## Introduction to Statistics

Statistics can be defined as the Science of collecting, displaying and analysing data.

We collect data from the observations less than whom they are representing, which is called the population.

A population is defined as the complete set of Objects of interest.

A Subset of a population usually chosen in such a way that it can be taken to represent the population. with respect to Some characteristics is called a Sample.

In many Situations. We collect data from a Sample. and generalise about the population based on this Sample data. Thus, Statistics refers to all the methods that are employed for collection, classification, presentation, analysis, interpretation and drawing conclusion from the Sample data.

Statistics is divided into two branches.

(1) Descriptive Statistics

(11) Inferential Statistics

Descriptive Statistics refers to the methods for Organising, displaying, and describing the data using tables, graphs, and summary measures.

it is the one that is generally used by common people in day to day life.

Deriving results about the population. based on the Sample data is called inferential Statistics

Here inference means a conclusion about a population based on logical reasoning from data collected from a Sample. in Other words. Descriptive Statistics. quantitatively describes a characteristics of the Sample data. Set. where as Inferential Statistics makes predictions or generalises about the population. based on the Samples.

Descriptive Statistics comprise the following measures.

Tabular representation :- table consisting of frequencies and. Sums.

eg:- absolute frequency, relative frequency, cumilative frequency, & Cross - tabulation

Graphical representation :- Graphs describing the data set.

eg:- Pie chart, bar chart, line chart, histogram. , box plot.

* Measure of Central Tendency :- which describe where a large part of the Sample is located.

    Eg:- Arithematic mean, geometric mean, median, [average] mode, etc.

* Measures of Dispersion :- Which describe how much is the spread of the Observations from the Central Value.

    Eg:- Standard deviation, Variance, Range.

## Levels of Measurements.

The statistical tools that can be applied for a given problem depend upon the level of measurement of the variables.

Level of measurement is a way of classifing the data values in terms of How they can be measured and compared to the other values.

The variables on which data is Collected can be classified into two namely Qualitative and Quantitative variables.

Qualitative variables are the Ones which cannot be measured in numbers.

Quantitative variables are the Ones which can be measured in numbers.

There are two levels of measurements for qualitative variables, nominal and ordinal. also there are two levels of measurement for Quantitative variable, interval and ratio.

Understanding these four levels of measurement is very important Since the level of measurement of the variable plays a keyrole. in deciding which statistical tool can apply to a given data.

## Nominal level.

Nominal variables are where the categories are just named.

The word Nomine means names. The variables falling into this category have no quantitative value.

In the nominal level only the name matters and order does not matter

eg:- eye colour, Gender.

## Ordinal Level.

This level is same as the nominal level. The only difference is that the order matters.

for eg:- assigning codes to household according to their levels of income. Say 0 to law, 1 to medium & 2 to high.

## Interval level of measurement

In this level of measurement both order and difference between the values of a variable matter but ratio does not matter.

for eg:-

City A has 10°c temperature and City B has 20°c temp. here we can interpret that City B is warmer than City B is warmer than city A. But it would be wrong to conclude that city B is twice as city A. also note that zero is an arbitary number in the interval level. ie it does not mean absence of that variable which means. if the temperature of a place is 0°c. then we cannot say that place has no temperature.

## Ratio or Level of Measurement

In this level of measurement name, order difference and ratio all are meaningful.

for eg:-

Suppose a box $B_1$ has 5 kg weight and box $B_2$ has 15 kg weight. Then we can say that box $B_2$ is thrice as heavy as box $B_2$.

## Descriptive Statistics

Quantitative data exhibit certain general characteristics in the following ways.

(1) They show a tendancy to concentrate at certain values usually some where in the centre of the distribution. measures of this tendancy are called measure of central tendancy or averages.

(2) The data vary about the measure. central tendancy which are called measures of dispersion or variation.

(3) The data in a frequency distribution may fall into symmetrical or assymmetrical. The measure of direction and degree of symmetrical or assymetry are known as measure of ~~distribution~~ Skewness.

4) The data in a frequency distribution exhibit flatness or peakdness of the frequency curve. This measure are known as measure of kurtosis.

## Frequency distribution

When observations are available on a single variable of a large no of individuals often it becomes necessary to summarise the data as far as possible:

To this end, a tabular representation which shows the distribution of the frequency in different classes may be used. The manner in the class frequency are distributed over the class intervals is called the grouped frequency distribution

eg:- Let us consider the marks in mathematics obtained by 50 students selected at random from a certain school. the following frequency table is obtained by computing the class intervals with corresponding frequencies.

| Class | 20-40 | 40-60 | 60-80 | 80-100 |
|-------|-------|-------|-------|--------|
| frequency | 4 | 20 | 20 | 6 |

The following points may be kept in mind for classification.

) The classes should be clearly defined and should not leave any ambiguity.

) the classes should be exhaustive, ie each of the given value should be included in one of the

(3) The classes should be mutually exclusive and non overlapping.

(4) The class should be of equal width, if the classes are of varying width. The different class frequencies will not be comparable.

(5) indeterminate classes: eg:- Open end class like $< a$; or $> b$. should be avoided as far as possible.

(6) The no of classes should neither be too large nor too small. The following formula may be used to determine an approximate no of class.

$$K = 1 + 3.322 \log_{10} N.$$

$N$ is the total frequency.

This equation is also known as Struges formulae.

## Class limit.

class limit should be chosen in such a way. that the midvalue of the class interval and actual average of the observation in that class interval are as near to eachother as possible.

## Continious frequency distribution

if we have a continiou variable, it is
not possible to arrange the data in the
class intervals of above type

for eg:- let us consider the distribution of age in years
if class intervals are ..15-19, 20-24, . . .
then the person with ages between 19 & 20
years are not taken into considiration.
In such a way that all the person with
any fraction of age are included in one
group. we may rewrite the above classes.
as shown in the following table.

| 0 - 5 | |
|-------|---|
| 5 - 10 | |
| 10 - 15 | |
| 15 - 20 | |
| 20 - 25 | |

The frequency distribution with such classes is known as Continous frequency distribution. In the above classes the upperlimit of each class are excluded from the respective classes and are included in the immediate next class.

Such classes are known as exclusive classes.

## Graphical representation of frequency distribution

it is often useful to represent a frequency distribution by means of a diagram. This representation facilitates the comparision of two or more frequency distributions.

## Histogram

To draw a histogram of a given continous frequency distribution. We first mark the class intervals in the $x$ axis.

In each class interval draw rectangles with height proportional to the frequency of the corresponding class interval. So that the area of the rectangle is proportional to the frequency of the class. Such diagrams of continous rectangles is called.

Histog
Remar
• if e
contin
distr
• To
dist
cor
ind

eg:

# Histograms

## Remarks

- If the grouped frequency distribution is not continuous, it should be converted to continuous distribution and then draw the histogram.

- To draw the histogram for ungrouped frequency distribution. We have to assume that the frequency corresponding to the value $x$ is spread over the interval $\left(x - \dfrac{h}{2}, \quad x + \dfrac{h}{2}\right)$

eg:- construct histogram for the following frequency table.

| Mark | No of Students |
|------|----------------|
| 15 - 19 | 9 |
| 20 - 24 | 11 |
| 25 - 29 | 10 |
| 30 - 34 | 44 |
| 35 - 39 | 45 |
| 40 - 44 | 54 |
| 45 - 49 | 37 |
| 50 - 54 | 26 |
| 55 - 59 | 8 |
| 60 - 64 | 5 |
| 65 - 69 | 1 |
| | 250 |

This frequency distribution is not continuous. So we first convert it into a continuous distribution with exclusive type classes.

Let $d$ be the gap between the upper limit of any class and lower limit of the succeeding class. then the new class boundaries are given by upperclass boundary = upper class unit + $\frac{d}{2}$.

lowerclass boundary = lower class unit $-\frac{d}{2}$.

So in our problem $d = 1$.

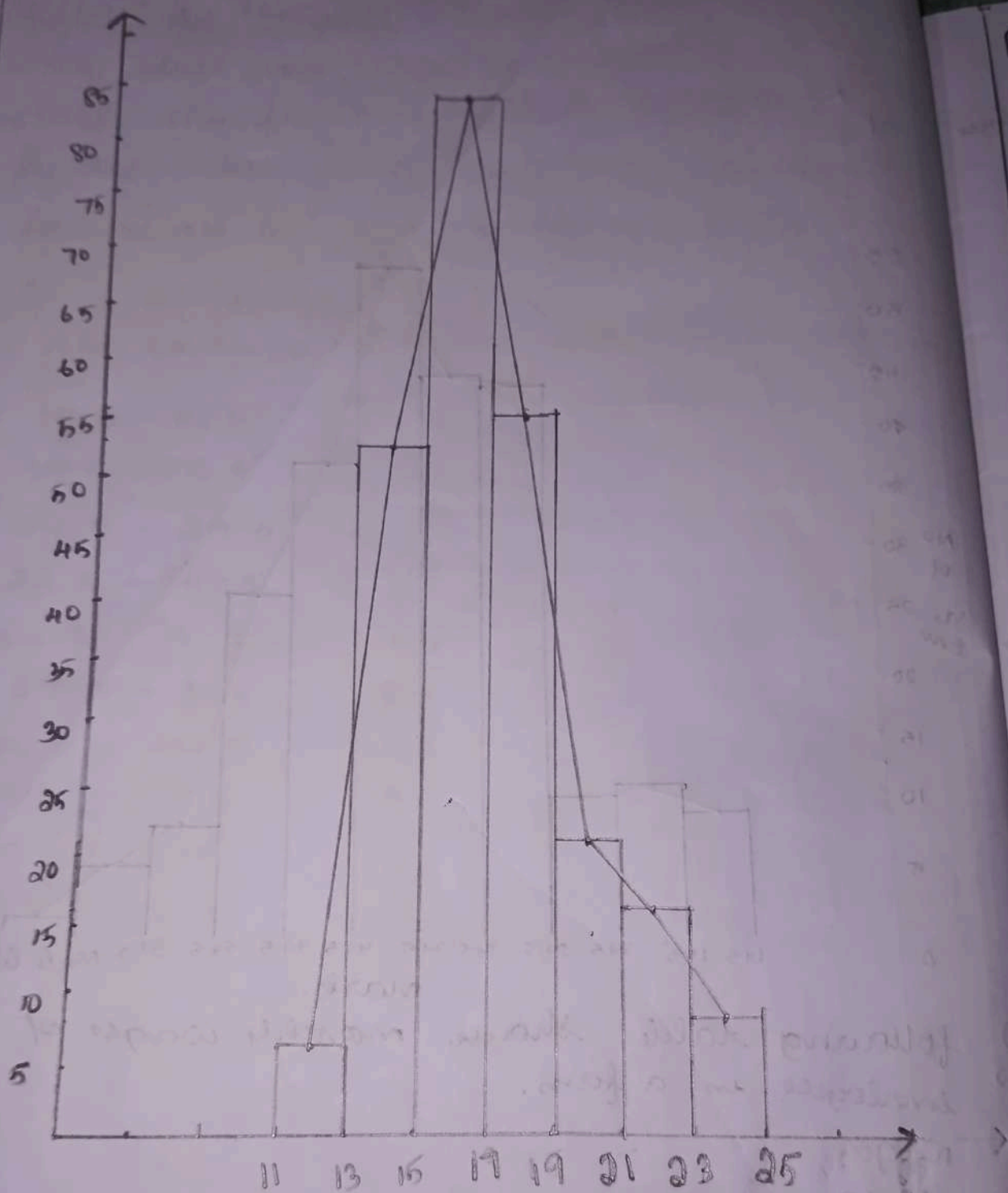∴ the continuous frequency distribution is given by

| | |
|---|---|
| 14.5 – 19.5 | 9 |
| 19.5 – 24.5 | 11 |
| 24.5 – 29.5 | 10 |
| 29.5 – 34.5 | 44 |
| 34.5 – 39.5 | 45 |
| 39.5 – 44.5 | 54 |
| 44.5 – 49.5 | 37 |
| 49.5 – 54.5 | 26 |
| 54.5 – 59.5 | 8 |
| 59.5 – 64.5 | 5 |
| 64.5 – 69.5 | 1 |

mark.

Q) following table shows monthly wages of employees in a firm.

wages

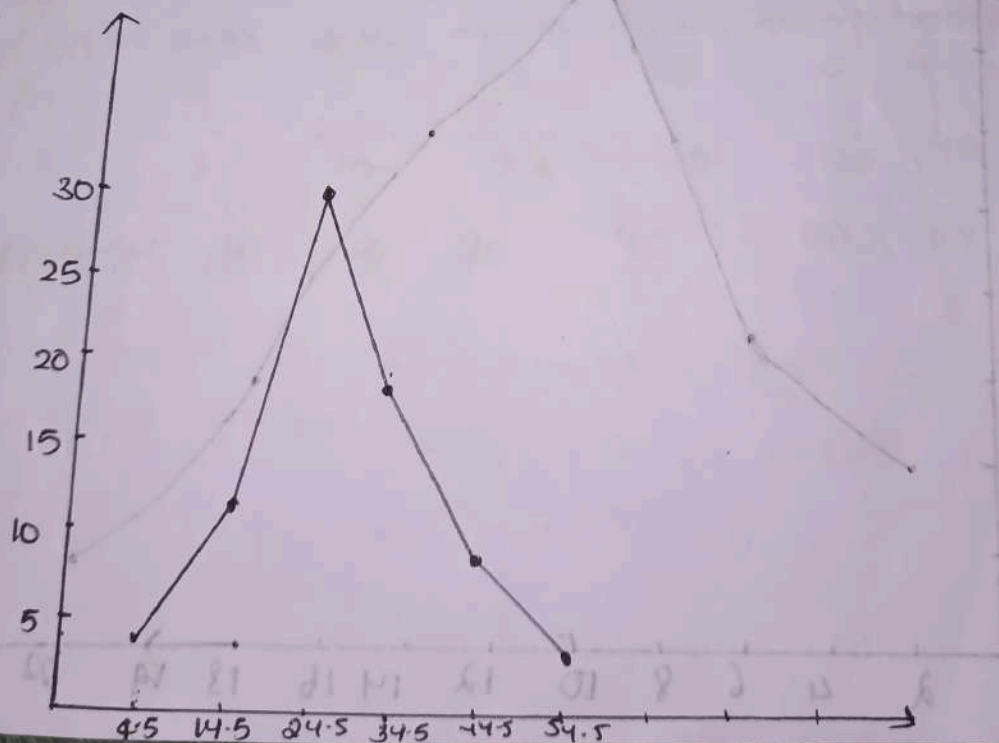| 11-13 | 13-15 | 15-17 | 17-19 | 19-21 | 21-23 | 23-25 |
|-------|-------|-------|-------|-------|-------|-------|
| 6 | 53 | 85 | 56 | 21 | 16 | 8 |

# Frequency Polygon.

In order to draw a frequency polygon plot in the x-y plain in such a way that the x axis values are equal to the mid values of the classes and Y axis values are equal to the corresponding frequencies. Then joint the adjacent points with straight lines which results a frequency polygon.

eg:- Represent the following frequency table by a frequency polygon.

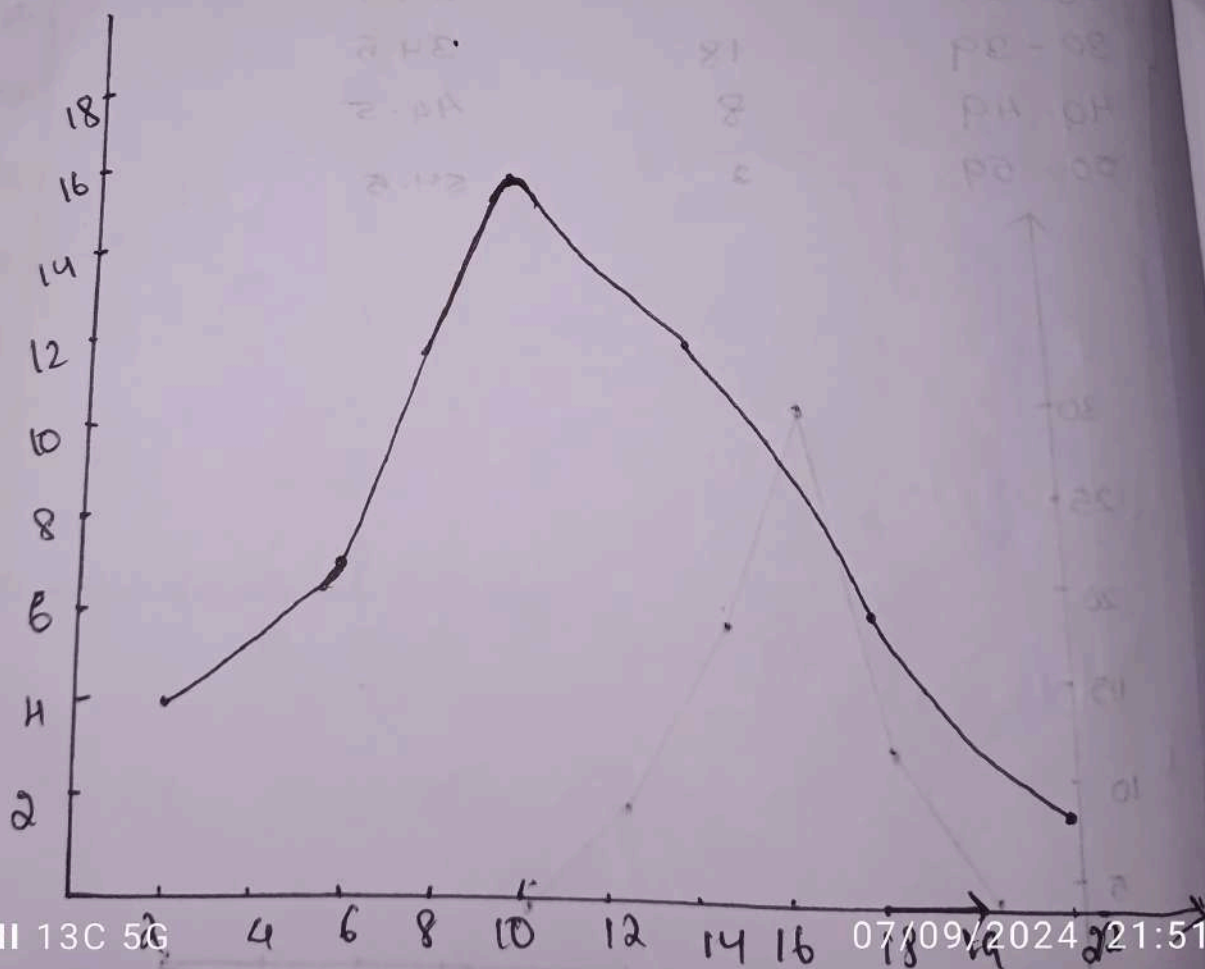| Class | frequency | mid value |
|-------|-----------|-----------|
| 0-9 | 4 | 4.5 |
| 10-19 | 12 | 14.5 |
| 20-29 | 30 | 24.5 |
| 30-39 | 18 | 34.5 |
| 40-49 | 8 | 44.5 |
| 50-59 | 2 | 54.5 |

## Frequency Curve

To draw a frequency curve from a frequency table. take the midvalues of classes on the x axis and the frequencies on Y axis.

Mark the points accordingly. next joint these plotted points with a smooth line, we get a frequency curve.

Eg:- draw a frequency curve from the following frequency table.

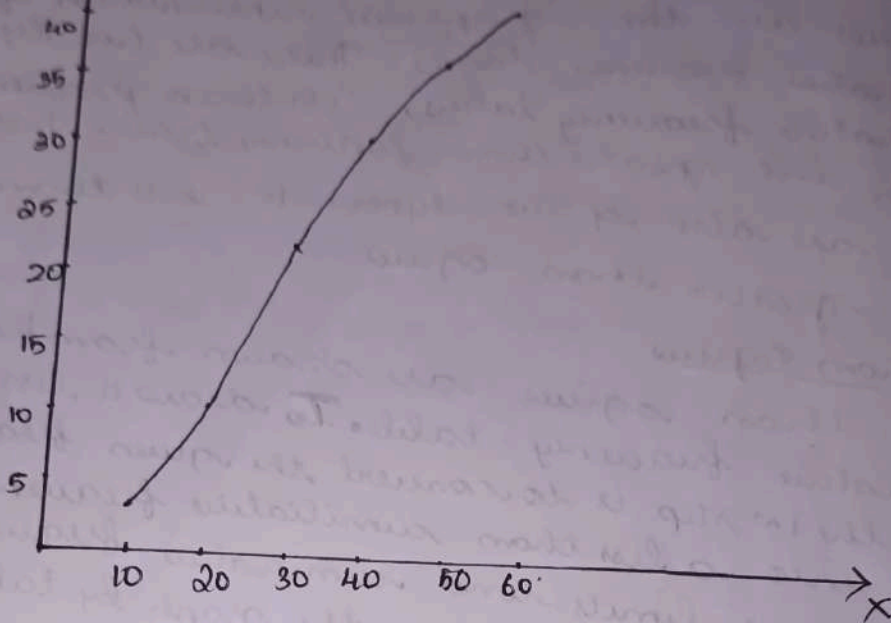| class | frequency | |
|-------|-----------|---|
| 0-4 | 4 | 2 |
| 4-8 | 7 | 6 |
| 8-12 | 16 | 10 |
| 12-16 | 12 | 14 |
| 16-20 | 6 | 18 |
| 20-24 | 2 | 22 |

# Ogives (Cumulative Frequency Curve)

Ogives are the graphical representation of cumulative frequency tables. There are two types of cumulative frequency tables, less than frequency table. less th and greater than frequency table. There fore. Ogives are also of two types. ie. less than ogives and greater than ogive.

## Less than Ogives

Less than ogives are drawn from less than cumulative frequency table. To draw a less than Ogive. the 1st step is to convert the given frequency table into a less than cumulative frequency table. with upper limits. and cumulative frequencies. Then plot the points on the graph by taking upper limits on the x axis and c.f on. Y axis. joining these points with a smooth line we get a less than ogive. Less than ogives are upward sloping curves from left to right.

eg:- draw Lessthan ogive curve from following table.

| Class | 0-10 | 10-20 | 20-30 | 30-40 | 40-50 | 50-60 |
|---|---|---|---|---|---|---|
| fre. | 3 | 7 | 12 | 8 | 6 | 4. |
| C.f | 3 | 10 | 22 | 30 | 36 | 40 |
| upperlims | 10 | 20 | 30 | 40 | 50 | 60 |

## Greater than Ogive

Graph of a Greater than Cumulative frequency table is called a greater than Ogive.
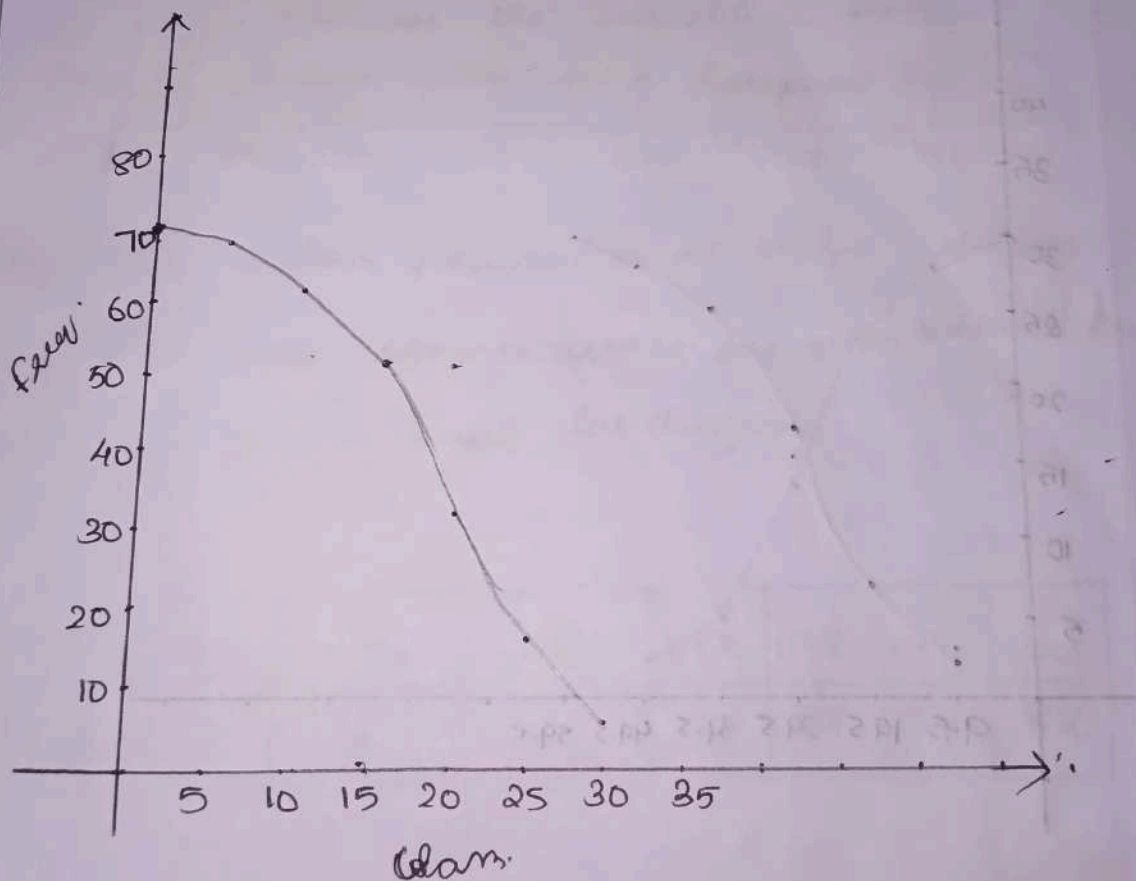
To draw a greater than ogive. marks the points on the graph in such a way that the x axis values are equal to the lower limits and y axis values are respective cumulative frequency. joining these points by a smooth line, we get a greater than Ogive. The greater than ogives are downward slopping curves. from left to right.

eg:-

represent the following frequency table in a more than ogive :

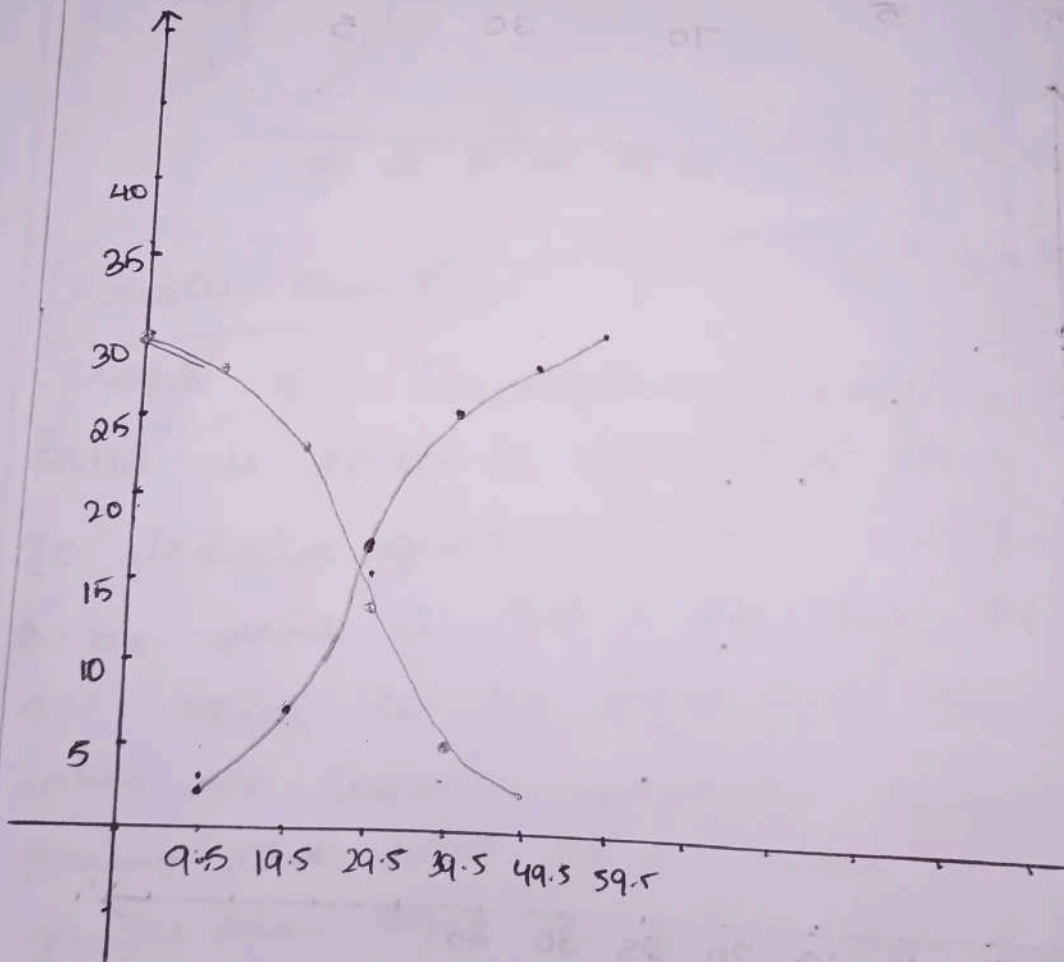| Class | Freq. | C.F. | lower | More than cf |
|-------|-------|------|-------|--------------|
| 0 – 5 | 2 | 2 | 0 | 70 |
| 5 – 10 | 5 | 7 | 5 | 68 |
| 10 – 15 | 12 | 19 | 10 | 63 |
| 15 – 20 | 20 | 39 | 15 | 51 |
| 20 – 25 | 16 | 55 | 20 | 31 |
| 25 – 30 | 10 | 65 | 25 | 15 |
| 30 – 35 | 5 | 70 | 30 | 5 |



**Remark**
it is possible to draw a less than ogive and a greater than ogive in the same graph.

Eg:- draw a less than ogive and a greater than ogive
on the same graph for the following data.

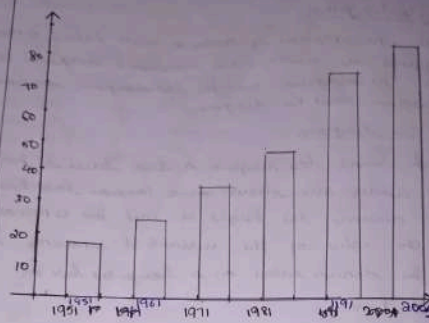| class | frequency | Class | L.S C.F. | G. C.F | Low | upper |
|-------|-----------|-------|----------|--------|------|-------|
| 0-9 | 2 | 0-9.5 | 2 | 30 | 0 | 9.5 |
| 10-19 | 5 | 9.5-19.5 | 7 | 28 | 9.5 | 19.5 |
| 20-29 | 10 | 19.5-29.5 | 17 | 23 | 19.5 | 29.5 |
| 30-39 | 8 | 29.5-39.5 | 25 | 13 | 29.5 | 39.5 |
| 40-49 | 3 | 39.5-49.5 | 28 | 5 | 39.5 | 49.5 |
| 50-59 | 2 | 49.5-59.5 | 30 | 2 | 49.5 | 59.5 |

An) Class.

## Type of diagrams

In the presentation of data a wide variety of diagrams diagrams are used. Some important diagrams are. Simple bar diagrams. multiple bar diagram, component bar diagrams and pie diagram.

## Simple bar diagram

In Simple bar diagram, a Suer Series of bars of equal width are drawn on a common base line of equal distances the height of each bar is proportional to the value of the variable it represents. it can be drawn either on a horizontal base or vertical base

eg) the urban population of india during the Cens. last six census years are given below. Present the data by a Simple bar diagram

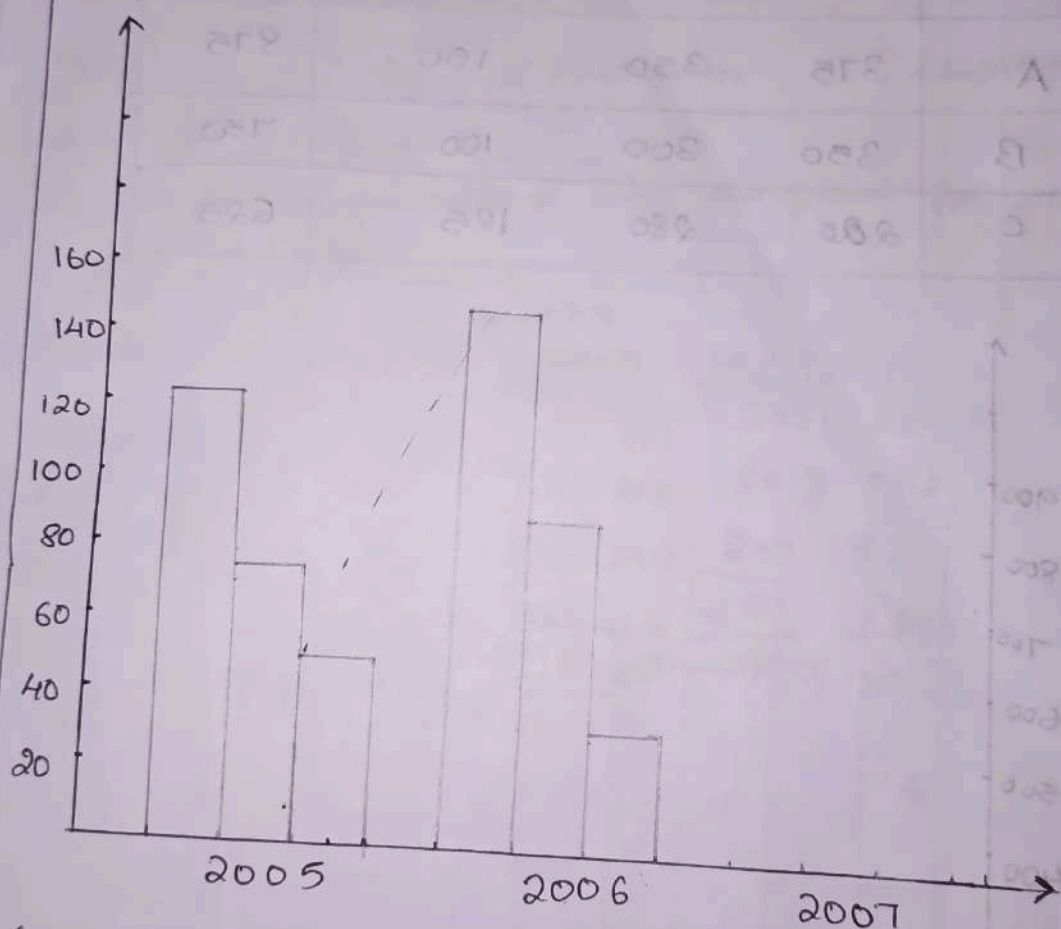| Census Year | 1951 | 1961 | 1971 | 1981 | 1991 | 2001 |
|---|---|---|---|---|---|---|
| Population | 18 | 26 | 35 | 48 | 76 | 82 |

## Multiple Bar diagrams

Multiple bar diagrams are used to compare two or more data sets.

In this type of diagrams bars representing diff values are drawn side by side and different sets of such bars are shown at equal distances. inorder to distinguish various bars different colours or shade are used.

eg:- Sales data of 3 companies a,b,c for 3 years are given below. represent the data by means of a multiple bar diagrams.

| Year | Sales in years | | |
|------|------|------|------|
| | A | B | C |
| 2005 | 122 | 75 | 50 |
| 2006 | 150 | 90 | 32 |
| 2007 | 130 | 81 | 49 |



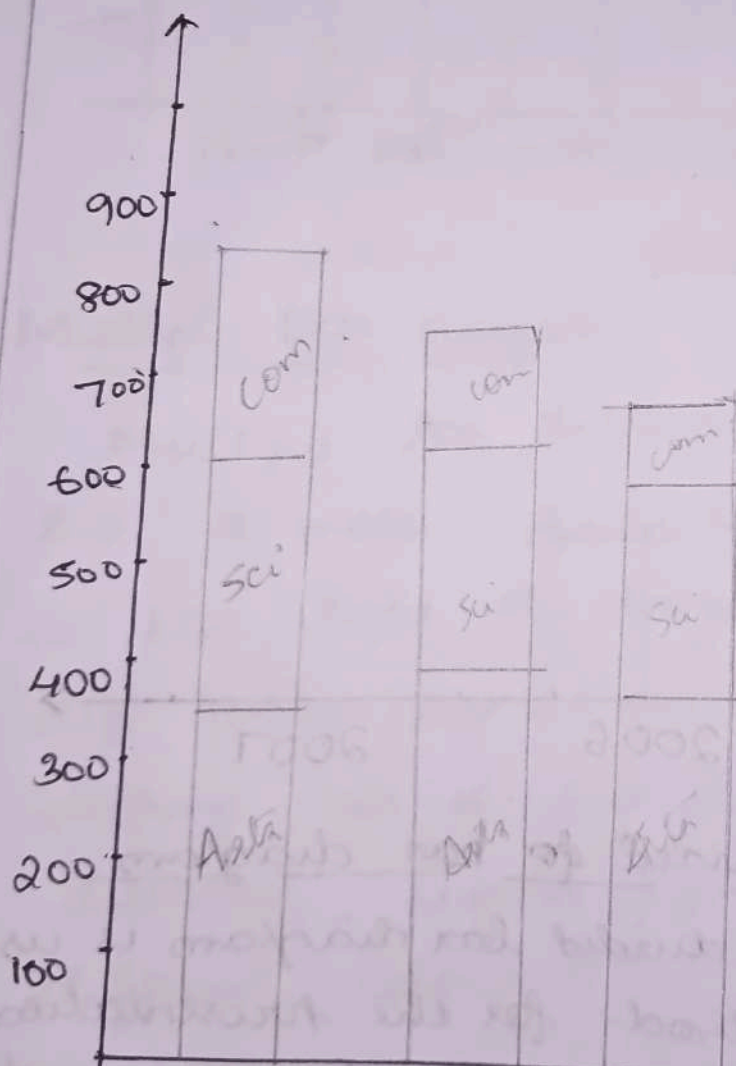## Component or Sub-divided bar diagram

A Component or Sub divided bar diagram is used as an effective method for the presentation of data. when the total value of t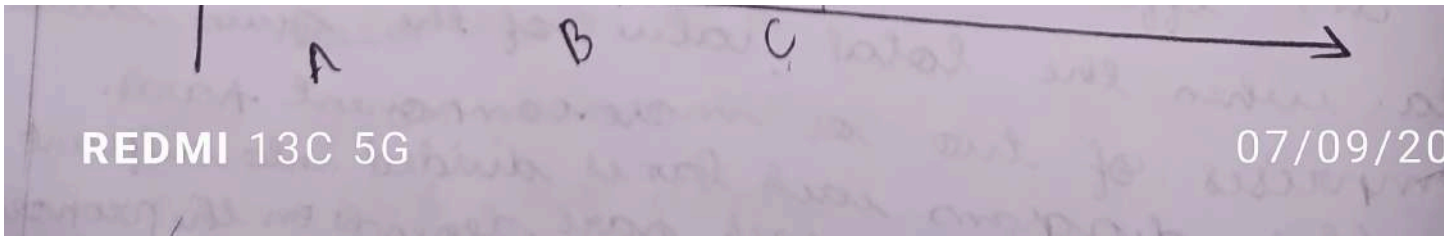he given data comprises of two or more component parts. in this diagram each bar is divided into different parts and size of each part depends on the proportional

Share of the component.

eg:- represent the following data by means of
component ~~mean ed~~ bar diagram.

| College | No of Students | | | Total |
|---------|------|--------|---------|-------|
| | arts | Science | Commerce | |
| A | 375 | 350 | 150 | 875 |
| B | 350 | 300 | 100 | 750 |
| C. | 220 | 280 | 125 | 625 |

A          B          C

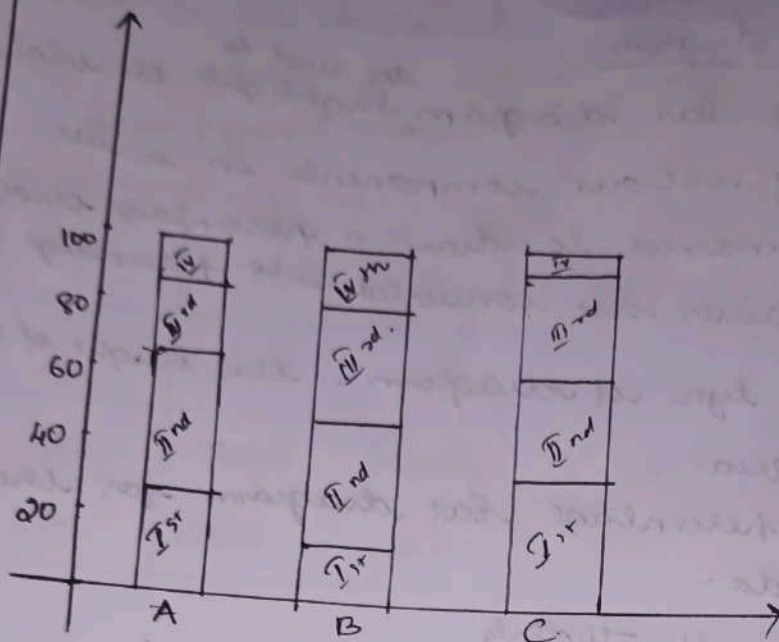## Percentage Bar diagram.

Percentage bar diagram are used to highlight the relative importance of various components in a bar diagram. inorder to draw a percentage diagram the given values are converted into percentage.

Thus in this type of diagram. the height of all bar will be equal.

eg: Draw a percentage bar diagram for the following data:

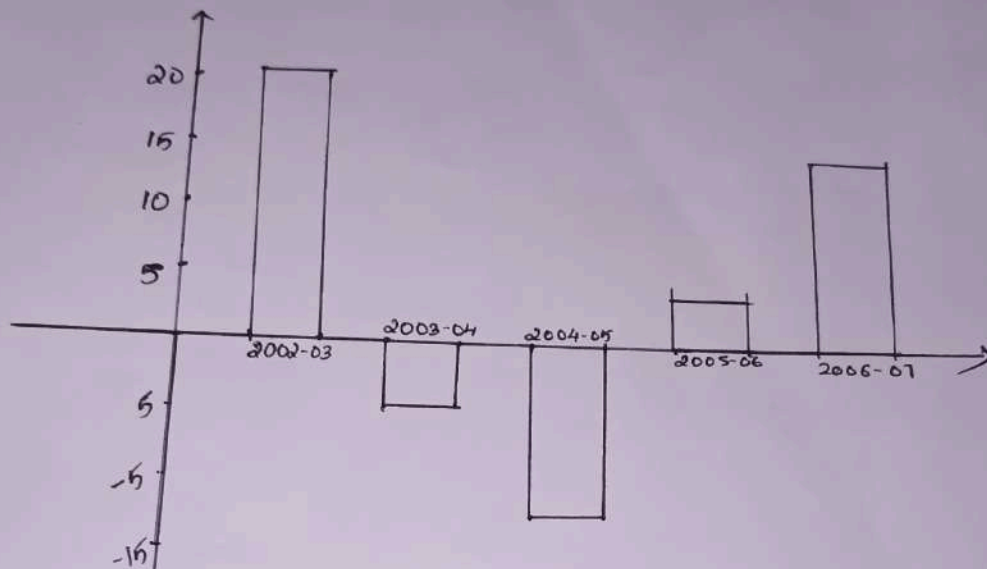| Result | NO of Students | | | | | |
|--------|----------|------|----------|------|----------|------|
| | college A % | | college B. % | | college C. % | |
| Ist | 250 | 29.3 | 243 | 15.3 | 342 | 38 |
| II | 318 | 37.2 | 530 | 33.5 | 258 | 28.6 |
| III | 210 | 24.6 | 600 | 37.9 | 250 | 27.7 |
| IV. | 75 | 8.85 | 208 | 13.15 | 60 | 6.666 |
| | 853 | 99.95 | 1581 | 99.85 | 900 | 100.9 |

## Deviation Bar Diagram:

Values of certain variables can be positive or negative. is Such cases deviation bar charts can be used to present the data.

A deviation chart is a particular type of bar diagram which shows both the +ve & -ve values like other bar diagrams it can be drawn either vertically or horizontally.

if bars are drawn vertically, negative values are marked below the base line. and if bars chart are drawn horizontally -ve values are marked in the left side of zero.

Eg:- Represent the following by a deviation chart

| Year. | |
|---|---|
| 2002-03 | 20 |
| 2003-04 | -5 |
| 2004-05 | -12. |
| 2005-06 | 4 |
| 2006-07. | 14. |



## Pie diagram

In a pie diagram total value of the given data is shown as the total area of a circle. This circle is then divided into different sectors in such a way that each sector represents a particular component. Inorder to design the area of such sectors, central angles proportionate to the corresponding values of the components are to be calculated. This is done by multiplying the ratio of each component to the total value. with 360.
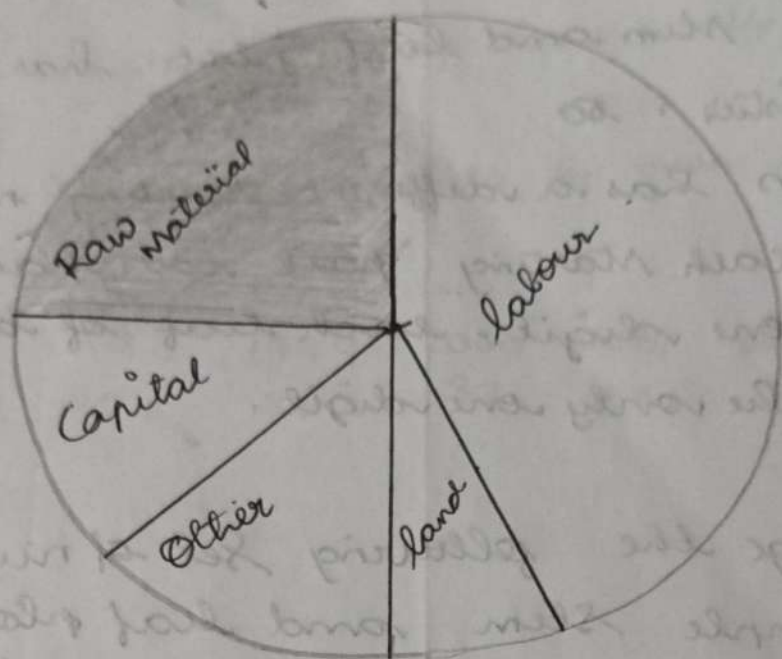
for eg:-

represent the following data using a pie diagram

| factor/inputs | cost in ₹ |
|---|---|
| Land | 44. |
| Labour. | 304. |
| Raw materials | 178. |
| Capital | 100. |
| Others | 94 |
| Total | 720. |

Solution.

The degree of angle representing each component is calculated as follows.

| Land | 44. | $\frac{44}{720} \times 360 = 22$ |
|---|---|---|
| Labour. | 304 | $\frac{304}{720} \times 360 = 152$ |
| Raw materials | 178 | $\frac{178}{720} \times 360 = 89.$ |
| Capital | 100 | $\frac{100}{720} \times 360 = 50$ |
| Other | 94. | $\frac{94}{720} \times 360 = 47$ |
| | 720. | |

# Stem And leaf plots.

Stem and leaf plots are techniques that allow rapid and informal exploration of the characteristics of a data d set. in a typical stem and leaf plot there is a vertical line of numbers called starting parts, and for each starting part there is a horizontal line of numbers called leaves.

Each complete horizontal line. ~~ie starting part plus~~ (ie starting part plus leaves) is called stem. every number in the data set being displayed has both a starting ~~part~~ ~~point~~ part and leaf.

The stem width determines which numbers in the data set are recorded on a given stem. a simple stem and leaf plot. has following characteristics : so

(1) each stem has a different starting part.

(2) while each starting part can have more than one digit, each leaf of a stem must be only one digit.

eg

eg:- Arrange the following set of numbers in a simple stem and leaf plot that has single digit starting parts & leafs. and stem width of 10.
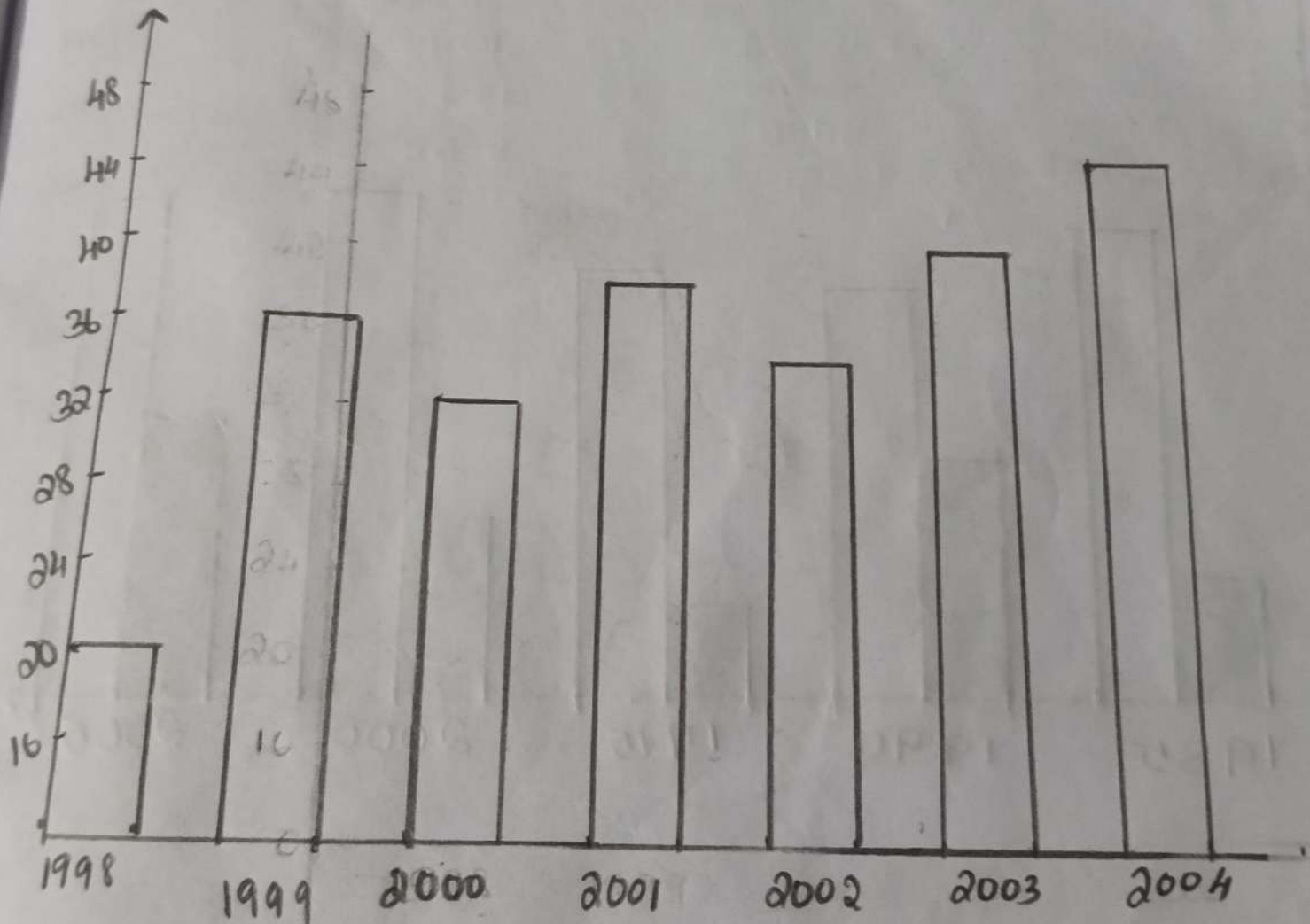
46, 35, 37, 20, 43, 15, 15, 26, 45, 25, 29, 13, 39, 44, 2?

24, 16, 40, 19, 45, 30, 34, 17, 39, 16, 40, 31, 21, 14,.

42, 16, 43, 22, 11, 24, 25, 31, 27, 40, 33.

| 1 | 5 5 3 6 9 7 6 4 6 1 | (10) |
|---|---|---|
| 2 | 0 6 5 9 1 4 1 2 4 5 7 | (11) |
| 3 | 5 7 9 0 4 9 1 1 3 | (9) |
| 4 | 6 3 5 4 0 5 0 2 3 0 | (10) |

starting parts are the vertial numbers (1, 2, 3, 4)
to the left of the vertial line and the leaves
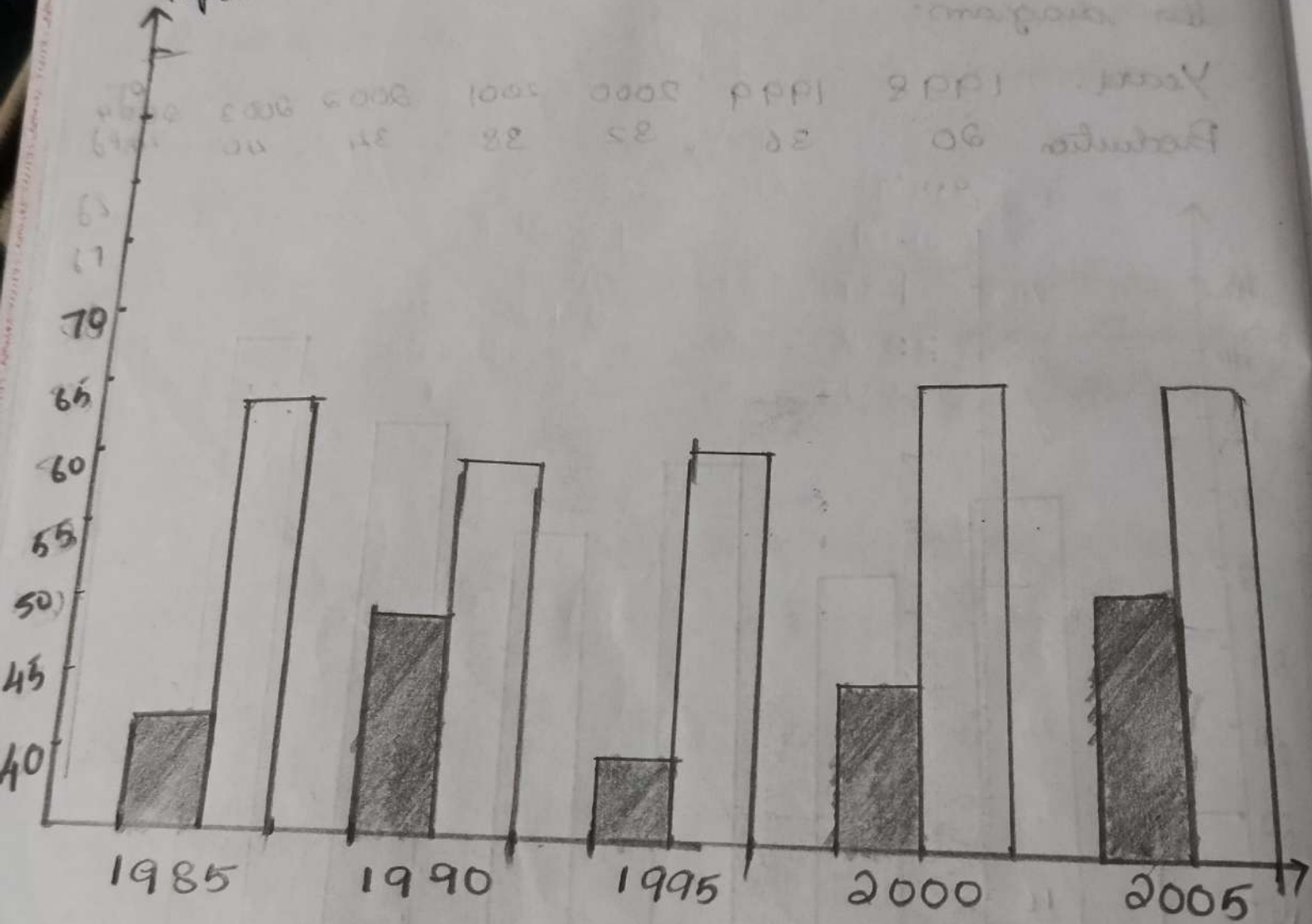are the numbers extending horizontally to the right.

Q). Present the following data by means of a Simple
bar diagram.

| Years | 1998 | 1999 | 2000 | 2001 | 2002 | 2003 | 2004 |
|---|---|---|---|---|---|---|---|
| Production | 20 | 36 | 32 | 38 | 34 | 40 | 45 |

2 draw a multiple bar diagram for the following data.

| year | 1985 | 1990 | 1995 | 2000 | 2005 |
|---|---|---|---|---|---|
| export | 42 | 48 | 40 | 45 | 52 |
| import | 64 | 60 | 62 | 70 | 65 |

3) following data shows the no of workers in 3 factories A,B C. draw a component bar diagram.

| Factory | No of workers. | | | Total |
| --- | --- | --- | --- | --- |
| | Male | female | childres | |
| A | 102 | 88 | 24 | 214. |
| B | 160 | 90 | 30 | 280. |
| C | 110 | 74 | 102 | 194. |

# Partition Values

partition values are the values which divide a frequency distribution into a no of equal parts.

the three points which divide the frequency distribution into 4 equal parts are called quartiles.

The first, second and third points are known as the 1st, 2nd and 3rd quartiles respectively.

The 1st quartile denoted by $Q_1$ is the value which exceed 25% of the Observations and is exceeded by 75% of the Observation

The 2nd quartile. $Q_2$ coincides with median the 3rd quartile $Q_3$ is the point which has .75% observations before it and .25% observations after it

The 9 points which divide the frequency distribution into 10 equal parts are called deciles.

whereas percentiles are the 99 points which divide the frequency distribution into 100 equal parts.

for eg:- The 7th decile, $D_7$ has 70% observation before i and 30% observation.

and 47th percentile. $P_{47}$ ~~~~~~~~~~~ is the point which exceeds 47% of the observation.

(class & frequency) ↓

$$Q_1 = L_1 + \frac{\left(\frac{N}{4} - m_1\right)}{f_1} \times C$$

L → Lower limit of $Q_1$ class.

$m_1$ → cumulative frequency just previous $Q_1$ class.

$f_1$ → frequency $Q_1$ class.

$c$ → class width

## Graphical Location of Partition Value:

The partition values can be located with the help of a cumulative frequency curve or ogive.

The procedure is given as follows.

①  first form the less than c.f. table. then plot the less than ogive curve. Similarly or plot the more than ogive curve.

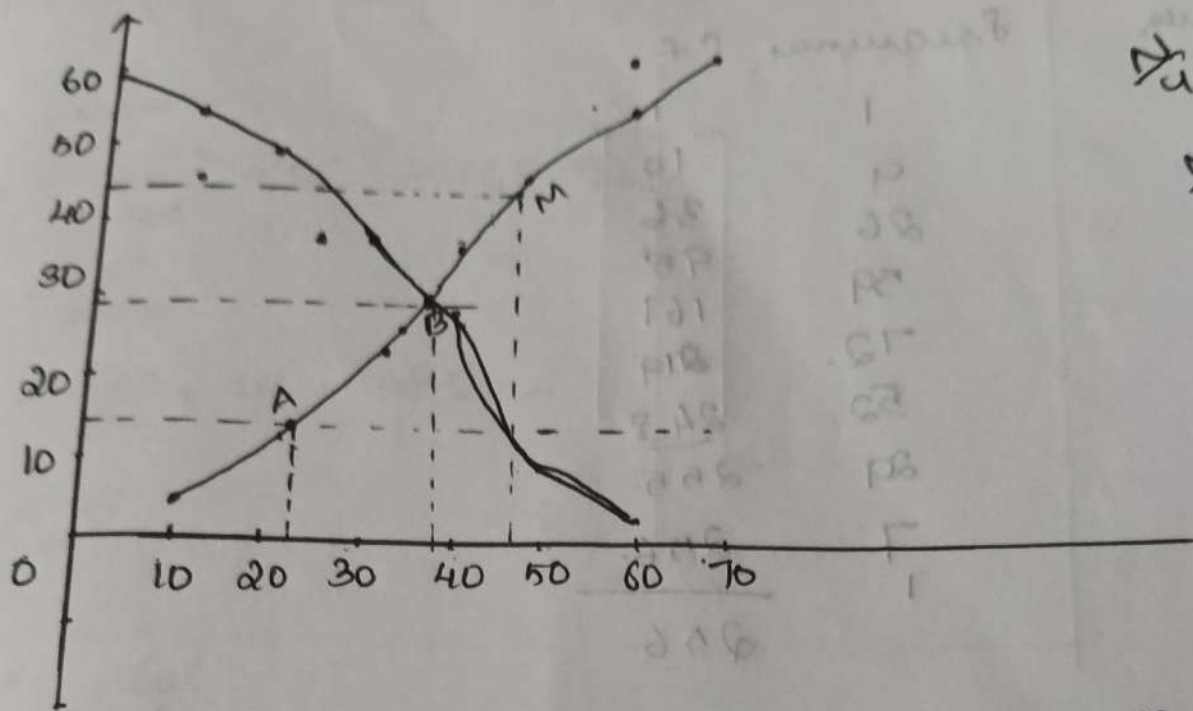To locate the value of $Q_2$. mark a point corresponding to $N/2$. where $N$ is the total freq... along $Y$ axis.

at this point draw a line parallel to $x$ axis meeting the ogive at the point A. draw a line from A perpendicular to $x$ axis meeting it in M. Then abscissa abscissa of M gives the value of median.

To locate the values of $Q_1$ and $Q_3$. mark the points corresponding to $N/4$. and $3N/4$.

— and proceed exactly as before.

eg:- locate $Q_1$, $Q_2$ & $Q_3$ for the following data

| Mark | No of std. | L. Than | M. than |
|------|-----------|---------|---------|
| 0-10 | 4 | 4 | 59 |
| 10-20 | 8 | 12 | 55 |
| 20-30 | 11 | 23 | 47 |
| 30-40 | 15 | 38 | 36 |
| 40-50 | 12 | 50 | 21 |
| 50-60 | 6 | 56 | 9 |
| 60-70 | 3 | 59 | 3 |

$\frac{3N}{4} = \frac{59}{4} = 44.$

$\frac{N}{2} = \frac{59}{2} = 29.5$

$\frac{3N}{4} = \frac{3 \times 59}{4}$

Q) 8 coins were tossed together and the no of hea resulting was noted. the experiment is repeate 256 times. and the frequencies that were ob obtained for different values are shown in the following tables.
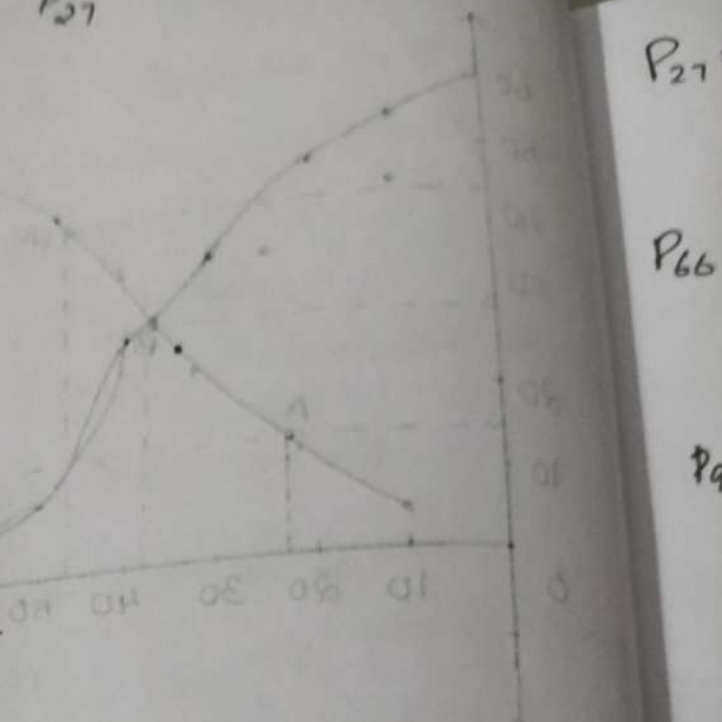
| Values | Frequencies |
|---|---|
| 0 | 1 |
| 1 | 9 |
| 2 | 26 |
| 3 | 59 |
| 4 | 72 |
| 5 | 52 |
| 6 | 29 |
| 7 | 7 |
| 8 | 1 |

calculate $Q_1$, $Q_3$, $D_4$, and $P_{27}$

| values | frequency | c.f. |
|--------|-----------|------|
| 0 | 1 | 10 |
| 1 | 9 | 36 |
| 2 | 26 | 95 |
| 3 | 59 | 167 |
| 4 | 72 | 219 |
| 5 | 52 | 248 |
| 6 | 29 | 255 |
| 7 | 7 | 256 |
| 8 | 1 | |
| | 256 | |

right margin labels
$P_{27}$

$P_{66}$

$P_9$

$$\frac{N}{4} = \frac{256}{4} = 64.$$

The Cumilative frequency. just greater than 64 is 95.

$Q_1$ = Value . Corresponding to the frequency 95.

$= 3$.

$Q_3$ . $\frac{3N}{4} = 192$

$Q_3 = 5$.

value corresponding
$Q_2 = 4$ to the $\frac{N}{2}$ th freq
$= 4$.

ⓓ. to find $D_4$.

$\frac{4}{10}N = \frac{4}{10} \times 256. = 102.4$.

$P_{27} = \frac{27}{100} \times 256 = 69.12$ .

$\underline{\underline{= 3}}$ .

$P_{66} = \frac{66}{100} \times 256. = 168.96$

$\underline{\underline{= 5}}$ .

$P_{90} = \frac{90}{100} \times 216 = 230.4$ .

$\underline{\underline{= 6}}$ .

## Box - plot.

box plot is a method for demonstrating graphically the spread of numerical data through their quartiles. in addition to the box on a box plot. there can be lines called whiskers. Hence the plot is sometimes called Box and whisker plot.

A box plot is a standarised way of displaying the data set based on the 5 number Summary; minimum, maximum, median, first quartile and 3rd quartile.

In addition to the minimum and maximum values, another important element can also be used to obtain a box plot is the inter quartile range. $IQR = Q_3 - Q_1$, A box is drawn from $Q_1$ to $Q_3$ with a horizontal line drawn inside to denote the median.