of weights, taking class length as 5.

Solution:

$$\begin{aligned} \text{We have, Range} &= \text{Maximum weight} - \text{Minimum weight} \\ &= 73 \text{ kgs.} - 44 \text{ kgs.} \\ &= 29 \text{ kgs.} \end{aligned}$$

$$\text{No. of class interval} \times \text{class lengths} \approx \text{Range}$$

$$\Rightarrow \text{No. of class interval} \times 5 \approx 29$$

$$\Rightarrow \text{No. of class interval} = \frac{29}{5} \approx 6.$$

(We always take the next integer as the number of class intervals so as to include both the minimum and maximum values).

Table 13.1.5

Frequency Distribution of weights of 36 BBA Students

Weight in kg (Class Interval)	Tally marks	No. of Students (Frequency)
44-48	III	3
49-53	IIII	4
54-58	III	5
59-63	III II	7
64-68	III III	9
69-73	III III	8
Total	-	36

Some important terms associated with a frequency distribution

Class Limit (CL)

Corresponding to a class interval, the class limits may be defined as the minimum value and the maximum value the class interval may contain. The minimum value is known as the lower class limit (LCL) and the maximum value is known as the upper class limit (UCL). For the frequency distribution of weights of BBA Students, the LCL and UCL of the first class interval are 44 kgs. and 48 kgs. respectively.

Class Boundary (CB)

Class boundaries may be defined as the actual class limit of a class interval. For overlapping classification or mutually exclusive classification that excludes the upper class limits like 10–20, 20–30, 30–40, etc. the class boundaries coincide with the class limits. This is usually done for a continuous variable. However, for non-overlapping or mutually inclusive classification that includes both the class limits like 0–9, 10–19, 20–29,..... which is usually applicable for a discrete variable, we have

$$LCB = LCL - \frac{D}{2}$$

$$\text{and } UCB = UCL + \frac{D}{2}$$

where D is the difference between the LCL of the next class interval and the UCL of the given class interval. For the data presented in table 10.5, LCB of the first class interval

$$= 44 \text{ kgs.} - \frac{(49 - 48)}{2} \text{ kgs.}$$

$$= 43.50 \text{ kgs.}$$

and the corresponding UCB

$$\begin{aligned} &= 48 \text{ kgs.} + \frac{49 - 48}{2} \text{ kgs.} \\ &= 48.50 \text{ kgs.} \end{aligned}$$

Mid-point or Mid-value or class mark

Corresponding to a class interval, this may be defined as the total of the two class limits or class boundaries to be divided by 2. Thus, we have

$$\begin{aligned} \text{mid-point} &= \frac{\text{LCL} + \text{UCL}}{2} \\ &= \frac{\text{LCB} + \text{UCB}}{2} \end{aligned}$$

Referring to the distribution of weight of BBA students, the mid-points for the first two class intervals are

$$\frac{44 \text{ kgs.} + 48 \text{ kgs.}}{2} \text{ and } \frac{49 \text{ kgs.} + 53 \text{ kgs.}}{2}$$

i.e. 46 kgs. and 51 kgs. respectively.

Width or size of a class interval

The width of a class interval may be defined as the difference between the UCB and the LCB of that class interval. For the distribution of weights of BBA students, C, the class length or width is 48.50 kgs. – 43.50 kgs. = 5 kgs. for the first class interval. For the other class intervals also, C remains same.

Cumulative Frequency

The cumulative frequency corresponding to a value for a discrete variable and corresponding to a class boundary for a continuous variable may be defined as the number of observations less than the value or less than or equal to the class boundary. This definition refers to the less than cumulative frequency. We can define more than cumulative frequency in a similar manner. Both types of cumulative frequencies are shown in the following table.

Table 13.1.6

Cumulative Frequency Distribution of weights of 36 BBA students

Weight in kg (CB)	Cumulative Frequency	
	Less than	More than
43.50	0	33 + 3 or 36
48.50	0 + 3 or 3	29 + 4 or 33
53.50	3 + 4 or 7	24 + 5 or 29
58.50	7 + 5 or 12	17 + 7 or 24
63.50	12 + 7 or 19	8 + 9 or 17
68.50	19 + 9 or 28	0 + 8 or 8
73.50	28 + 8 or 36	0

Frequency density of a class interval

It may be defined as the ratio of the frequency of that class interval to the corresponding class length. The frequency densities for the first two class intervals of the frequency distribution of weights of BBA students are $3/5$ and $4/5$ i.e. 0.60 and 0.80 respectively.

Relative frequency and percentage frequency of a class interval

Relative frequency of a class interval may be defined as the ratio of the class frequency to the total frequency. Percentage frequency of a class interval may be defined as the ratio of class frequency to the total frequency, expressed as a percentage. For the last example, the relative frequencies for the first two class intervals are $3/36$ and $4/36$ respectively and the percentage frequencies are $300/36$ and $400/36$ respectively. It is quite obvious that whereas the relative frequencies add up to unity, the percentage frequencies add up to one hundred.



13.1.5 GRAPHICAL REPRESENTATION OF A FREQUENCY DISTRIBUTION

We consider the following types of graphical representation of frequency distribution :

- (i) Histogram or Area diagram;
- (ii) Frequency Polygon;
- (iii) Ogives or cumulative Frequency graphs.

(i) Histogram or Area diagram

This is a very convenient way to represent a frequency distribution. Histogram helps us to get an idea of the frequency curve of the variable under study. Some statistical measure can be obtained using a histogram. A comparison among the frequencies for different class intervals is possible in this mode of diagrammatic representation.

In order to draw a histogram, the class limits are first converted to the corresponding class boundaries and a series of adjacent rectangles, one against each class interval, with the class

interval as base or breadth and the frequency or frequency density usually when the class intervals are not uniform as length or altitude, is erected. The histogram for the distribution of weight of 36 BBA students is shown below. The mode of the weights has also been determined using the histogram.

i.e. Mode = 66.50 kgs.

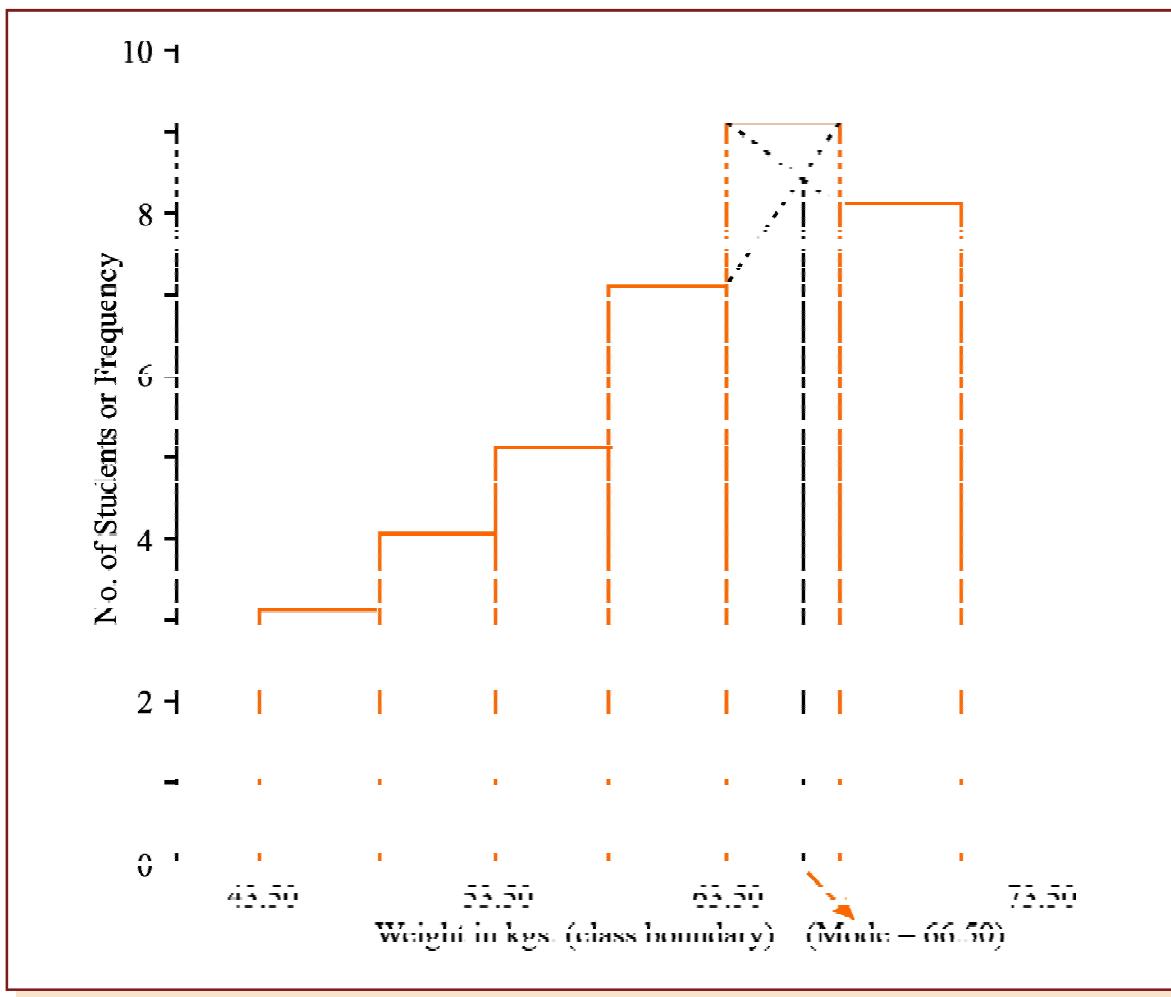


Figure 13.1.6

Showing histogram for the distribution of weight of 36 BBA students

(ii) Frequency Polygon

Usually frequency polygon is meant for single frequency distribution. However, we also apply it for grouped frequency distribution provided the width of the class intervals remains the same. A frequency curve can be regarded as a limiting form of frequency polygon. In order to draw a frequency polygon, we plot (x_i, f_i) for $i = 1, 2, 3, \dots, n$ with x_i denoting the mid-point of its class interval and f_i , the corresponding frequency, n being the number of class intervals. The plotted points are joined successively by line segments and the figure, so drawn, is given the shape of a polygon, a closed figure, by joining the two extreme ends of the drawn figure to two additional points $(x_0, 0)$ and $(x_{n+1}, 0)$.

The frequency polygon for the distribution of weights of BBA students is shown in Figure 13.7. We can also obtain a frequency polygon starting with a histogram by adding the mid-points of the upper sides of the rectangles successively and then completing the figure by joining the two ends as before.

Mid-points	No. of Students (Frequency)
46	3
51	4
56	5
61	7
66	9
71	8

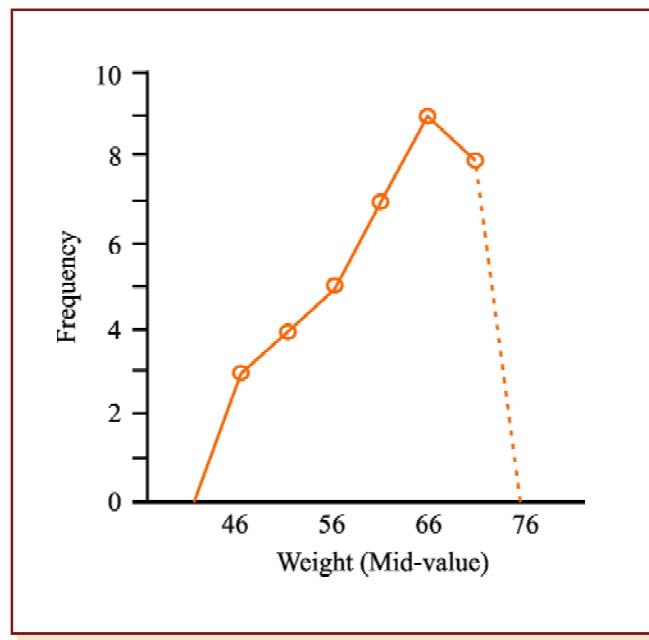


Figure 13.1.7

Showing frequency polygon for the distribution of height of 36 BBA students

(iii) Ogives or Cumulative Frequency Graph

By plotting cumulative frequency against the respective class boundary, we get ogives. As such there are two ogives – less than type ogives, obtained by taking less than cumulative frequency on the vertical axis and more than type ogives by plotting more than type cumulative frequency on the vertical axis and thereafter joining the plotted points successively by line segments. Ogives may be considered for obtaining quartiles graphically. If a perpendicular is drawn from the point of intersection of the two ogives on the horizontal axis, then the x-value of this point gives us the value of median, the second or middle quartile. Ogives further can be put into use for making short term projections.

Figure 13.8 depicts the ogives and the determination of the quartiles. This figure give us the following information.

1st quartile or lower quartile (Q_1) = 55 kgs.

2nd quartile or median (Q_2 or Me) = 62.50 kgs.

3rd quartile or upper quartile (Q_3) = 68 kgs.

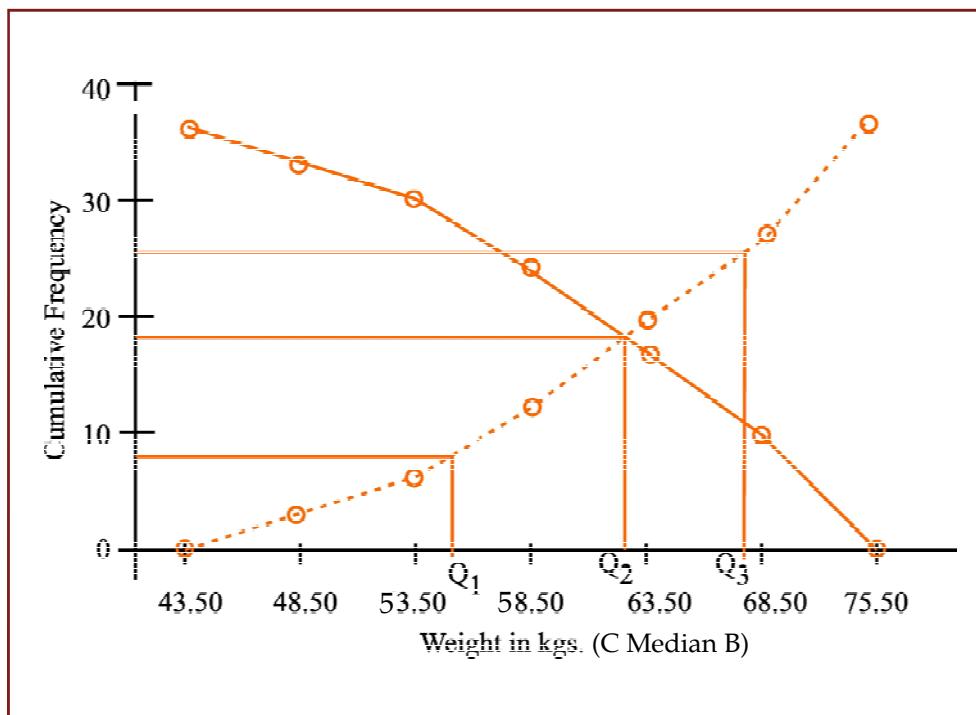


Figure 13.1.8

Showing the ogives for the distribution of weights of 36 BBA students

We find $Q_1 = 55$ kgs.

$Q_2 = Me = 62.50$ kgs.

$Q_3 = 68$ kgs.

Frequency Curve

A frequency curve is a smooth curve for which the total area is taken to be unity. It is a limiting form of a histogram or frequency polygon. The frequency curve for a distribution can be obtained by drawing a smooth and free hand curve through the mid-points of the upper sides of the rectangles forming the histogram.

There exist four types of frequency curves namely

- (a) Bell-shaped curve;
- (b) U-shaped curve;
- (c) J-shaped curve;
- (d) Mixed curve.

Most of the commonly used distributions provide bell-shaped curve, which, as suggested by the name, looks almost like a bell. The distribution of height, weight, mark, profit etc. usually belong to this category. On a bell-shaped curve, the frequency, starting from a rather low value, gradually reaches the maximum value, somewhere near the central part and then gradually decreases to reach its lowest value at the other extremity.

For a U-shaped curve, the frequency is minimum near the central part and the frequency slowly but steadily reaches its maximum at the two extremities. The distribution of Kolkata bound commuters belongs to this type of curve as there are maximum number of commuters during the peak hours in the morning and in the evening.

The J-shaped curve starts with a minimum frequency and then gradually reaches its maximum frequency at the other extremity. The distribution of commuters coming to Kolkata from the early morning hour to peak morning hour follows such a distribution. Sometimes, we may also come across an inverted J-shaped frequency curve.

Lastly, we may have a combination of these frequency curves, known as mixed curve. These are exhibited in the following figures.

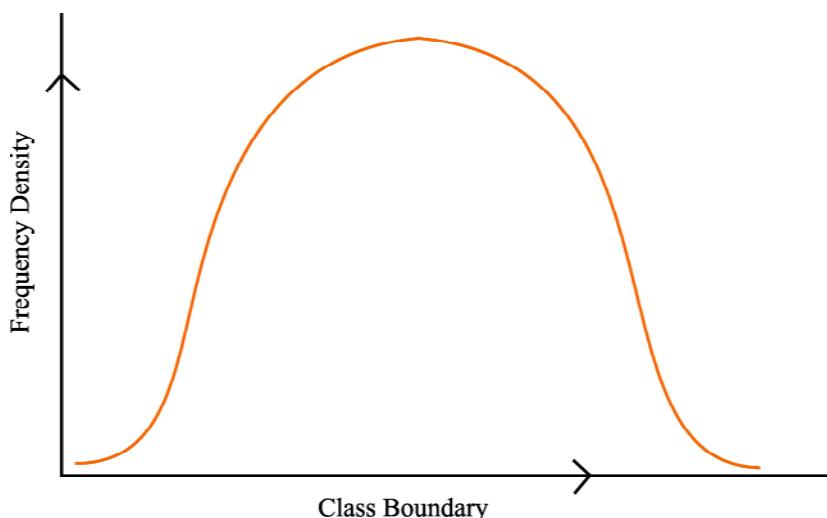


Figure 13.1.9
Bell-shaped curve

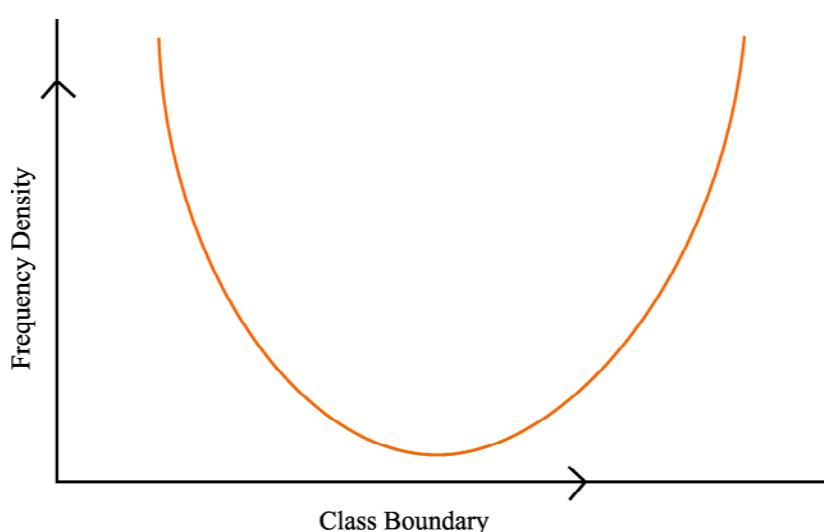


Figure 13.1.10
U-shaped curve

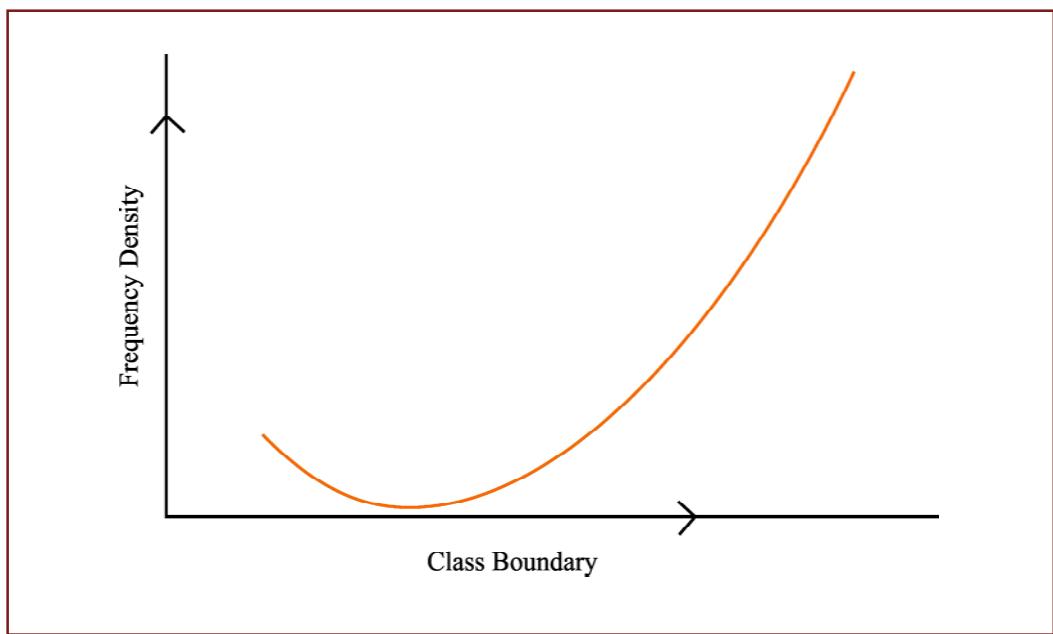


Figure 13.1.11
J-shaped curve

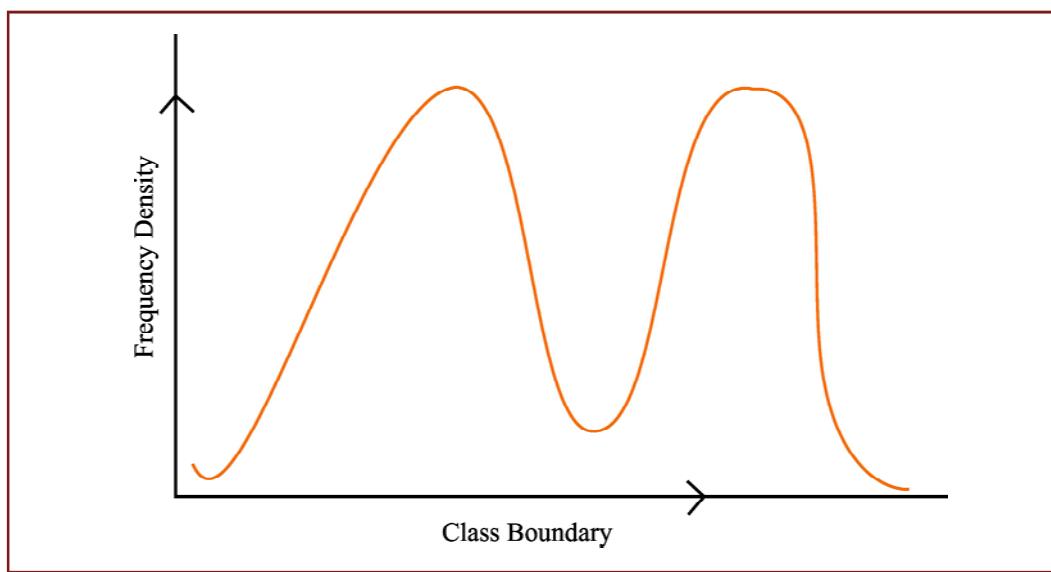


Figure 13.1.12
Mixed curve



SUMMARY

- ◆ Statistics deals with the aggregates. An individual, to a statistician has no significance except the fact that it is a part of the aggregate.
- ◆ Statistics is concerned with quantitative data. However, qualitative data also can be converted to quantitative data by providing a numerical description to the corresponding qualitative data.
- ◆ The theory of statistical inferences is built upon random sampling. If the rules for random sampling are not strictly adhered to, the conclusion drawn on the basis of these unrepresentative samples would be erroneous.
- ◆ We can broadly classify data as
 - (a) Primary;
 - (b) Secondary.
- ◆ Mode of Presentation of Data
 - (a) Textual presentation;
 - (b) Tabular presentation or Tabulation;
 - (c) Diagrammatic representation.
- ◆ The types of diagrams:
 - (a) Line diagram or Historiogram;
 - (b) Bar diagram;
 - (c) Pie chart.
- ◆ Frequency Distribution of a Variable
 - (a) Find the largest and smallest observations and obtain the difference between them, known as Range, in case of a continuous variable.
 - (b) Form a number of classes depending on the number of isolated values assumed by a discrete variable. In case of a continuous variable, find the number of class intervals using the relation, No. of class Interval \times class length \equiv Range.
 - (c) Present the class or class interval in a table known as frequency distribution table.
 - (d) Apply 'tally mark' i.e. a stroke against the occurrence of a particular value in a class or class interval.
 - (e) Count the tally marks and present these numbers in the next column, known as frequency column, and finally check whether the total of all these class frequencies tally with the total number of observations.

 **UNIT I EXERCISE****Set A**

Answer the following questions. Each question carries 1 mark.

1. Which of the following statements is false?
 - (a) Statistics is derived from the Latin word 'Status'
 - (b) Statistics is derived from the Italian word 'Statista'
 - (c) Statistics is derived from the French word 'Statistik'
 - (d) None of these.
2. Statistics is defined in terms of numerical data in the
 - (a) Singular sense
 - (b) Plural sense
 - (c) Either (a) or (b)
 - (d) Both (a) and (b).
3. Statistics is applied in
 - (a) Economics
 - (b) Business management
 - (c) Commerce and industry
 - (d) All these.
4. Statistics is concerned with
 - (a) Qualitative information
 - (b) Quantitative information
 - (c) (a) or (b)
 - (d) Both (a) and (b).
5. An attribute is
 - (a) A qualitative characteristic
 - (b) A quantitative characteristic
 - (c) A measurable characteristic
 - (d) All these.
6. Annual income of a person is
 - (a) An attribute
 - (b) A discrete variable
 - (c) A continuous variable
 - (d) (b) or (c).
7. Marks of a student is an example of
 - (a) An attribute
 - (b) A discrete variable
 - (c) A continuous variable
 - (d) None of these.
8. Nationality of a student is
 - (a) An attribute
 - (b) A continuous variable
 - (c) A discrete variable
 - (d) (a) or (c).
9. Drinking habit of a person is
 - (a) An attribute
 - (b) A variable
 - (c) A discrete variable
 - (d) A continuous variable.

29. In tabulation source of the data, if any, is shown in the
(a) Footnote (b) Body
(c) Stub (d) Caption.

30. Which of the following statements is untrue for tabulation?
(a) Statistical analysis of data requires tabulation
(b) It facilitates comparison between rows and not columns
(c) Complicated data can be presented
(d) Diagrammatic representation of data requires tabulation.

31. Hidden trend, if any, in the data can be noticed in
(a) Textual presentation (b) Tabulation
(c) Diagrammatic representation (d) All these.

32. Diagrammatic representation of data is done by
(a) Diagrams (b) Charts
(c) Pictures (d) All these.

33. The most accurate mode of data presentation is
(a) Diagrammatic method (b) Tabulation
(c) Textual presentation (d) None of these.

34. The chart that uses logarithm of the variable is known as
(a) Line chart (b) Ratio chart
(c) Multiple line chart (d) Component line chart.

35. Multiple line chart is applied for
(a) Showing multiple charts
(b) Two or more related time series when the variables are expressed in the same unit
(c) Two or more related time series when the variables are expressed in different unit
(d) Multiple variations in the time series.

36. Multiple axis line chart is considered when
(a) There is more than one time series (b) The units of the variables are different
(c) (a) or (b) (d) (a) and (b).

37. Horizontal bar diagram is used for
(a) Qualitative data (b) Data varying over time
(c) Data varying over space (d) (a) or (c).

38. Vertical bar diagram is applicable when
- (a) The data are qualitative
 - (b) The data are quantitative
 - (c) When the data vary over time
 - (d) (b) or (c).
39. Divided bar chart is considered for
- (a) Comparing different components of a variable
 - (b) The relation of different components to the total
 - (c) (a) or (b)
 - (d) (a) and (b).
40. In order to compare two or more related series, we consider
- (a) Multiple bar chart
 - (b) Grouped bar chart
 - (c) (a) or (b)
 - (d) (a) and (b).
41. Pie-diagram is used for
- (a) Comparing different components and their relation to the total
 - (b) Representing qualitative data in a circle
 - (c) Representing quantitative data in circle
 - (d) (b) or (c).
42. A frequency distribution
- (a) Arranges observations in an increasing order
 - (b) Arranges observation in terms of a number of groups
 - (c) Relates to a measurable characteristic
 - (d) All these.
43. The frequency distribution of a continuous variable is known as
- (a) Grouped frequency distribution
 - (b) Simple frequency distribution
 - (c) (a) or (b)
 - (d) (a) and (b).

44. The distribution of shares is an example of the frequency distribution of
- (a) A discrete variable
 - (b) A continuous variable
 - (c) An attribute
 - (d) (a) or (c).
45. The distribution of profits of a blue-chip company relates to
- (a) Discrete variable
 - (b) Continuous variable
 - (c) Attributes
 - (d) (a) or (b).
46. Mutually exclusive classification
- (a) Excludes both the class limits
 - (b) Excludes the upper class limit but includes the lower class limit
 - (c) Includes the upper class limit but excludes the lower class limit
 - (d) Either (b) or (c).
47. Mutually inclusive classification is usually meant for
- (a) A discrete variable
 - (b) A continuous variable
 - (c) An attribute
 - (d) All these.
48. Mutually exclusive classification is usually meant for
- (a) A discrete variable
 - (b) A continuous variable
 - (c) An attribute
 - (d) Any of these.
49. The LCB is
- (a) An upper limit to LCL
 - (b) A lower limit to LCL
 - (c) (a) and (b)
 - (d) (a) or (b).

59. Most of the commonly used frequency curves are

 - (a) Mixed
 - (b) Inverted J-shaped
 - (c) U-shaped
 - (d) Bell-shaped.

60. The distribution of profits of a company follows

 - (a) J-shaped frequency curve
 - (b) U-shaped frequency curve
 - (c) Bell-shaped frequency curve
 - (d) Any of these.

Set B

Answer the following questions. Each question carries 2 marks.

How many students got marks more than 30?

7. Find the number of observations between 250 and 300 from the following data :

Value	: More than 200	More than 250	More than 300	More than 350
No. of observations :	56	38	15	0
(a) 56	(b) 23	(c) 15	(d) 8	

Set C

Answer the following questions. Each question carries 5 marks.

1. In a study about the male and female students of commerce and science departments of a college in 5 years, the following datas were obtained :

1995	2000
70% male students	75% male students
65% read Commerce	40% read Science
20% of female students read Science	50% of male students read Commerce
3000 total No. of students	3600 total No. of students.

After combining 1995 and 2000 if x denotes the ratio of female commerce student to female Science student and y denotes the ratio of male commerce student to male Science student, then

- (a) $x = y$ (b) $x > y$ (c) $x < y$ (d) $x \geq y$

2. In a study relating to the labourers of a jute mill in West Bengal, the following information was collected.

'Twenty per cent of the total employees were females and forty per cent of them were married. Thirty female workers were not members of Trade Union. Compared to this, out of 600 male workers 500 were members of Trade Union and fifty per cent of the male workers were married. The unmarried non-member male employees were 60 which formed ten per cent of the total male employees. The unmarried non-members of the employees were 80'. On the basis of this information, the ratio of married male non-members to the married female non-members is

- (a) 1 : 3 (b) 3 : 1 (c) 4 : 1 (d) 5 : 1

3. The weight of 50 students in pounds are given below :

82,	95,	120,	174,	179,	176,	159,	91,	85,	175
88,	160,	97,	133,	159,	176,	151,	115,	105,	172
170,	128,	112,	101,	123,	117,	93,	117,	99,	90
113,	119,	129,	134,	178,	105,	147,	107,	155,	157
98,	117,	95,	135,	175,	97,	160,	168,	144,	175

If the data are arranged in the form of a frequency distribution with class intervals as 81-100, 101-120, 121-140, 141-160 and 161-180, then the frequencies for these 5 class intervals are

- (a) 6, 9, 10, 11, 14 (b) 12, 8, 7, 11, 12 (c) 10, 12, 8, 11, 9 (d) 12, 12, 6, 9, 11

4. The following data relate to the marks of 48 students in statistics :

56,	10,	54,	38,	21,	43,	12,	22
48,	51,	39,	26,	12,	17,	36,	19
48,	36,	15,	33,	30,	62,	57,	17
5,	17,	45,	46,	43,	55,	57,	38
43,	28,	32,	35,	54,	27,	17,	16
11,	43,	45,	2,	16,	46,	28,	45

What are the frequency densities for the class intervals 30-39, 40-49 and 50-59

- (a) 0.20, 0.50, 0.90
 (b) 0.70, 0.90, 1.10
 (c) 0.1875, 0.1667, 0.2083
 (d) 0.90, 1.1, 0.7

5. The following information relates to the age of death of 50 persons in an area :

36,	48,	50,	45,	49,	31,	50,	48,	42,	57
43,	40,	32,	41,	39,	39,	43,	47,	45,	52
47,	48,	53,	37,	48,	50,	41,	49,	50,	53
38,	41,	49,	45,	36,	39,	31,	48,	59,	48
37,	49,	53,	51,	54,	59,	48,	38,	39,	45

If the class intervals are 31-33, 34-36, 37-39, Then the percentage frequencies for the last five class intervals are

- (a) 18, 18, 10, 2 and 4. (b) 10, 15, 18, 4 and 2. (c) 14, 18, 20, 10 and 2.
 (d) 10, 12, 16, 4 and 6.

ANSWERS

Set A

- | | | | | | |
|---------|---------|---------|---------|---------|---------|
| 1. (c) | 2. (b) | 3. (d) | 4. (d) | 5. (a) | 6. (b) |
| 7. (b) | 8. (a) | 9. (a) | 10. (c) | 11. (b) | 12. (a) |
| 13. (d) | 14. (c) | 15. (a) | 16. (c) | 17. (b) | 18. (a) |
| 19. (a) | 20. (d) | 21. (c) | 22. (a) | 23. (b) | 24. (c) |

- | | | | | | |
|----------------|----------------|----------------|----------------|----------------|----------------|
| 25. (d) | 26. (d) | 27. (c) | 28. (a) | 29. (a) | 30. (b) |
| 31. (c) | 32. (d) | 33. (b) | 34. (b) | 35. (b) | 36. (d) |
| 37. (d) | 38. (d) | 39. (d) | 40. (c) | 41. (a) | 42. (d) |
| 43. (a) | 44. (a) | 45. (b) | 46. (b) | 47. (a) | 48. (d) |
| 49. (b) | 50. (a) | 51. (a) | 52. (a) | 53. (b) | 54. (a) |
| 55. (a) | 56. (c) | 57. (b) | 58. (c) | 59. (d) | 60. (c) |

Set B

- | | | | | | |
|---------------|---------------|---------------|---------------|---------------|---------------|
| 1. (a) | 2. (b) | 3. (d) | 4. (d) | 5. (a) | 6. (c) |
| 7. (b) | | | | | |

Set C

- | | | | | |
|---------------|---------------|---------------|---------------|---------------|
| 1. (b) | 2. (c) | 3. (d) | 4. (d) | 5. (a) |
|---------------|---------------|---------------|---------------|---------------|

ADDITIONAL QUESTION BANK

1. Graph is a

(a) Line diagram	(b) Bar diagram	(c) Pie diagram	(d) Pictogram
------------------	-----------------	-----------------	---------------
2. Details are shown by

(a) Charts	(b) Tabular presentation
(c) both	(d) none
3. The relationship between two variables are shown in

(a) Pictogram	(b) Histogram	(c) Bar diagram	(d) Line diagram
---------------	---------------	-----------------	------------------
4. In general the number of types of tabulation are

(a) two	(b) three	(c) one	(d) four
---------	-----------	---------	----------
5. A table has

(a) four	(b) two	(c) five	(d) none parts.
----------	---------	----------	-----------------
6. The number of errors in Statistics are

(a) one	(b) two	(c) three	(d) four
---------	---------	-----------	----------
7. The number of “Frequency distribution” is

(a) two	(b) one	(c) five	(d) four
---------	---------	----------	----------
8. $(\text{Class frequency}) / (\text{Width of the class})$ is defined as

(a) Frequency density	(b) Frequency distribution
(c) both	(d) none

9. Tally marks determines
(a) class width (b) class boundary (c) class limit (d) class frequency
10. Cumulative Frequency Distribution is a
(a) graph (b) frequency (c) Statistical Table (d) distribution
11. To find the number of observations less than any given value
(a) Single frequency distribution (b) Grouped frequency distribution
(c) Cumulative frequency distribution (d) None is used.
12. An area diagram is
(a) Histogram (b) Frequency Polygon
(c) Ogive (d) none
13. When all classes have a common width
(a) Pie Chart (b) Frequency Polygon
(c) both (d) none is used.
14. An approximate idea of the shape of frequency curve is given by
(a) Ogive (b) Frequency Polygon
(c) both (d) none
15. Ogive is a
(a) Line diagram (b) Bar diagram (c) both (d) none
16. Unequal widths of classes in the frequency distribution do not cause any difficulty in the construction of
(a) Ogive (b) Frequency Polygon
(c) Histogram (d) none
17. The graphical representation of a cumulative frequency distribution is called
(a) Histogram (b) Ogive (c) both (d) none.
18. The most common form of diagrammatic representation of a grouped frequency distribution is
(a) Ogive (b) Histogram (c) Frequency Polygon (d) none
19. Vertical bar chart may appear somewhat alike
(a) Histogram (b) Frequency Polygon
(c) both (d) none
20. The number of types of cumulative frequency is
(a) one (b) two (c) three (d) four

21. A representative value of the class interval for the calculation of mean, standard deviation, mean deviation etc. is
 - (a) class interval
 - (b) class limit
 - (c) class mark
 - (d) none
22. The number of observations falling within a class is called
 - (a) density
 - (b) frequency
 - (c) both
 - (d) none
23. Classes with zero frequencies are called
 - (a) nill class
 - (b) empty class
 - (c) class
 - (d) none
24. For determining the class frequencies it is necessary that these classes are
 - (a) mutually exclusive
 - (b) not mutually exclusive
 - (c) independent
 - (d) none
25. Most extreme values which would ever be included in a class interval are called
 - (a) class limits
 - (b) class interval
 - (c) class boundaries
 - (d) none
26. The value exactly at the middle of a class interval is called
 - (a) class mark
 - (b) mid value
 - (c) both
 - (d) none
27. Difference between the lower and the upper class boundaries is
 - (a) width
 - (b) size
 - (c) both
 - (d) none
28. In the construction of a frequency distribution, it is generally preferable to have classes of
 - (a) equal width
 - (b) unequal width
 - (c) maximum
 - (d) none
29. Frequency density is used in the construction of
 - (a) Histogram
 - (b) Ogive
 - (c) Frequency Polygon
 - (d) none when the classes are of unequal width.
30. "Cumulative Frequency" only refers to the
 - (a) less-than type
 - (b) more-than type
 - (c) both
 - (d) none
31. For the construction of a grouped frequency distribution
 - (a) class boundaries
 - (b) class limits
 - (c) both
 - (d) none are used.
32. In all Statistical calculations and diagrams involving end points of classes
 - (a) class boundaries
 - (b) class value
 - (c) both
 - (d) none are used.
33. Upper limit of any class is _____ from the lower limit of the next class
 - (a) same
 - (b) different
 - (c) both
 - (d) none
34. Upper boundary of any class coincides with the Lower boundary of the next class.
 - (a) true
 - (b) false
 - (c) both
 - (d) none.

ANSWERS

- | | | | | |
|---------|---------|---------|---------|---------|
| 1. (a) | 2. (b) | 3. (d) | 4. (a) | 5. (c) |
| 6. (b) | 7. (a) | 8. (a) | 9. (d) | 10. (c) |
| 11. (c) | 12. (a) | 13. (b) | 14. (b) | 15. (a) |
| 16. (c) | 17. (b) | 18. (b) | 19. (a) | 20. (b) |
| 21. (c) | 22. (b) | 23. (b) | 24. (a) | 25. (c) |

- | | | | | |
|----------------|----------------|----------------|----------------|----------------|
| 26. (c) | 27. (c) | 28. (a) | 29. (a) | 30. (a) |
| 31. (b) | 32. (a) | 33. (b) | 34. (a) | 35. (a) |
| 36. (b) | 37. (a) | 38. (b) | 39. (a) | 40. (c) |
| 41. (a) | 42. (a) | 43. (a) | 44. (b) | 45. (c) |
| 46. (b) | 47. (b) | 48. (c) | 49. (b) | 50. (a) |
| 51. (c) | 52. (d) | 53. (a) | 54. (a) | 55. (b) |
| 56. (a) | 57. (a) | 58. (a) | | |