

Stat 101

* Nature of Statistic its definition, Importance & Limitation

Nature (origin and development) of statistic, in a sense, is as old as the human society itself. It has been in existence from the time of life on this earth. Its origin can be traced to the old days when it was regarded as the science of statecraft and was the by-product of administrative activity of the state. The word "statistic" seems to have been derived from the Latin word "status" or the Italian word "statista" or the German word "statistik" each of which means a "political state".

In ancient times, the government used to collect this information regarding the population "property or wealth" of the country. The former enabling the government to have an idea of the man power of the country and the latter providing its basis for introducing new taxes and levies. It is therefore very important to start a course of this kind with an explanation of the purpose and concept of statistical science.

Statistic is basically a science of

Collecting, organizing, analysing, presenting and interpreting **data**. These five purposes are intertwined and any one is usually controlled by another to a large extent, but until the beginning of the 19th century far more attention was paid to the collection and presentation of data than to their interpretation. Large volume of data were usually collected and frequently misinterpreted if indeed interpretation was attempted.

However, since that time the importance of a scientific approach in the interpretation of data has been realized and great steps have been made in the development of appropriate methods. Correspondingly, methods of collecting and presentation of data have been altered to keep pace with the new methods of analysis and interpretation especially in this era of Information and Communication Technology where computers formed the center stage of data analysis until recently. The interpretation of data hold a central position in statistic and method of collection and presentation are hinged J method

of Interpretation.

The growth of statistic especially in the methods of data interpretation during the William Sealy, Karl Pearson, Francis Yates, Charles Edward Spearman past 60y may be linked with the names of and in particular Sir Ronald Aylmer Fisher. This growth started in the field of biological research, but the nature of the problems in Counter has caused the new method to be applied to medical, physiological and economic data and to some extend in physics and engineering. For example the biologist in his routine work is confronted with the difficulty that the measurement after identical treatment, of two animals or plants apparently similar in all assets, can give widely different result.

Again the background, statistic has been defined as the science which deals with the collection, organisation, presentation, analysing and interpretation of numerical data with aim of drawing a logical conclusion. Statistic is the most verifiable tool for planning policy formulation and decision making.

every life, in every sector, decision are taken and plans are made but the success of every plan depends on the accuracy, analysis and interpretation of the available data. This shows that statistic placed a vital role in all planning and decision making process.

Typically, there are two main branches of statistic which are:

Descriptive statistic and Inferential statistic both of these are employed in scientific analysis of data and both are equally important.

* Descriptive statistic - these simply deals with the collection and presentation of data which is usually the first part of statistical analysis.

* Inferential statistics - as the name suggest this involve drawing the right conclusion from the statistical analysis that has been performed using descriptive statistic.

IMPORTANT OF STATISTIC

In modern times, statistic is the back bone of the decision making and it been introduced in to all sector of life, medical, Industrial, economics, political, among others.

statistic is no longer consider as a tool mere device for collecting numerical data but as a tool for developing sound techniques for data handling, analysis and making viable Inference for them, some of the importance and uses of statistic are as follows :-

1. statistic and planning - statistic, is indispensable to planning all over the world government especially of the developing economic are taking for planning for economic and technological aspect.
2. statistic and economics - statistic and techniques of statistical analysis have proved immensely important in solving the varieties of economics problems such as wages prices, time series analysis, in fact wide application of mathematics and statistic in the study of economic have lead to the development of new discipline called the economic statistic and econometrics
3. statistic and business - statistic is an indispensable tool in product marketing, inventory analysis and

production Control. Business executives are relying more and more on statistical techniques for studying the needs and desires of the consumer and for many other business purposes. The success of a business man depends so much on the accuracy and precision of his statistical forecasting.

4. Statistics and Industry - In the Industrial sector, statistic is widely used in quality control, in production engineering as to know whether the product is conform to specification or not. Statistical tools like Inspection plant, control chart, acceptance inspection and process control of extreme importance in the industry. Statistical data analysis has also proved useful in research and scientific field like biology, astronomy and medical science.

LIMITATIONS OF STATISTIC

- 1. If sufficient care is not taken in collecting, analysing and interpreting data, statistical might be misleading.
- 2. Statistical method are only based applicable to

quantity data and not suited for quality phenomenon.

3. statistic can not be applied to heterogeneous data.

4. statistic doesn't study individuals and therefore has no recognition to these individuals if any of a series are individual.

5. unlike the laws of physical and natural sciences, statistical laws are only approximation and not exact.

Types of statistical data

There are two general types of statistical data. These are:

1. Qualitative data

2. Quantitative data.

* Qualitative data - In certain investigation, we are only concern with the present or absent of some characteristics in a set of object or individual. For example if we have record of birth, we are interested only to know whether is male or female and count the number of male babies. Similarly, if a coin is

tossed a number of times we may only record a number of heads in a given set of tosses.

This type of data is called qualitative data. The characteristics used to classify individual into different categories is called attributes. Typically examples of attributes are eye colour, religion, gender, nationality.

* Quantitative data - when we are interested in a variable we either measure the magnitude of each of the individuals or unit under consideration. This type of data is called quantitative data. Example height of a person, the speed of a car, weight of a person, age etc.

The primary and secondary data

Data may be gathered from 2 main sources namely:

1. primary data - primary data are those data gathered directly from the source; that is, when data is collected directly from the observe (values for e.g., data collected from questionnaire).

Interviews, observations or experiment.

2. Secondary data - are those collected ~~data~~ gathered directly from the other source, but essentially not directly from observed values for example data collection through transcription; or documentation from existing record data collected from magazines, journals, library, archives among others.

Note

The major difference between primary and secondary data is that secondary data have to be extracted from data that was probably collected from a different purpose whereas primary data is collected purposely to fulfill user's practical needs.

METHOD OF DATA COLLECTION

Data are generally needed, collected and analysed to provide useful and meaningful information for planning and execution of a survey or experiment. It depends on the types of data needed which is great influence by the method through with the data are -

collected. The decision and choice of method of the data collection should be arrived at after considering the aims and objectives of the survey or experiment, the nature of information needed, the population and study the degree of accuracy desired, practicability of result, time and cost.

Thus, the method are discussed as follows

1. mail or postal Questionnaires method - in this method, prepared questionnaires on the subject matter of the survey are sent through an agent such as post office or email to the respondent requesting them back by posting online or directly mailing.
2. personal Interview method - this method requires the interviewer to ask prepared set of questions from the respondent and record the answer. This method allowed the interviewer and the respondent to meet face to face this method is widely use in journalism and population census.
3. Telephone interview method - with the modern telephone system, an interviewer can obtain

Information from respondent by making use of a telephone. Some decade ago, telephone interview was most common in advance country where there is good No of telephone substance and good telecommunication system. Now with the advent of mobile phones the telephone interview could be used in the development country.

4. Direct observation method - this is the method of data collection by physical observation or measurement of the unit, item or respondent under study. The researcher record the result using the appropriate measuring instruments.

5. Experimental method - In this method the desired data are collected through in the design of an and in the success or failure of this method depends on the skills of the experimenter as well as the quality of the instrument used for the experiment.

6. Method of registration - In this method, the respondent are required to registered the required

also the information at some designated places
and the vital statistic registration in many
countries provides and illustration of this
method in the life of a nation

7) Transcription or Documentation method -
this method is used when the data
is needed for a specified survey - data
already in registration files, document
among others are usually collected by
this method this is also method of collecting
secondary data in forms.

Graphical and Diagrammatical Representation
of data

Graphical representation of data is the most
convenient and popular way of describing
data is using graphical presentation.
It's easier to understand and interpret
data when they are present
graphically than using words or a
frequency table. a graph can present
data in a simple and clear way.
also it can illustrate the important

analytical aspect of the data - these lead to better analysis and presentation of the data.

Statistical graphs or chart serve two purposes: the first purpose is the presentation of statistical information on graphs and chart is to make the reader appreciate a simple, evocative display. The second main use is as a private aid to statistical analysis. The statistician will often have recourse to graphs and chart to gain insight into the structure of the data and to take assumptions which might be made in an analysis. Data can repr be presented in any of the following forms;

1. pie chart
2. Bar chart
3. Histogram
4. frequency polygon
5. Cummulative frequency Curve (ogive)

Pie chart

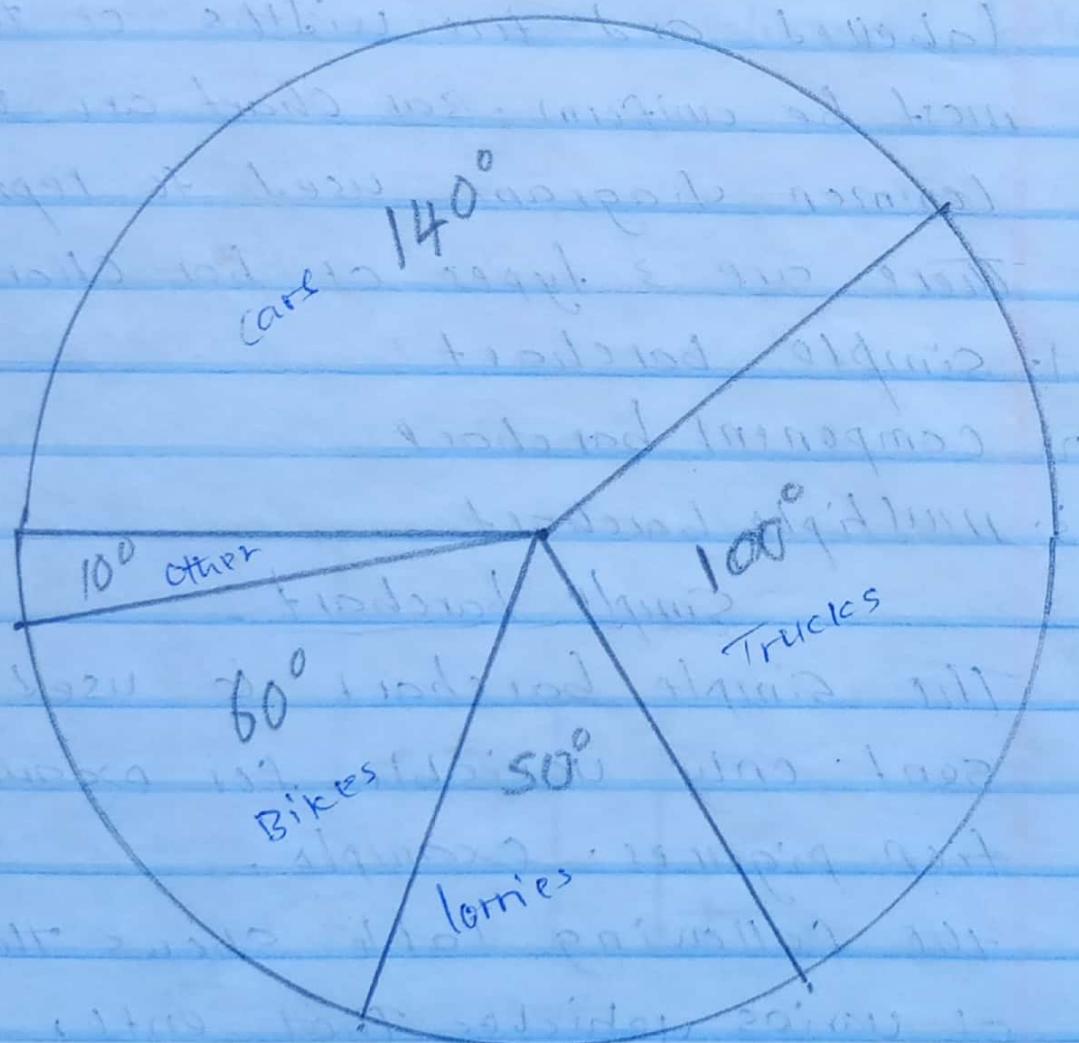
With pie chart data are represented which is divided into sectors with angles in each sector proportional to the frequency in that particular group. Angle subtended at the center of the circle is 360° the angle in each group is computed by the formula below:

$$\text{Angle} = \frac{\text{frequency}}{\text{total}} \times 360^\circ$$

The following table shows the distribution of various vehicles that enter a teaching hospital in a day. Represent the information in a pie chart.

Vehicle	number of vehicle
Cars	70
Trucks	50
Cotties	25
Bikes	30
(Others)	5
Total	180

<u>Vehicle</u>	<u>Frequency</u>	<u>Angle</u>
Cars	70	$\frac{70}{180} \times 360^\circ = 140^\circ$
Trucks	50	$\frac{50}{180} \times 360^\circ = 100^\circ$
Lorries	25	$\frac{25}{180} \times 360^\circ = 50^\circ$
Bikes	30	$\frac{30}{180} \times 360^\circ = 60^\circ$
Others	5	$\frac{5}{180} \times 360^\circ = 10^\circ$
Total	180	360°



pie chart

Bar chart

With bar chart data are represented with a series of equally spaced rectangles called bar. The bars are drawn either vertically or horizontally such that their height correspond or are proportional to the frequency in each group. For clarity, both axis of the chart must be properly labelled and the widths of the bars must be uniform. Bar chart are the most common diagram used to represent data. There are 3 types of bar chart namely;

1. simple barchart
2. component barchart
3. multiple barchart

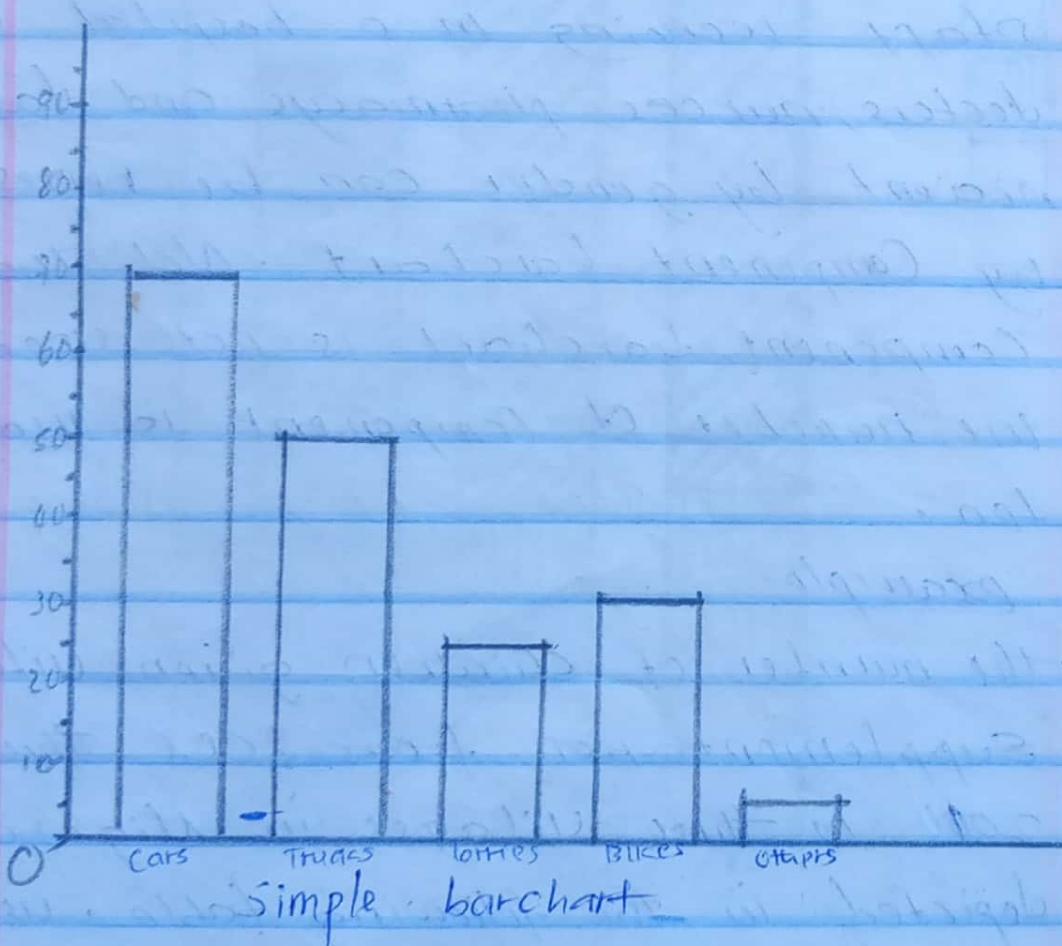
Simple barchart

The simple barchart is used to represent only variable, for example population figures. Example -

The following table shows the distribution of various vehicles that enter a teaching hospital in a day. Represent the information in a bar chart.

<u>Vehicles</u>	<u>Number of vehicles</u>
Cars	70
Trucks	50
Lorries	25
Bikes	30
others	5
total	180

The simple barchart in diagram below



Component bar chart

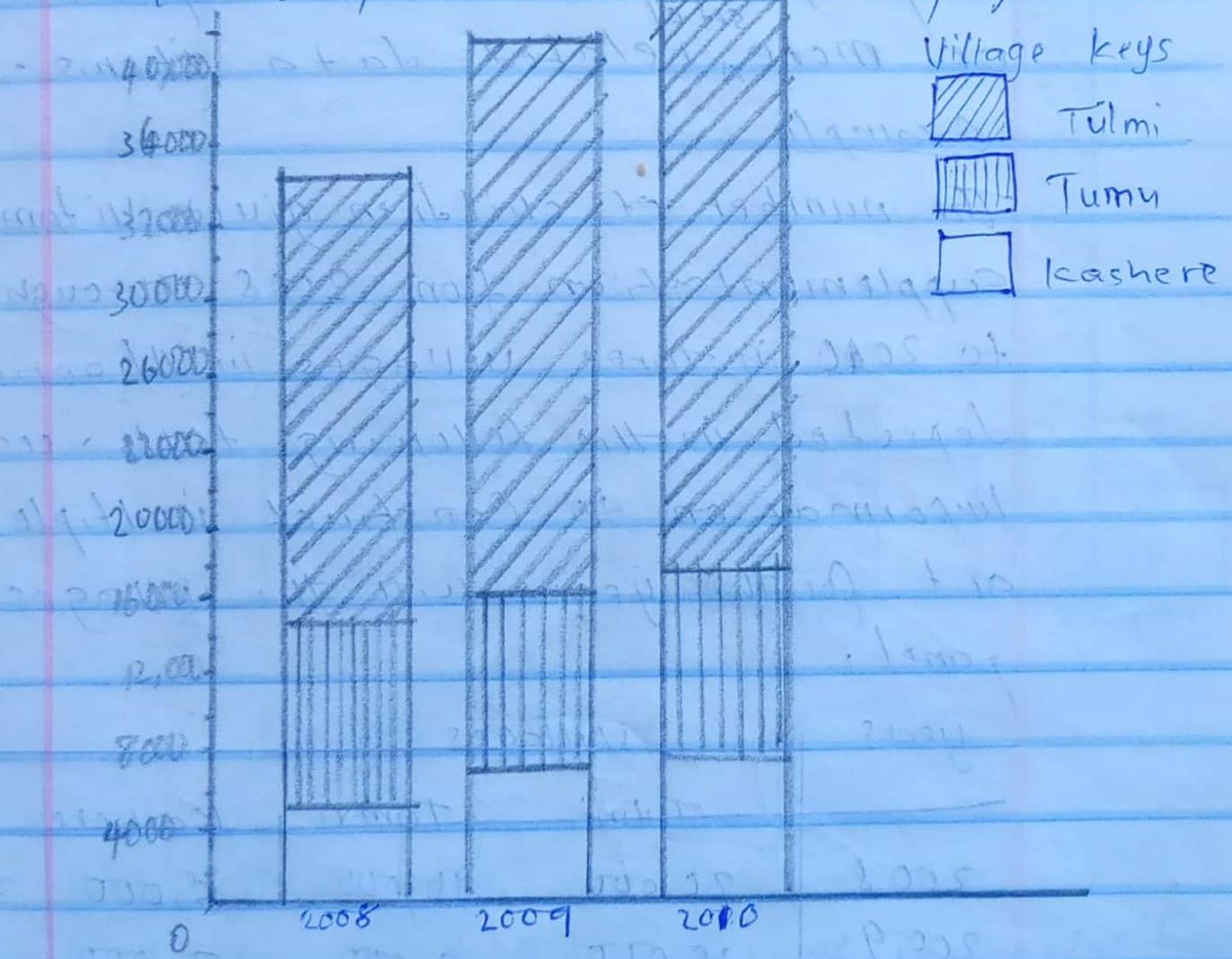
The component bar chart is used to represent two or more variable, for example population figures by gender where male and female population are displayed, it's used to represent the magnitude of a giving phenomena which is further subdivided into its various component. for example the number of professionals staff working in a hospital like doctors, nurses, pharmacists and labo-technician by gender can be represented by Component bar chart. Note that the Component bar chart is not used when the number of component is more than ten.

example

- the number of children given vitamin A supplementation from 2008 through to 2010 in three villages in Nigeria is depicted in the following table. use the information to construct Component bar chart for the years with the villages as panel.

years	Tulmi	Tumu	Kashere	total
2008	20,000	10,000	5,000	35,000
2009	28,000	9,000	7,000	44,000
2010	31,000	9,500	7,500	48,000
total	79,000	28,500	19,500	127,000

With respect to this component bar chart, is displayed below



Component bar chart

Multiple bar chart

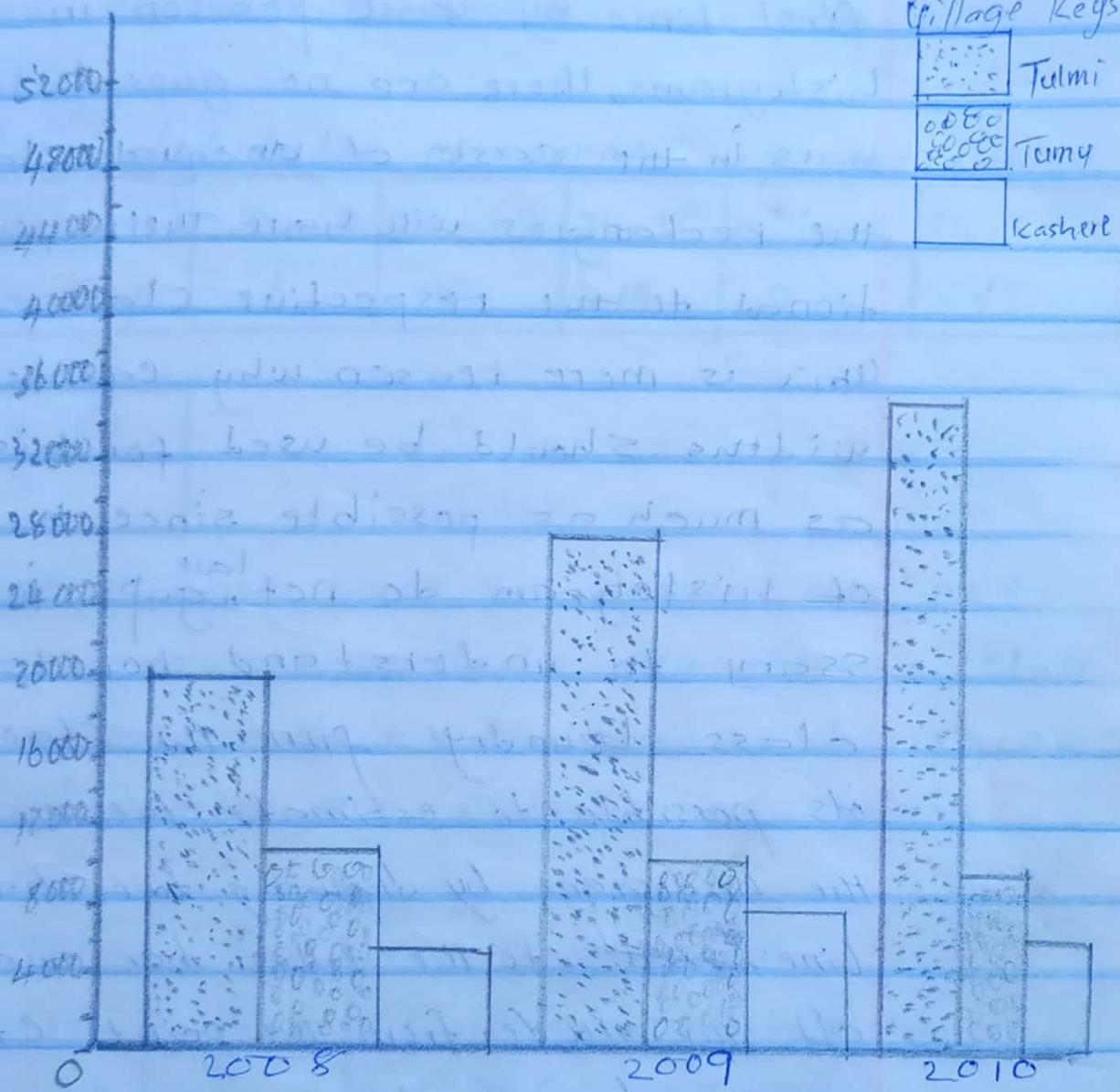
The multiple bar chart are also used to represent two or more variable where different bars side by side, with different shades or colours. For example population figures by gender where male and female population are displayed. It's use to represent two or more ^{set of} related data items.

Example

The number of children given Vitamin A supplementation from 2008 through the to 2010 in three villages in Nigeria is depicted in the following table. Used the information to construct multiple bar chart for the years with the villages as panel.

Years	Villages			Total
	Tulmi	Tumu	Kashere	
2008	20,000	10,000	5,000	35,000
2009	28,000	9,000	7,000	44,000
2010	31,000	9,500	7,500	48,000
Total	78,000	28,500	19,500	127,000

the multiple Bar chart displayed below



multiple bar chart

Histogram

A histogram is the series of rectangles having areas proportional to the frequency of each class. The term histogram was used for the

first time by Karl Pearson in 1895. In histograms, there are no gaps between the bars. In the case of unequal class width, the rectangles will have their width proportional to the respective class sizes. This is more reason why equal class widths should be used for histogram as much as possible since the bars of histogram do not have gaps, it is necessary to understand how to obtain class boundaries from the class limits. It's possible to estimate the mode from the histogram by drawing two diagonal lines ~~parallel~~ to the taller bar and reading off the mode from the point of intercept of the diagonal.

Example

An ophthalmologist from the National Eye Center observed the following data on the distribution of 50 people by age who suffer from eye problem in a particular village. Use the data to draw a histogram.

and estimate the modal age from your histogram.

Data distribution of people with eye problem

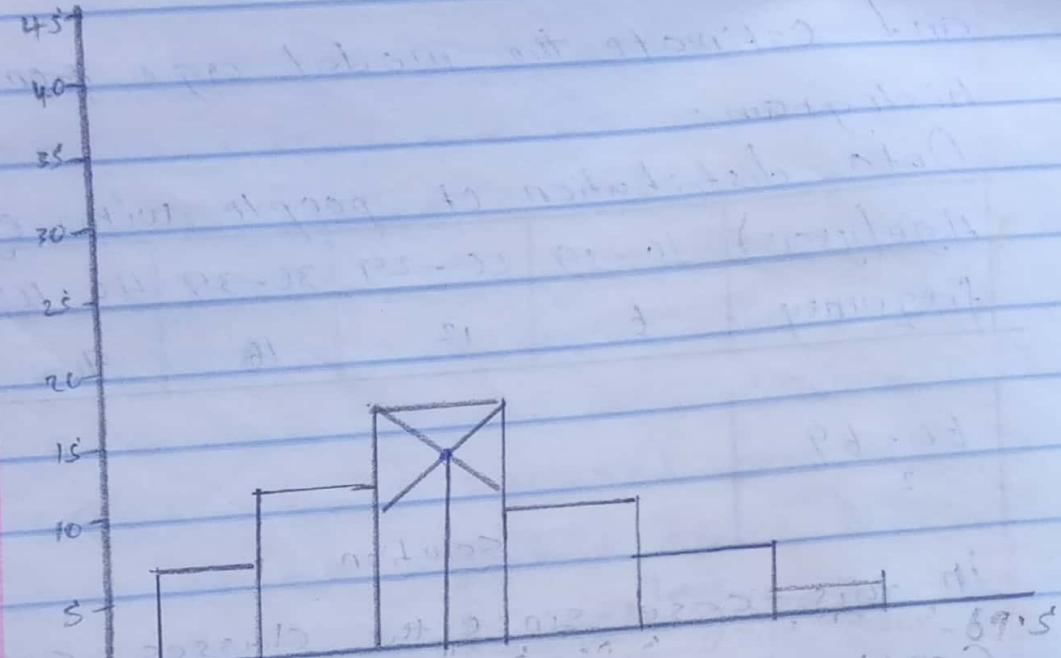
Age (years)	10 - 19	20 - 29	30 - 39	40 - 49	50 - 59
frequency	6	12	10	9	5

60 - 69	2
---------	---

Solution

In this case since the classes are not continuous, we need to compute the class boundaries before drawing the histogram. The boundaries are obtained as follows.

Age	f	class boundary
10 to 19	6	9.5 - 19.5
20 to 29	12	19.5 - 29.5
30 - 39	10	29.5 - 39.5
40 - 49	9	39.5 - 49.5
50 - 59	5	49.5 - 59.5
60 - 69	2	59.5 - 69.5



the modal age is 33.5 approximately as estimated from the histogram.

frequency polygon

A frequency polygon is a series of connected line segments in which the points of connection represent the values in the classes against their respective frequencies for group data, the points of connection are the mid point of each class against their respective frequency. Hence, the frequency polygon is a diagram drawn

(using mid point of histogram)

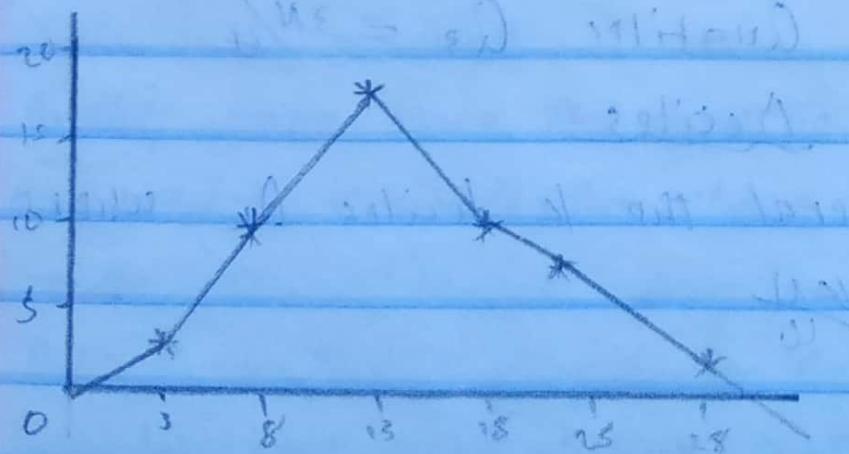
• by joining of the mid point of the tops of histogram bars with straight lines.

Example (bar chart)

An ophthalmologist from the national eye centre observed the following data on the distribution of 50 people by age who suffered from eye problem in a particular village we have the data to draw a) frequency polygon.

Age	frequency	No of patient mid point
1 - 5	3	3
6 - 10	10	8
11 - 15	18	13
16 - 20	10	18
21 - 25	7	23
26 - 30	2	28

the frequency polygon is displayed below.



$P = \frac{N}{100}$
 $P_k = \frac{kN}{100}$
 Deciles = $\frac{N}{10}$
 $D_8 = \frac{8N}{10}$
 $Q_1 = \frac{N}{4}$
 $Q_3 = \frac{3N}{4}$
 $Q_2 = \frac{2N}{4} = \frac{N}{2}$
Cummulative frequency Curve (ogive)
 When frequency are added successively, they are called cumulative frequency.
 The curve obtained by plotting the cumulative frequency against the upper class boundary is called the cumulative frequency curve(ogive) the curve can be used to estimate partition values such as Quartiles, Deciles, percentiles. The position of the Quartiles, deciles and percentiles are first located and subsequently estimated from the ogive those positions can be located by the following formulae.

Quartiles

Lower Quartiles $Q_1 = \frac{N}{4}$

Middle Quartiles $Q_2 = \frac{2N}{4} = \frac{N}{2}$

Upper Quartiles $Q_3 = \frac{3N}{4}$

Deciles

In general the k deciles D_k where

$$D_k = \frac{kN}{10}$$

Percentiles

In general the k^{th} percentile is P_k where
 $P_k = \frac{kN}{100}$ and

where N represents the total frequency

Example

The table below shows the distribution of marks obtained by 200 candidates in a practical biology examination. Make a cumulative frequency table for the data and construct the cumulative frequency curve. Hence, use the curve to estimate the median marks.

marks (100%)	1-10	11-20	21-30	31-40	41-50	51-60	61-70
frequency	2	5	11	20	26	40	40
71-80	81-90	91-100					
36	9	3					

Solution

In this case since the classes are not continuous, we need to compute the class boundary, as well as the cumulative frequency before drawing the cumulative frequency curve. The class boundary and cumulative frequency are obtained

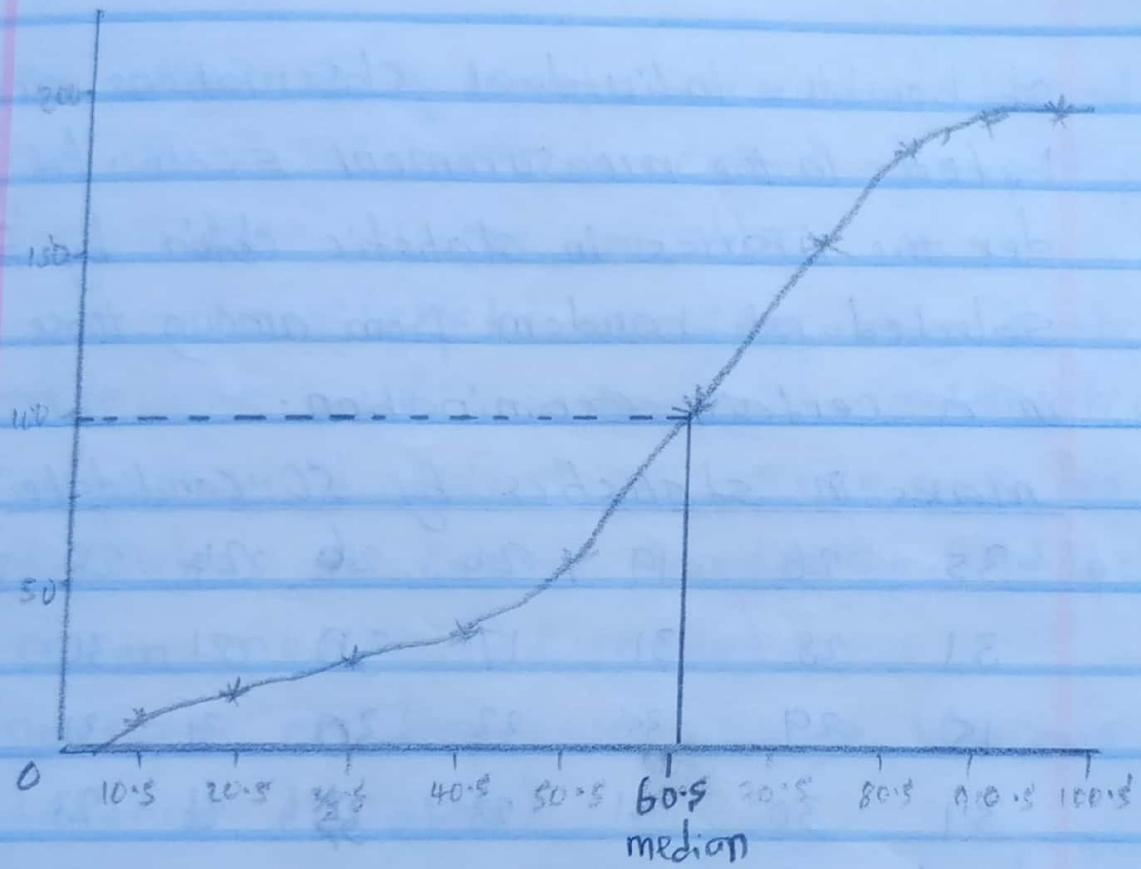
as follows

Intervals	frequency	Class boundary	C.F
1 - 10	2	0.5 - 10.5	2
11 - 20	5	10.5 - 20.5	7
21 - 30	11	20.5 - 30.5	18
31 - 40	20	30.5 - 40.5	38
41 - 50	26	40.5 - 50.5	64
51 - 60	42	50.5 - 60.5	106
61 - 70	46	60.5 - 70.5	152
71 - 80	36	70.5 - 80.5	188
81 - 90	9	80.5 - 90.5	197
91 - 100	3	90.5 - 100.5	200

The median is located as follows

$$\text{Median } Q_2 = \frac{2N}{4} = \frac{2+200}{4} = Q_2 = \underline{\underline{100}}$$

The median can now be estimated
the curve at the point where C.F is 100
The Cumulative frequency curve is drawn
below:



frequency distribution

frequency distribution - is an organized tabulation/graphical representation of the number of individuals in each categories on the scale of measurement. It allows the researcher to have a glance at the entire data conveniently. It shows whether the observations are high or low and also whether they are concentrated in one area or spread out across the entire scale. Thus, Frequency distribution present a picture

of how the individual observations are distributed. In the measurement scale, let us consider the marks in statistic obtain by 50 candidate selected at random from among those appear in a certain examination.

Marks in statistics by 50 candidate.

35	26	19	24	26	24	28	24	22
31	28	31	17	30	21	30	31	19
15	29	35	22	30	31	31	30	78
31	30	21	23	35	26	23	31	35
17	35	18	31	31	27	17	24	15

marks (x)	tally	frequency
15 - 19		9
20 - 24		11
25 - 29		10
30 - 34		15
35 - 39		5
		50

Such a table showing the distribution frequencies in the different class is called frequency table and the manner in which

the class frequencies are distributed over the class intervals is called grouped frequency distribution of the variable.

Characteristic of frequency distribution
there are four important characteristic of frequency distribution there are as follows

1. measures of central tendency and location (mean, median, and mode).
2. measures of dispersion (range, Varians and standard deviation)
3. the extent of symmetry/asymmetry (~~kurtosis~~ ~~skewness~~)
4. the flatness or peakedness (kurtosis)

measures of central tendency

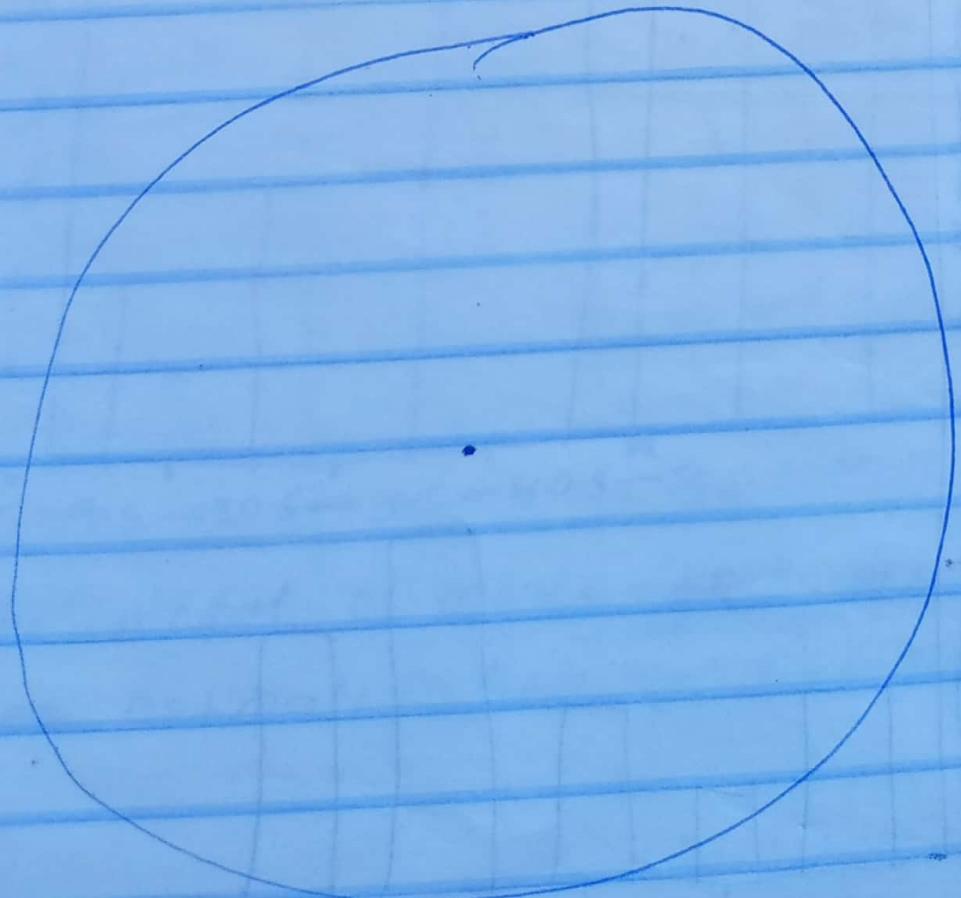
The measures of central tendency are otherwise called measures of location, they are scores or numerical values in central part of the distribution. They are values somewhere midway between the highest and the lowest observation which are used to represent the entire information. The following are the five measures of central tendency that are in common use :-

1. Arithmetic Mean or Simply mean
2. median
3. mode
4. Geometric mean
5. Harmonic mean

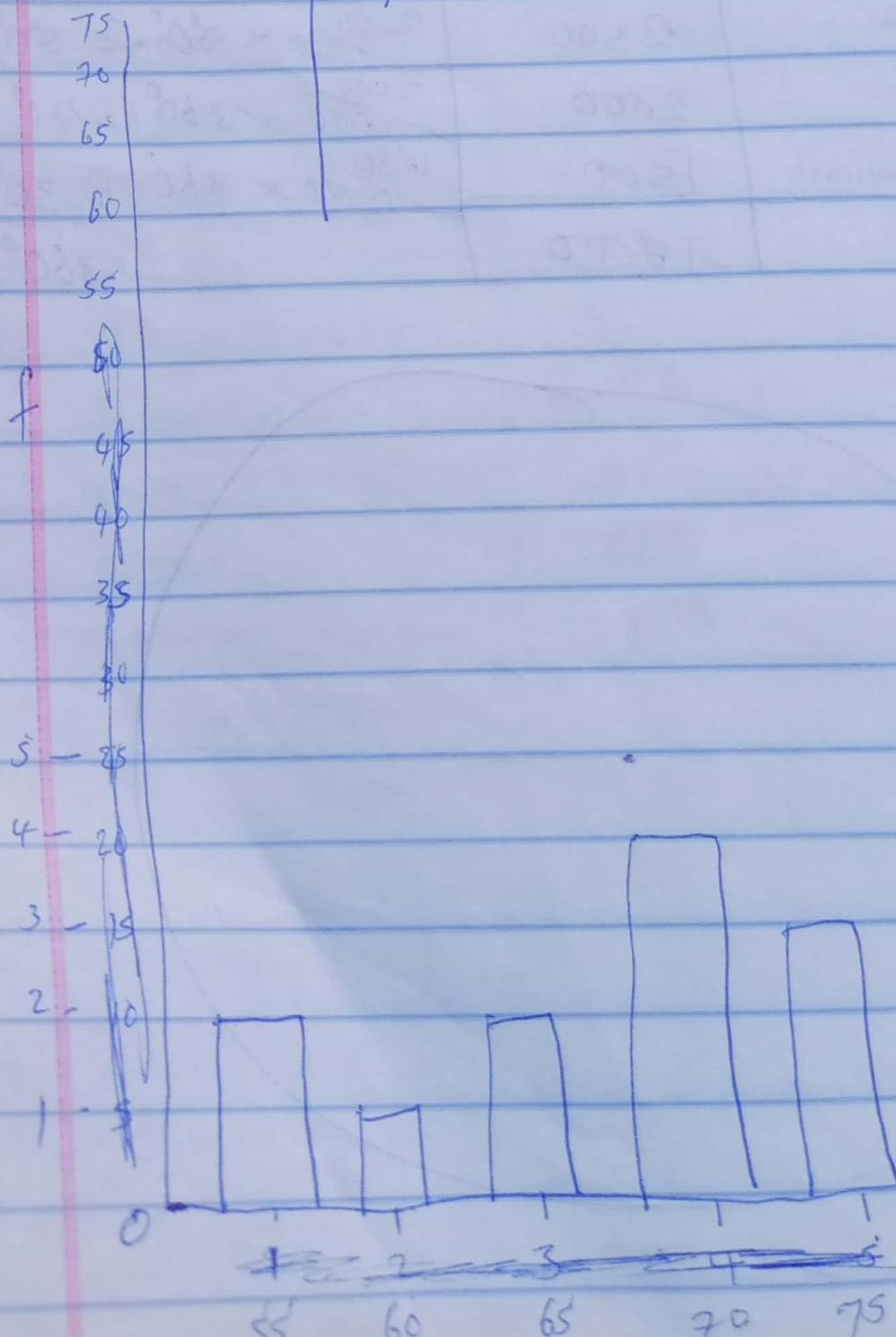
1. food items	Amounts (#)
Rice	3,000
meat	6,000
oil	500
vegetables	500
Beans	2,000
yam	2,500
fish	2,000
miscellaneous	1,800

Soln

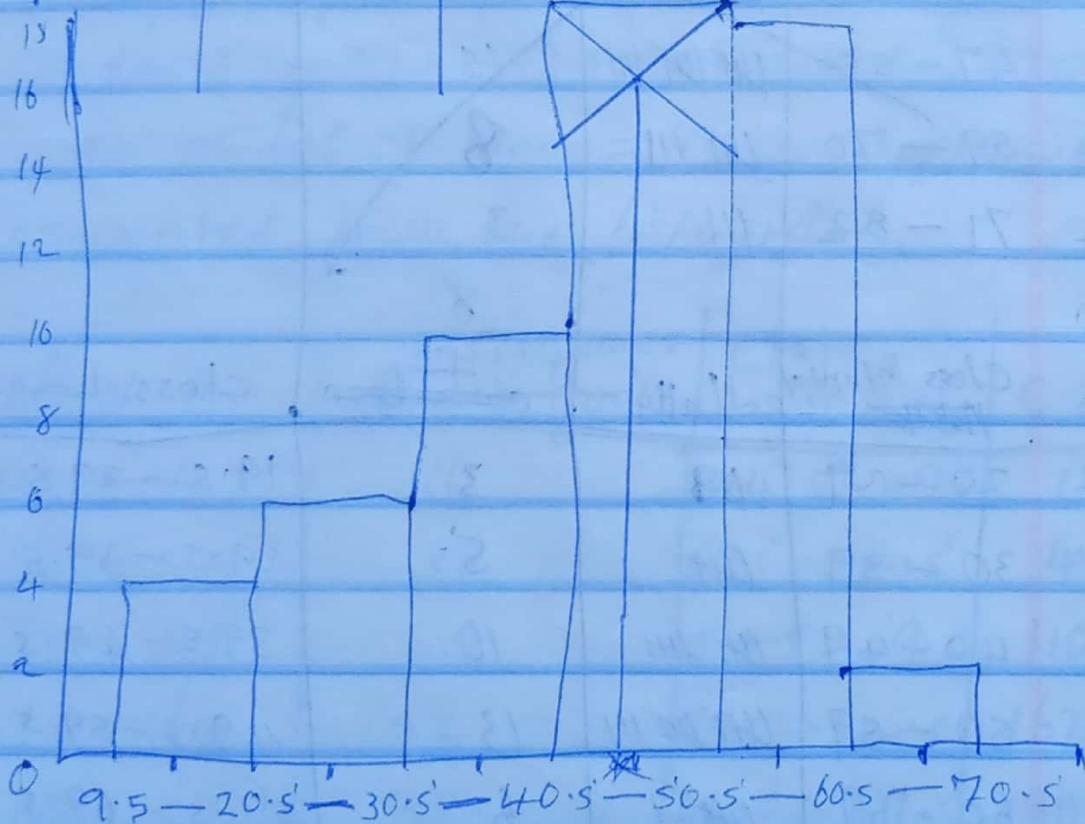
1. food items	Amounts	Angle $^{\circ}$
Rice	3,000	$\frac{3000}{18,000} \times 360^{\circ} = 60^{\circ}$
meat	6,000	$\frac{6000}{18,000} \times 360^{\circ} = 120^{\circ}$
oil	500	$\frac{500}{18,000} \times 360^{\circ} = 10^{\circ}$
vegetables	500	$\frac{500}{18,000} \times 360^{\circ} = 10^{\circ}$
Beans	2,000	$\frac{2000}{18,000} \times 360^{\circ} = 40^{\circ}$
yam	2,500	$\frac{2500}{18,000} \times 360^{\circ} = 50^{\circ}$
fish	2,000	$\frac{2000}{18,000} \times 360^{\circ} = 40^{\circ}$
miscellaneous	1,500	$\frac{1500}{18,000} \times 360^{\circ} = 30^{\circ}$
total	18,000	360°



frequency	kg
2	55
1	60
2	65
4	70
3	75



marks	f.	class boundary
10 - 20	4	9.5 - 20.5
20 - 30	6	20.5 - 30.5
30 - 40	10	30.5 - 40.5
40 - 50	20	40.5 - 50.5
50 - 60	18	50.5 - 60.5
60 - 70	2	60.5 - 70.5



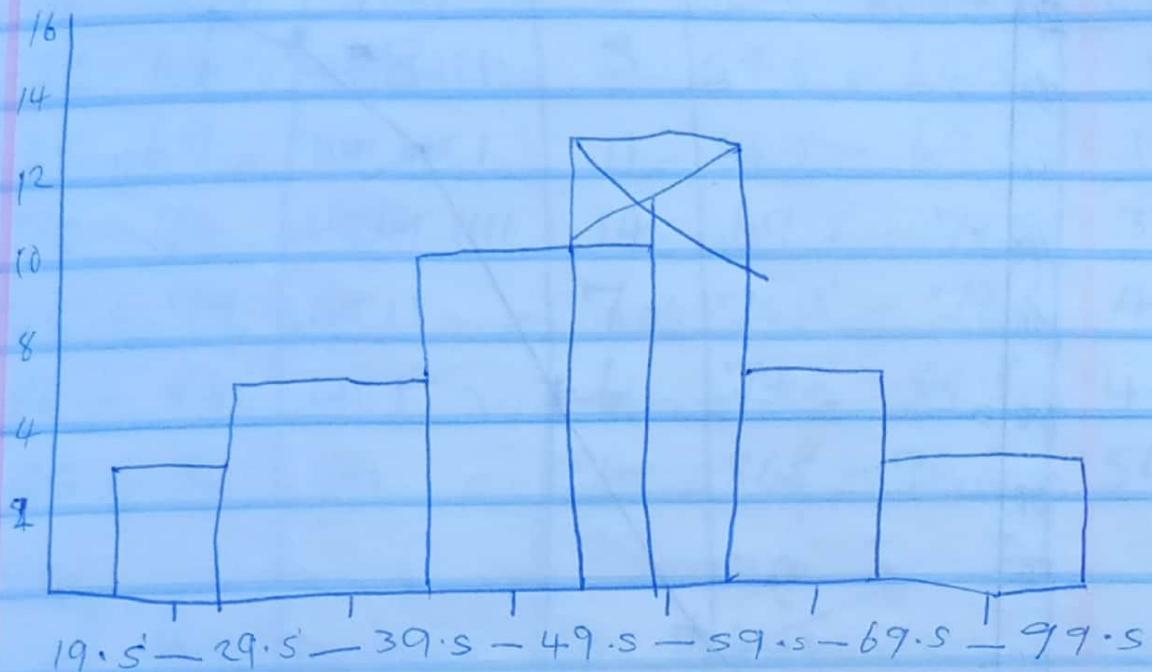
the modal of marks = 45.5 approximately
from estimated from the histogram

4. marks in statistic of 40 student

41	52	37	56	63	48	55	46	79	71
54		32	51	66	74	23	35	61	63 55
58		44	49	53	45	57	56	38	43 47
59		28	50	49	67	56	36	45	56 26

marks	tally	frequency
23 - 34		4
35 - 46		10
47 - 58		15
59 - 70		8
71 - 82		3

Class Interval Marks	Tally	Frequency	Class Boundry
20 - 29		3	19.5 - 29.5
30 - 39		5	29.5 - 39.5
40 - 49		10	39.5 - 49.5
50 - 59		13	49.5 - 59.5
60 - 69		6	59.5 - 69.5
70 - 79		3	69.5 - 79.5

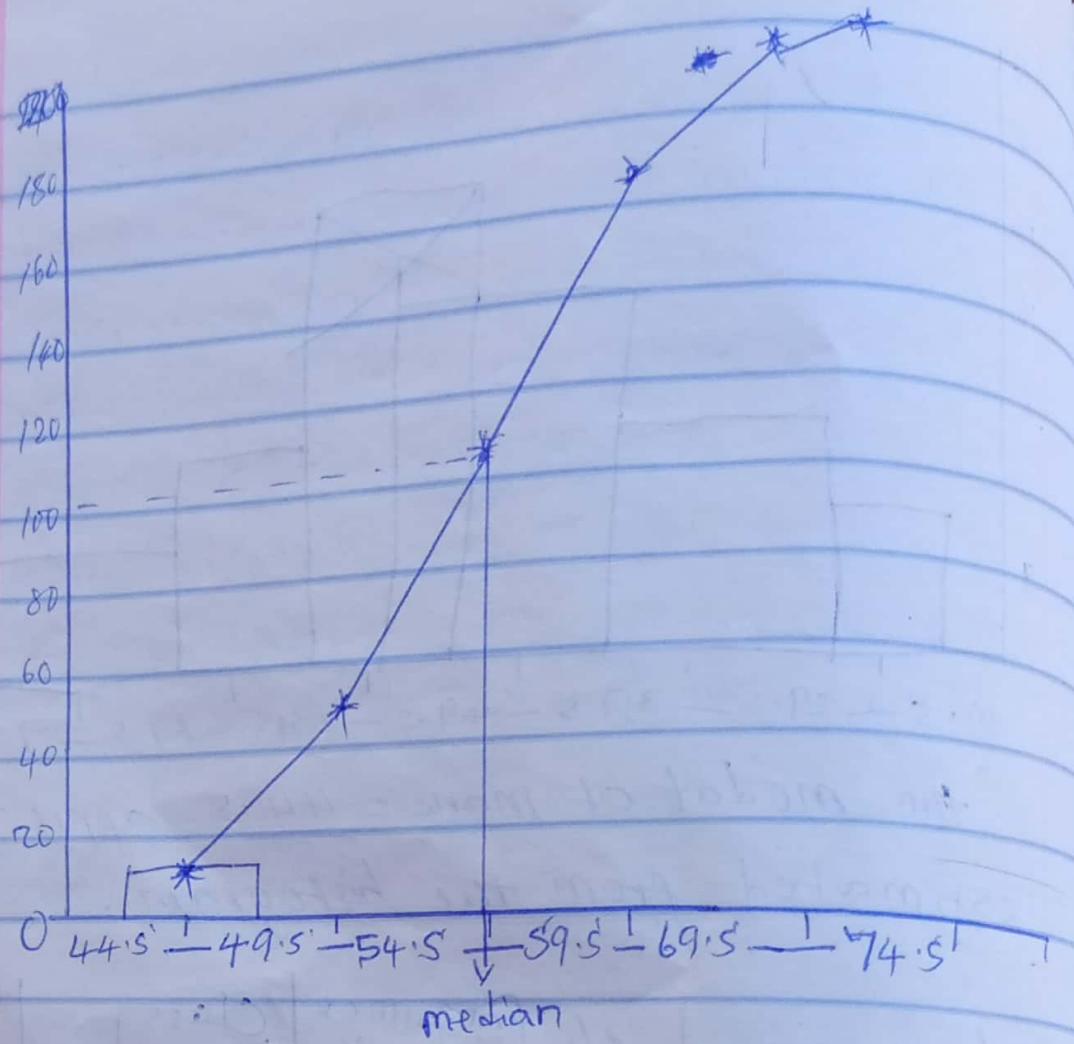


the modal of manc = 44.5 approximately estimated from the histogram.

S. No. of eggs	No. of quails	frequency	Class boundary	C. f
45 - 49	10	10	44.5 - 49.5	10
50 - 54	36	46	49.5 - 54.5	46
55 - 59	64	110	54.5 - 59.5	110
60 - 64	52	162	59.5 - 64.5	162
65 - 69	28	190	64.5 - 69.5	190
70 - 74	10	200	69.5 - 74.5	200

the median is located as follows.

$$Q_{12} = \frac{2N}{4} Q_2 = \frac{200 \times 20}{4} = \frac{400}{4} Q_2 = \underline{\underline{100}}$$



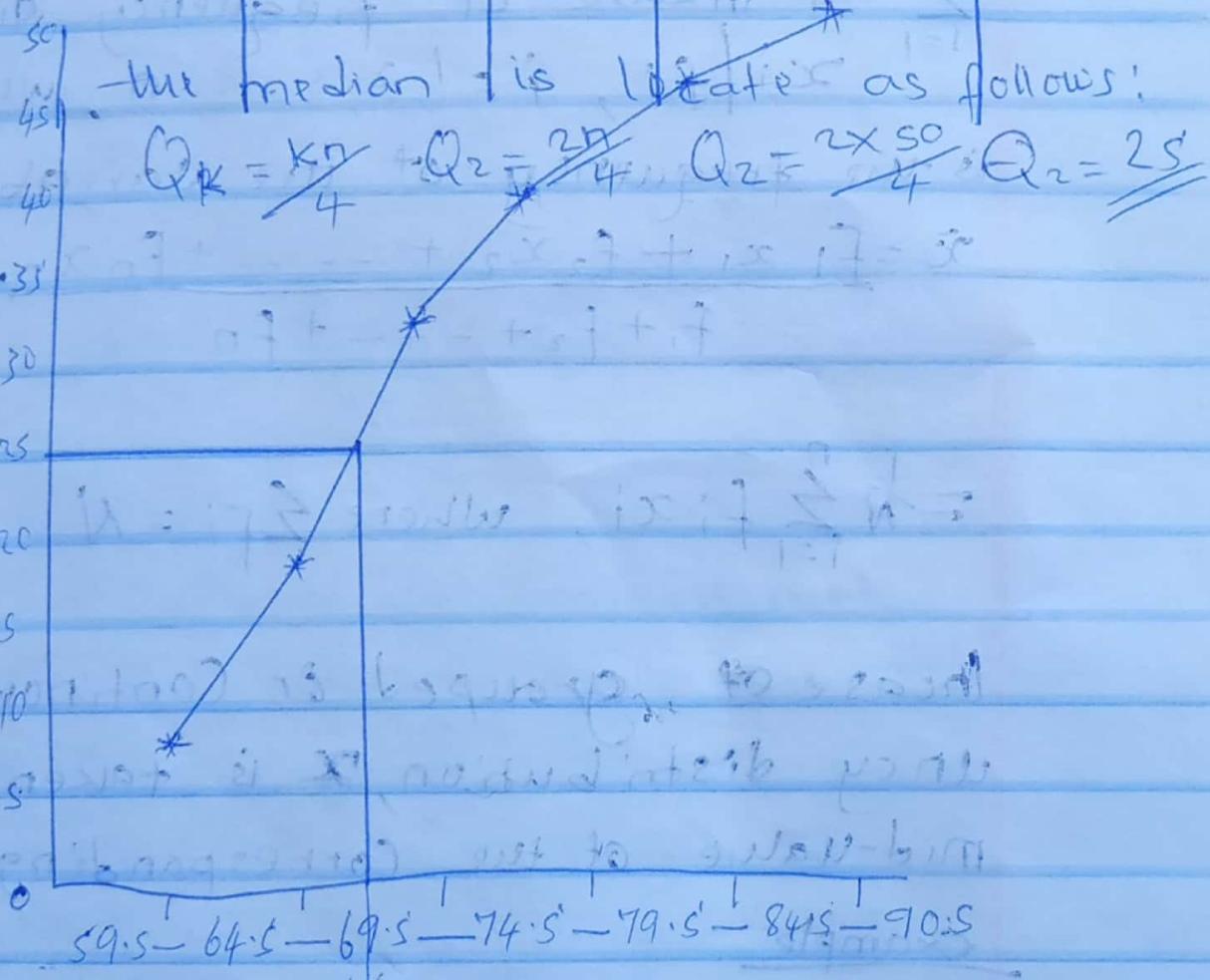
for 7th deciles $D_7 = \frac{140}{10}$

$$D_7 = \frac{7 \times 200}{10} \quad D_7 = \frac{1400}{10} \quad D_7 = 140$$

b. Weight of ⁵⁰ pregnant Women in Statistic

74	84	65	86	71	78	66	82	64
78	79	81	69	88	67	82	73	70
60	65	68	61	74	69	73	77	62
64	75	71	62	87	73	67	72	64
82	80	70	73	71	74	67	71	85

Class Interval	F tally	f.	class boundary	C.F
60 - 64	8	8	59.5 - 64.5	8
65 - 69	11	11	64.5 - 69.5	19
70 - 74	11	14	69.5 - 74.5	33
75 - 79	11	7	74.5 - 79.5	40
80 - 84	1	6	79.5 - 84.5	46
85 - 90	4	4	84.5 - 90.5	50



for 6th deciles $D_6 = \frac{6 \times 50}{10} = 30$

for 30th percentile $P_{30} = \frac{30 \times 50}{100} = 15$

Arithmetic mean

Arithmetic mean of set of observation is their sum divided by the number of observation.

Example - the arithmetic mean \bar{x} of n of observation x_1, x_2, \dots, x_n is given by $\bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n) = \frac{\sum x_i}{n}$ in case of frequency distribution $\sum_{i=1}^n x_i / f_i, i = 1, 2, 3, \dots, n$, where f_i is the frequency of the variable x_i .

$$\bar{x} = \frac{f_1 x_1 + f_2 x_2 + \dots + f_n x_n}{f_1 + f_2 + \dots + f_n} = \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i}$$

$$= \frac{1}{N} \sum_{i=1}^n f_i x_i, \text{ where } \sum f_i = N$$

In case of grouped or continuous frequency distribution, x is taken as the mid-value of the corresponding class.

Example

Find the arithmetic mean of the following frequency distribution.

$x: 1, 2, 3, 4, 5, 6, 7$

$f: 5, 9, 12, 17, 14, 10, 6$

solution

x	f	fixation formula	approximate value
1	5	5	01 → 0
2	9	18	05 → 01
3	12	36	08 → 08
4	17	68	04 → 08
5	14	70	02 → 04
6	10	60	06 → 02
7	6	42	09 → 38
$\sum f = 73$		$\sum fx = 299$	

$$\therefore \bar{x} = \frac{1}{N} \sum fx = \frac{299}{73} = \underline{\underline{4.09}}$$

example 2

calculate the arithmetic mean of the marks from the following table.

Marks: 0-10 - 10-20 - 20-30 - 30-40

No of stu: 12 18 27 20
 40-50 50-60

marks	No. of student	mid point x_c	$f x$
0 - 10	12	5	60
10 - 20	18	15	270
20 - 30	27	25	675
30 - 40	20	35	700
40 - 50	17	45	765
50 - 60	6	55	330
	$\sum f = 100$		$\sum f x = 2,800$

Arithmetic mean

$$\bar{x} = \frac{1}{N} \cdot \sum f x = \frac{2,800}{100} = 28$$

$$\bar{x} = 28$$

Assume mean method

Assume mean — if the values of x or (and) f are large, calculation of mean by the formula above is quite time-consuming and tedious. The arithmetic is reduced to a great extent by taking deviations of the given values from any arbitrary point "A" as explained below.

Let $d_i = x_i - A$, then $f_i d_i = f_i (x_i - A)$
 $f_i d_i = f_i x_i - A f_i$

Summing both sides over i from 1 to n
we get

$$\sum_{i=1}^n f_i d_i = \sum_{i=1}^n f_i x_i - A \sum_{i=1}^n f_i$$

$$\frac{\sum_{i=1}^n f_i d_i}{N} = \frac{\sum_{i=1}^n f_i x_i}{N} - \frac{A \sum_{i=1}^n f_i}{N}$$

$$\frac{\sum_{i=1}^n f_i d_i}{N} = \bar{x} - A$$

$$\therefore \bar{x} = A + \frac{\sum_{i=1}^n f_i d_i}{N}$$

$$\text{or } \boxed{\bar{x} = A + \frac{1}{N} \sum f_i d_i}$$

Any number can solve the purpose of
arbitrarily point ' A ' but usually the value
of x corresponding to the middle part of
the distribution will be much more convenient
In case of the grouped or continuous

frequency distribution the arithmetic
reduced to a still greater extend by
taking $d_i = \frac{x_i - A}{h}$

Where A is an arbitrary point and h
is the common magnitude of class
interval in this case we have
 $h d_i = x_i - A$ and proceeding ex-

actly as above, we get

$$\bar{x} = A + \frac{h}{N} \sum_{i=1}^n f_i d_i$$

Example

calculate the mean for the following

frequency distribution

class Interval : 0-8 8-16 16-24 24-32 32-

frequency : 8 7 16 24 15
40-48

Solution

Take $A = 28$

$h = 8$

~~X/R~~

class Interval	mid point (x)	frequency (f)	$d = \frac{(x-A)}{h}$	$f \cdot d$
0 - 8	4	8	$\frac{4-28}{8} = -3$	-24
8 - 16	12	7	$\frac{12-28}{8} = -2$	-14
16 - 24	20	6	$\frac{20-28}{8} = -1$	-6
24 - 32	28	24	$\frac{28-28}{8} = 0$	0
32 - 40	36	15	$\frac{36-28}{8} = 1$	15
40 - 48	44	7	$\frac{44-28}{8} = 2$	14
		$\sum f = 77$	$\sum fd = -25$	

$$\bar{x} = A + \frac{h}{N} \sum_{i=0}^7 f_i d_i$$

$$A = 28 \quad N = \sum f = 77 \quad h = 8 \quad \sum fd = -25$$

$$\bar{x} = 28 + \frac{8}{77} \times (-25)$$

$$\bar{x} = 28 + (-2.59)$$

$$\bar{x} = 28 - 2.59$$

$$\bar{x} = 25.41$$

Advantages of Arithmetic Mean

1. It is rigidly defined
2. It is easy to calculate and widely understood
3. It can be determined when only the total

Value and the number of observations are available

4. It is liable to mathematical precision disadvantages

1. It is affected by extreme values

2. It may not correspond to any value in the data set

3. It requires assumption when Open Intervals values are involved

4. Arithmetic mean is not be used if we dealing with qualitative characteristics which can not be measured quantitatively. Intelligent, honesty etc in such cases median is the only average to be used

Median =

L = lower class boundary of median class

h = magnitude of the class interval

f = frequency of the median class

N = total frequency = $\sum f$

Cf = cumulative frequency of the class before the median

$$\text{Median} = L + \frac{h}{f} \left(\frac{N}{2} - cf \right)$$

class
L = lower class boundary of the median

median - the median is the middle number in a given set of data or a frequency distribution when the numbers are arranged in order of magnitude ascending or descending order. Thus for example, for an orderly arranged set $(x_1, x_2, x_3, x_4, x_5, x_6, x_7)$, the median is x_4 . On the other hand, the median in the orderly arranged set, $(x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8)$ is $\frac{(x_4 + x_5)}{2}$.

frequency distribution median is obtained by considering the cumulative frequency. The steps for calculating median are given below:

1. find $N/2$, where $N = \sum f_i$
2. See the cumulative frequency just greater than $N/2$.
3. The corresponding value of x_c is the median.

Example

Obtain the median for the following frequency distribution:

x_i	1	2	3	4	5	6	7	8	9
f_i	10	11	16	20	15	15	9	6	

x	f	CF
1	8	8
2	10	18
3	11	29
4	16	45
5	20	65
6	25	90
7	15	105
8	9	114
9	6	120

$$\sum f = 120$$

$$\text{Hence } N = 120 \Rightarrow N/2 = \frac{120}{2} = 60$$

Cumulative frequency just greater than $N/2$ is 65 and the value of x corresponding to 65 is 5. Therefore median is 5.

In the case of continuous f.d., the class corresponding to the C.F. just greater than $N/2$ is called the median class and the value of median is obtained by the formula:

$$\text{median} = L + \frac{h}{f} [N/2 - c.f_b]$$

where L = is the lower class boundary of the median class.

f = is the frequency of the median class.

h = is the magnitude of the median class.

$C.F_b$ = is the cumulative frequency before the median class.

$$N = \sum f$$

Example:

find the median wage of following distribution.

wages (#) = 20 - 30 30 - 40 40 - 50 50 - 60

No of labourers = 3 5 20 10

60 - 70 70 - 80 80 - 90 90 - 100

5

solution

wages (#)	No of Lab	C.F	C.b - c.b
20 - 30	3	3	19.5 - 30.5
30 - 40	5	8	29.5 - 40.5
40 - 50	20	28	39.5 - 50.5
50 - 60	10	38	49.5 - 60.5
60 - 70	5	43	59.5 - 70.5

$$\sum f = 43$$

$$\text{Hence } \frac{N}{2} = \frac{43}{2} = 21.5$$

$$m = L + \frac{h}{f} \left(\frac{N}{2} - C.F_b \right)$$

$$L = 39.5, h = 10, f = 20, C.F_b = 8, N = 43$$

$$m = 39.5 + \frac{10}{20} (21.5 - 8) =$$

$$m = 39.5 + 6.75$$

$$\underline{\underline{m = 46.25}}$$

Example 2 Exercise

Compute the median of the data below, which depicts the hourly wage (₹) of the 60 junior workers in a pharmaceutical firm.

Hourly Wage (₹) : 10-19 20-29 30-39 40-49

No of workers :	4	6	10	20
-----------------	---	---	----	----

50-59	60-69		
-------	-------	--	--

18	12		
----	----	--	--

Solution

H.W (#)	N ₀ of W	C.F	C.B
10-19	4	4	9.5 - 19.5
20-29	6	10	19.5 - 29.5
30-39	10	20	29.5 - 39.5
40-49	20	40	39.5 - 49.5
50-59	18	58	49.5 - 59.5
60-69	2	60	59.5 - 69.5

$$\sum f = 60$$

Hence $N/2 = \frac{60}{2} = 30$

$$m = L + \frac{h}{f} (N/2 - C.F_b)$$

Advantages of median

1. It is not affected by extreme values.
2. Computation is very easy.
3. Since it is an actual value, it is always representable and realistic.

Disadvantages

1. For group data, the median is only an estimate.
2. It does not take all values into account in its computation.
3. It is not useful or practicable in large data.

$$\text{mean Deviation} = \frac{\sum_{i=1}^n (x_i - \bar{x})}{N}$$

$$\text{Variance } \sigma_w^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{N}$$

$$\text{Standard deviation } \sigma_w = \sqrt{\text{Variance}} = \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{let } w = \{x_1, x_2, x_3, x_4, x_5\} \\ \{2, 6, 5, 7, 5\}$$

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{(w)} = \frac{2+6+5+7+5}{5} = \frac{25}{5} = 5$$

$$\begin{aligned} \text{mean deviation} &= \frac{(2-5)+(6-5)+(5-5)+(7-5)+(5-5)}{5} \\ &= \frac{(-3)+(1)+(0)+(2)+(0)}{5} = \frac{0}{5} = 0 \end{aligned}$$

$$\text{Variance } \sigma_w^2 = \frac{(2-5)^2+(6-5)^2+(5-5)^2+(7-5)^2+(5-5)^2}{5}$$

$$= \frac{(-3)^2+(1)^2+(0)^2+(2)^2+(0)^2}{5} = \frac{9+1+4}{5}$$

$$= \frac{14}{5}$$

$$\text{Standard deviation } \sigma_w = \sqrt{\frac{14}{5}} = 0.7$$

$$(x - \bar{x})^2 = (1) > \text{normal}$$

$$(x - \bar{x})^2 = (1)^2 = 1 \rightarrow \text{sub back}$$

$$1 + 2 + 3 + 5 = 11 \rightarrow \text{total}$$

$$(2-3)^2 = 2+1+2+4+5 = (n) \cdot \frac{1}{n} = \bar{x}^2$$

$$(2-3)(3-3) + (2-3)(2-3) + (2-3)(1-3) = 10 \rightarrow \text{sum}$$

$$(2-3)(3-3) + (2-3)(2-3) + (2-3)(1-3) = 10$$

$$(2-3)^2 + (2-3)^2 + (2-3)^2 + (2-3)^2 + (2-3)^2 = 10 \rightarrow \text{sum}$$

$$4+1+4+1+1 = 11 \rightarrow \text{sum}$$