

DISCRIMINATIVE STATISTICS

Statistics

Descriptive

Central Tendency

Mean

Median

Mode

Graphic Display

Dot Plot

Frequency Distribution

Histogram

Frequency Polygon

Bar chart

Pie chart

Box and Whiskers

Inferrential

Hypothesis Testing

Estimation

Point

Interval

Regression Analysis

Descriptive Statistics

Enroll
Page No: 21
Date: / /

Statistics \Rightarrow used as singular or plural sense

Singular— In singular case statistics refers to the procedures used to organize and interpret observed data.

In this statistics is defined as branch of mathematics that deals with collecting, organizing and summarizing the data and drawing conclusions about the environment from which the data was collected.

Plural Sense—

Statistics are quantitative values that are used to describe a set of observed data.

Thus some time ~~states~~ "Statistics is" and "Statistics are" are used ~~in~~ depending on context

Statistics Vs. Probability

Probability:

- Fully defined probability problems have unique and precise solutions.
- Probability laws apply across an entire population of interest.

Statistics:

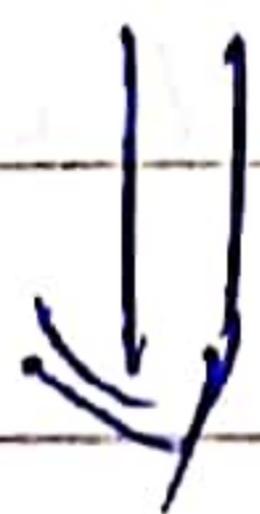
Statistics is concerned with the relationship between an observed segment of a population and the entire population.



In statistics we are interested in understanding an observed observation that is based on a segment of a population of interest and how the observation can apply to the entire population.

Statistician works by formulating the data in ways
that make sense ↓
~~start~~

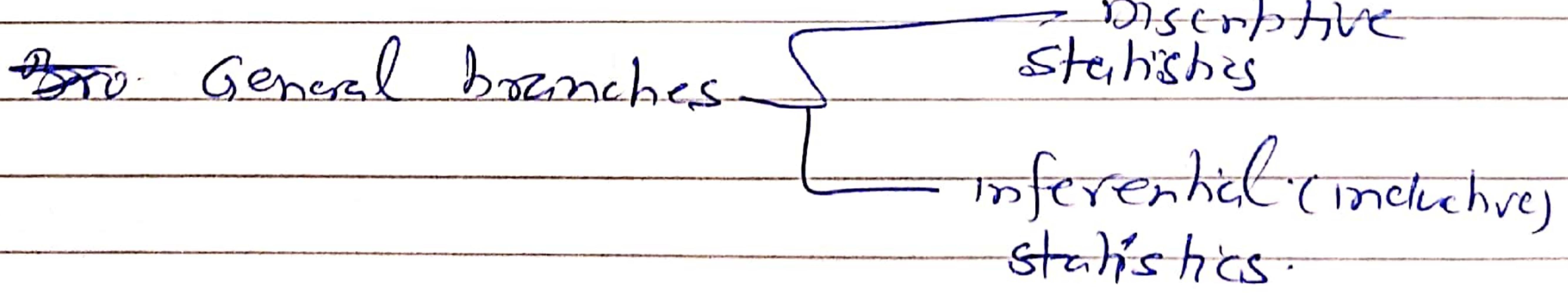
Postulate the Probabilistic model.
for system under investigation based
on the physical mechanisms involved in the
system and on personal experience



Statistician expects that the model

to exhibit a probabilistic behavior
that is similar to that of the physical
system.

Branches of Statistics



Descriptive Statistics :-

Concerned with collecting, organizing, and
summarizing the raw scores in the more
meaningful ways.

Data = raw score are measurements or
Observed ~~values~~ Values

Inferential Statistics

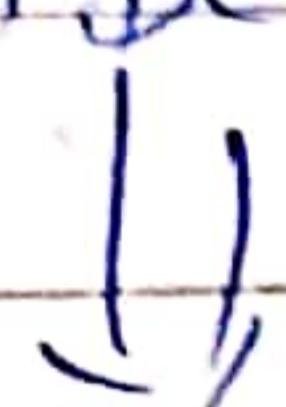
In inferential statistics deals with procedures or techniques that can be used to study a segment of the population called a sample and make generalizations about the population from which the sample was obtained with the help of probability.

Descriptive Statistics

→ This deals with collecting, grouping, and presenting data in a way that can be easily understood.



make sense of observed data.



present data the data in a more meaningful way, which allows simpler interpretation of ~~data~~

the data.

e.g Grade of two students
= (i) distribution of grade
= (ii) spread of grade.

organizing the data by a set of graphs, bar charts,

tables, or frequency distributions.



Descriptive statistics does not, however, allow us to make conclusions beyond the data we have

analyzed; it is simply a way to describe a set of observed data by providing

simple summaries about the sample and measures.

Methods that are used to describe the set of observed data.

a) Measures of central tendency: —

→ Ways of describing the central position of a frequency distribution for a group of data.

b) Measures of Spread: —

Ways to summarize a group of data by describing how spread out the data values

c) Graphical Displays: —

Ways to visually visualize the data to see how it is distributed and if any patterns emerge from the data.

Measures of Central Tendency

describe the "center" of the data set.

→ central tendency describes the tendency of the observations to bunch around a particular value.

(→ measure of central tendency are numerical)

summaries that are used to summarize a data set ~~+~~ with a single "typical" number.

Three ~~main~~ main measures of central tendency:

(i) Mean

(ii) Median

(iii) Mode

} measures of the "average" of the distribution of the data ~~set~~ set.

Mean! — Mean \Leftrightarrow Average

Let $x_1, x_2, x_3, \dots, x_N$ be data set, the mean

is given by:

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_N}{N} = \frac{1}{N} \sum_{i=1}^N x_i$$

e.g. 66, 54, 00, 56, 34, 12, 40, 50, 00, 50, 90, 65
Q1

Then, the mean is given by:

$$\bar{x} = \frac{66 + 54 + 00 + 56 + 34 + 12 + 40 + 50 + 00 + 50 + 90 + 65}{12}$$

$$\bar{x} = 59.75$$

Note:- Mean may not be a member of the data set.

~~Median:-~~

Median

Divide the data into two equal halves

where each half contains 50% of the data.

The numerical value where the data set is divided is called the median.

For median computation:

(i) ordered rank by arranging them

in increasing order of magnitude.

(ii) Median data position is calculated

Even no of data values

odd number of data values

R_m : rank-position of the median in the rank-ordered data set.

No. Number of data values.

(iii) If N is odd:

$$R_m = \frac{N+1}{2}$$

And median is data value at the position R_m .

If N is even then the value of

the median is given by

$$M = \frac{d_{\frac{N}{2}} + d_{\frac{N}{2}+1}}{2}$$

$$M = \frac{d_{\frac{N}{2}} + d_{\frac{N}{2}+1}}{2}$$

d_k : data value at position k .

Ex:-

Arrange the data (D) in average section in

increasing order

$$\begin{array}{cccccccccccccc} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 \\ 12, 34, 40, 50, 50, 54, 56, 65, 66, 80, 80, 90 \end{array}$$

(ROD)

$N=12$ data values. \Rightarrow Even \Rightarrow median

is not member of data set
 \Downarrow

It is average of data values at $\frac{12}{2}=6^{\text{th}}$ and
 $\frac{12}{2}+1=7^{\text{th}}$ position of sorted data.

Therefore, median:

$$M = \frac{54 + 56}{2} = 55$$

Note:- Append 95 to the rank data set (ROD).

$N=13 = \text{odd}$.

$$\therefore R_m = \frac{13+1}{2} = 7$$

$M=56$ is the member of the data set.

Mode

→ Mode is data value that occurs the most frequently in the data set.

12, 34, 40, 50, 50, 54, 56, 65, 66, 80, 88, 90, 95 → RD2

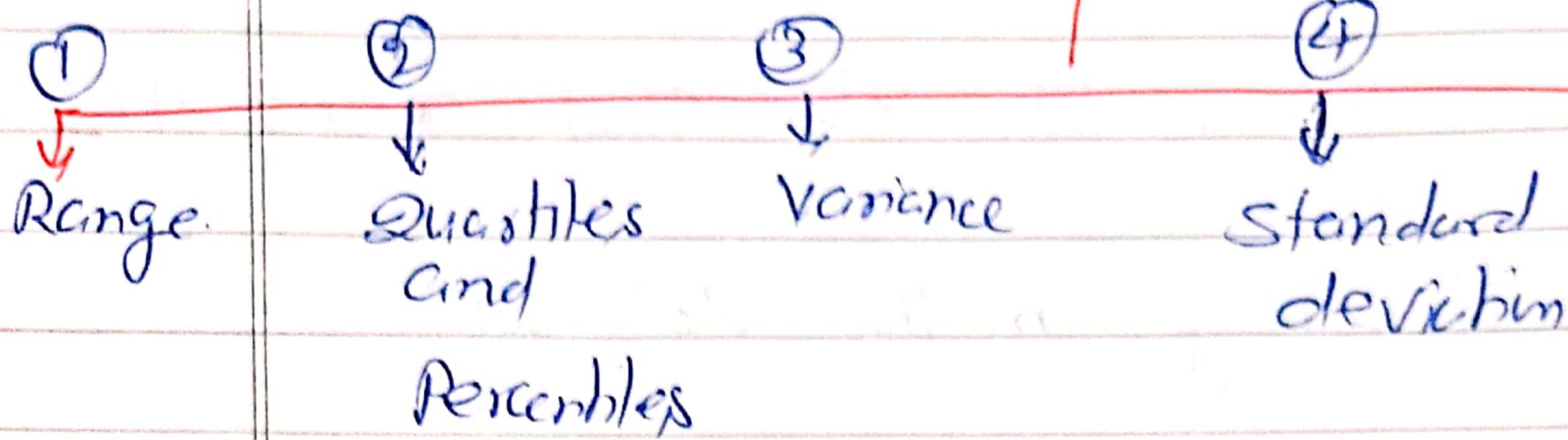
In RD1 and RD2 50 occurs twice while every other value occurs once.

Thus mode of RD1 and RD2 is 50.

and the data sets are said to be unimodal.

→ Sometimes data set can have more than one mode, in this case it is said to be a multi-modal set.

→ ? Measure of Dispersion



→ Range:

↳ measure of statistical dispersion or spread.

Range = Maximum data value - Minimum Data Value

Range = Largest data value - Smallest data value

→ Quartiles And Percentiles:

↳ divide the ordered data sets into quarters.

↳ b th percentile of an ordered data set is value such that at least $100p\%$ of the observations are at or below this value

and at least $100(1-p)\%$ are at or above

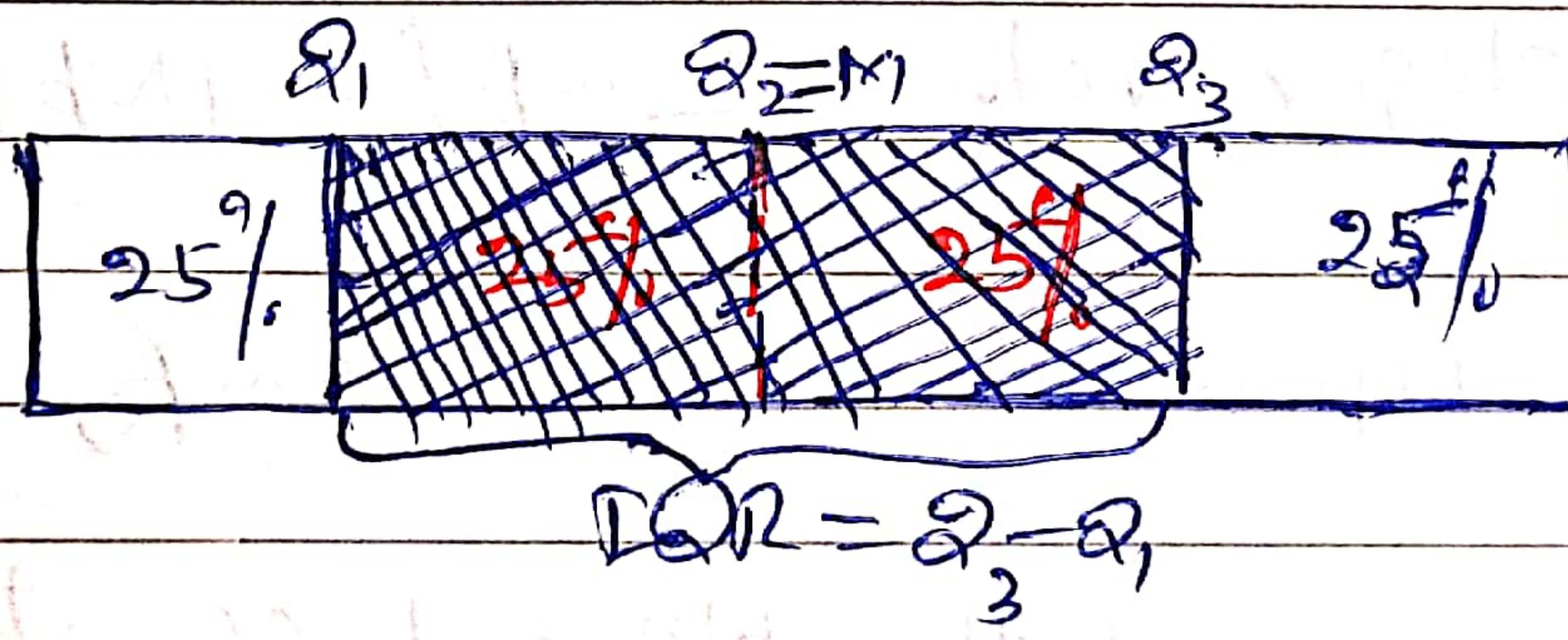
this value, where $0 < p < 1$.

- (i) First Quartile $Q_1 = 25^{\text{th}}$ percentile and $p = 0.25$
- (ii) Second Quartile $Q_2 = 50^{\text{th}}$ percentile and $p = 0.50$
- (iii) Third Quartile $Q_3 = 75^{\text{th}}$ percentile and $p = 0.75$

Second quartile is median

Interquartile range (IQR) ? This is difference between the third quartile and the first quartile.

$$\text{IQR} = Q_3 - Q_1$$



$\text{IQR} \equiv$ middle one-half (or 50%) of an ordered data set.

Quartile computation depends on the ordered data set i.e. data position in a rank-ordered data set and not on the data value itself.

Procedure For Combining the 100th percentile is as follows:

(i) Order the N observations from the smallest to the largest.

(ii) Compute the product Np . If Np is not integer, round it to next integer and find the corresponding ordered value

i.e. $\lceil Np \rceil$ i.e. compute $\lceil Np \rceil$

Integer value $\lceil Np \rceil$ or $\lfloor Np \rfloor + 1$

(iii) If Np is an integer, say k , then calculate the mean of the k^{th} and $(k+1)^{\text{th}}$ ordered value.

e.g.: 1 2 3 4 5 6 7 8 9 10 11 12
 66, 54, 80, 56, 34, 12, 40, 50, 80, 50, 90, 65
 Given ordered

1 2 3 4 5 6 7 8 9 10 11 12
 12, 34, 40, 50, 50, 54, 56, 65, 66, 80, 80, 90

$N=12$, for first quartile $p=0.25$.

Therefore first quartile $Np = 12 \times 0.25 = 3$

$\therefore Q_1 = \text{mean of } 3^{\text{rd}} \text{ and } 4^{\text{th}} \text{ element of the ordered data}$

$$\text{First quartile } Q_1 = \frac{48 + 50}{2} = 49$$

For Third quartile ~~p=0.75~~ $p = 0.75$

$$Np = 12 \times 0.75 = 9$$

$\therefore \text{Third quartile } Q_3 = \text{mean of the } 9^{\text{th}} \text{ and } 10^{\text{th}} \text{ element of the ordered data}$

$$\text{Third quartile } Q_3 = \frac{66 + 70}{2} = \frac{136}{2}$$

$$\text{Third quartile } Q_3 = 73.$$

Inter-quartile range (IQR) = $Q_3 - Q_1$,

$$= 73 - 49$$

$$\boxed{IQR = 24}$$

For second quartile $p = 0.50$

$$Np = 12 \times 0.50 = 6$$

Median = mean of 6^{th} and 7^{th} elts value in ordered data set

$$\text{Median} = \frac{54 + 56}{2} = 55 \quad \boxed{\text{Median} = 55}$$

Alternative Method for computing LQR!

(i) Find the median: it divides the data, Q_2 , into two halves.

lower half < median

upper half > median

(ii) → Find the median of lower half, which is Q_1 .

Find the median of upper half, which

is Q_3

In previous example

$$\text{Median } Q_2 = \frac{54+56}{2} = 55$$

lower half of ordered data,

12, 34, 48, 50, 56, 54.

In this total number of data values = 6

Which even therefore median of lower

half = average of 3rd and 4th elements

$$\begin{aligned} &= \frac{48+50}{2} \\ &= \frac{98}{2} \\ &= 49 \end{aligned}$$

Thus first quartile $Q_1 = 49$.

Upper half of the data value 56, 65, 66, 80, 88, 90

Here $N=6$, which is even. Therefore third quartiles Q_3 is the average of 3rd and 4th

$$Q_3 = \frac{66 + 80}{2} = 73$$

$$\therefore IQR = 73 - 49$$

$$\boxed{IQR = 24}$$

~~Median~~

Variance

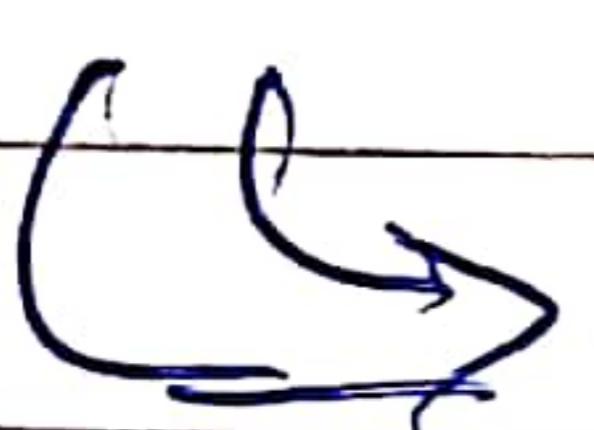
Deviation Score

Subtract each data value

from the mean, the obtained value
called a deviation score.



deviation score gives the numerical
distance between data value and
data set's "typical" value.



The sum of all the deviation scores
equals zero as shown in table below.

Data Value x	Deviation $x - \bar{x}$ $(\bar{x} = 59.75)$	squared Deviation $(x - \bar{x})^2$
66	8.25	68.0625
54	-3.75	14.0625
88	30.25	915.0625
56	-1.75	3.0625
34	-23.75	564.0625
12	-45.75	2093.0625
48	-9.75	95.0625
50	-7.75	60.0625
80	22.25	495.0625
50	-7.75	60.0625
90	32.25	1040.0625
65	7.25	52.5625
Sum	$\sum(x - \bar{x}) = 0$	$\frac{1}{12} \sum(x - \bar{x})^2 = 455.021$

→ The data values above and below the mean have positive and negative deviation scores, respectively, that cancel each other out.

Dr. Manoj Kumar Singh
Associate Professor
abharac.in

Remove the negative score \Rightarrow square the deviation scores \Rightarrow Sum of squared deviation scores

\Rightarrow The average of square deviation ~~last~~ gives ^{score} the value of Variance.

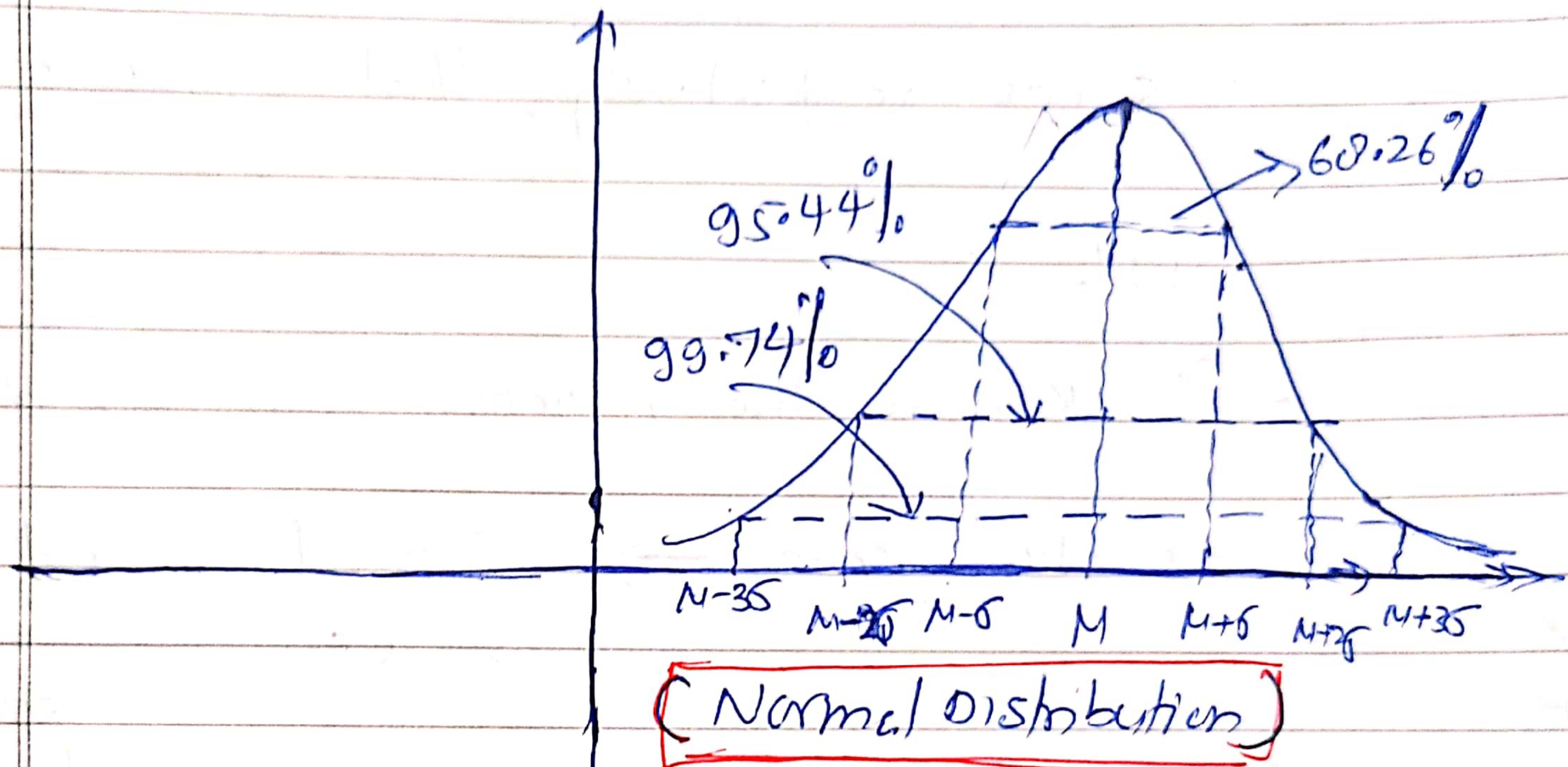
$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$

For example ~~given~~ discussed above:

$$\sigma^2 = \frac{1}{12} \sum_{i=1}^{12} (x_i - \bar{x})^2 = 455.02$$

Standard Deviation:

Standard deviation (σ) = square root of Variance



Probability that observation lies within :

$$P[M_x - \sigma_x \leq X \leq M_x + \sigma_x] = 0.6826$$

$$P[M_x - 2\sigma_x \leq X \leq M_x + 2\sigma_x] = 0.9544$$

$$P[M_x - 3\sigma_x \leq X \leq M_x + 3\sigma_x] = 0.9974$$

Graphical And Tabular Displays

organize the data into a graphical or tabular

form so that ~~a trend~~ a trend, if any, emerging
out of the data can be seen easily

Dr. Me
ASD
@bhuvan

(i) Dot plots (ii) Frequency distribution (iii) Bar

charts (iv) Histograms, (v) Frequency Polygon.

(vi) Pie charts, (vii) Box plots and ~~(viii) Whiskers plot~~

(viii) Whiskers plot

Dot plots

Dot plots = dot chart

A dot plot uses dots to show where the data values (or scores) in a distribution.

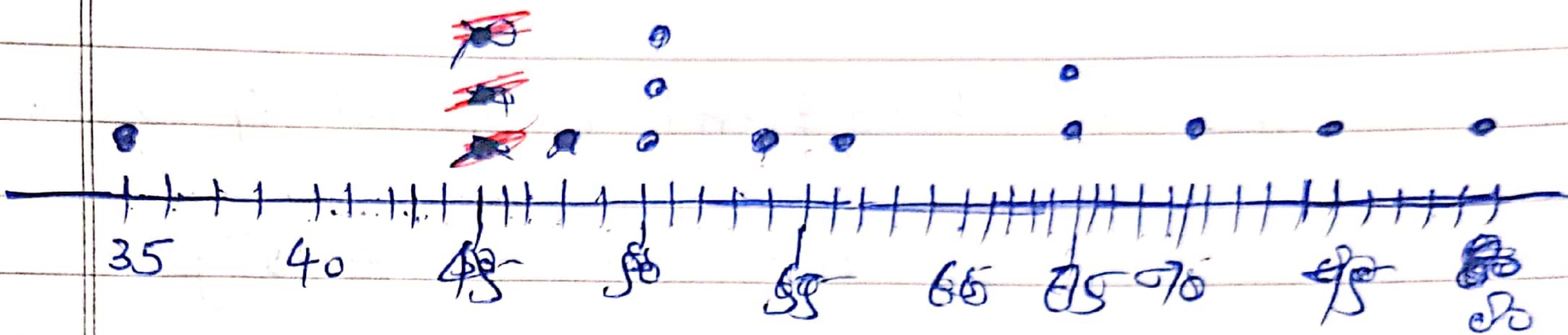
~~uses~~ dot dots are plotted against their

actual data values that are on the horizontal
~~scale~~ scale.

→ If there are identical data values, the dots
are "piled" on top of each other.

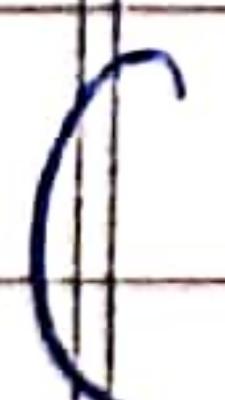
35, 40, 50, 50, 50, 54, 56, 65, 65, 70, 75, 80

Date: 11



Dot plot is useful when number of data is small in size.

Frequency distribution



- Table that lists the set of values in a dataset and their frequency frequencies
- frequency \equiv number of times each occurs in the data sets

e.g.

35, 40, 50, 50, 50, 54, 56, 65, 65, 70, 75, 80.

In a frequency distribution we find the
relative frequency of the different values

Table	Frequency
X	
36	1
48	1
50	3
54	1
56	1
65	2
70	1
75	1
80	1

When data covers ~~large~~ wide range of values that a list of all ~~the~~ the X values would be too long to be a "simple" presentation of ~~data~~ data.

→ In this case grouped frequency distribution

table in which the X column lists groups of data values, called class intervals, rather than individual data values.

→ The width of each class can be determined by dividing ^{the} range of observations by the number ~~of~~ of classes.

→ Have equal class widths, and ~~the~~ the class interval ~~should~~ should be mutually exclusive and non-overlapping.

Class width = difference between the lower limit of two consecutive classes
 = difference between the upper limit of two consecutive classes

Class mark = is the ~~middle~~ number in the middle of the class

= Average of lower and upper limit

Group frequency distribution

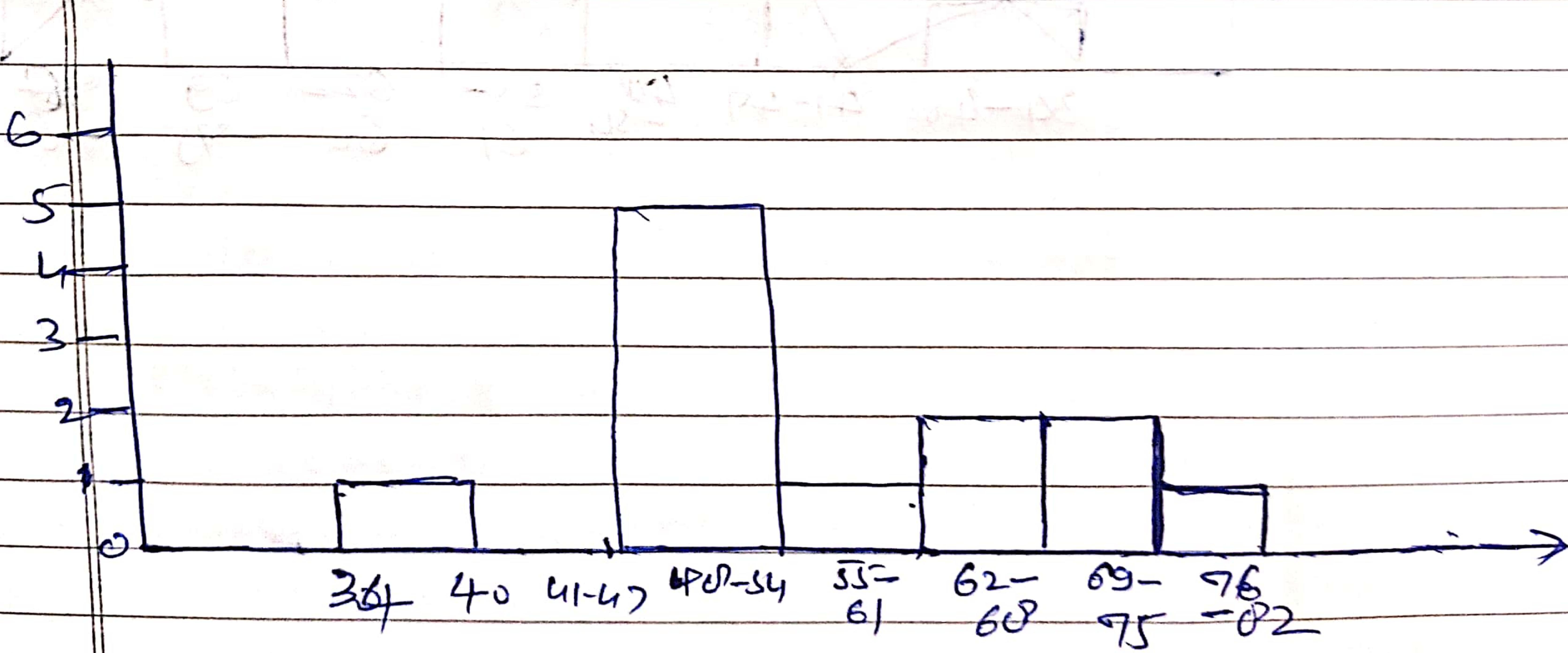
class	Frequency	class mark
30 - 40	1	39
41 - 47	0	44
48 - 54	5	51
55 - 61	1	58
62 - 68	2	65
69 - 75	2	72
76 - 82	1	79

Histogram :-

Frequency histogram (or simply histogram) is used to graphically display the grouped frequency distribution.

- Draw vertical bars above the classes so that that the height of a bar corresponds to the frequency of the class that it represents
- width of the bar extends to the real limits of the score class. Thus, the columns are of equal width, and there are no spaces between

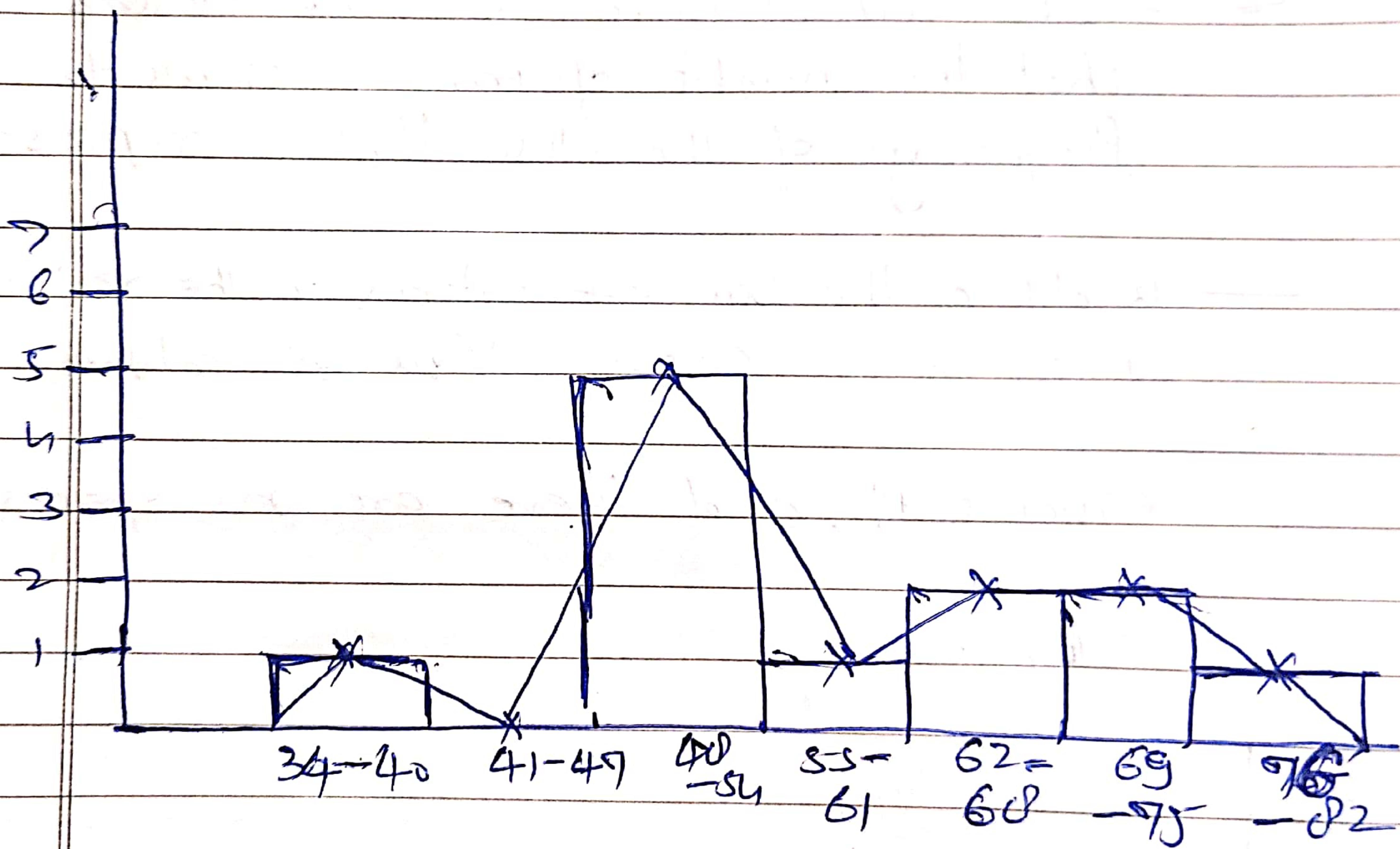
Column:



Frequency Polygons:

This gives idea about shape of the distribution.

Obtained by joining the class marks of a histogram with ^{the two} end points lying on the horizontal axis.



: Bar Graphs :

Used for ~~present~~ graphical representation of

the categorical Variable

→ As the variable that has two or more categories with no intrinsic ordering to the categories.

→ e.g. gender is categorical variable with two categories : Male - female.

→ length (horizontal), height (vertical) shows

the frequency for each category or

characteristics.

Eg. survey of 100 computer science students

to determine that their interest in subjects

(Prob) Probability : 30

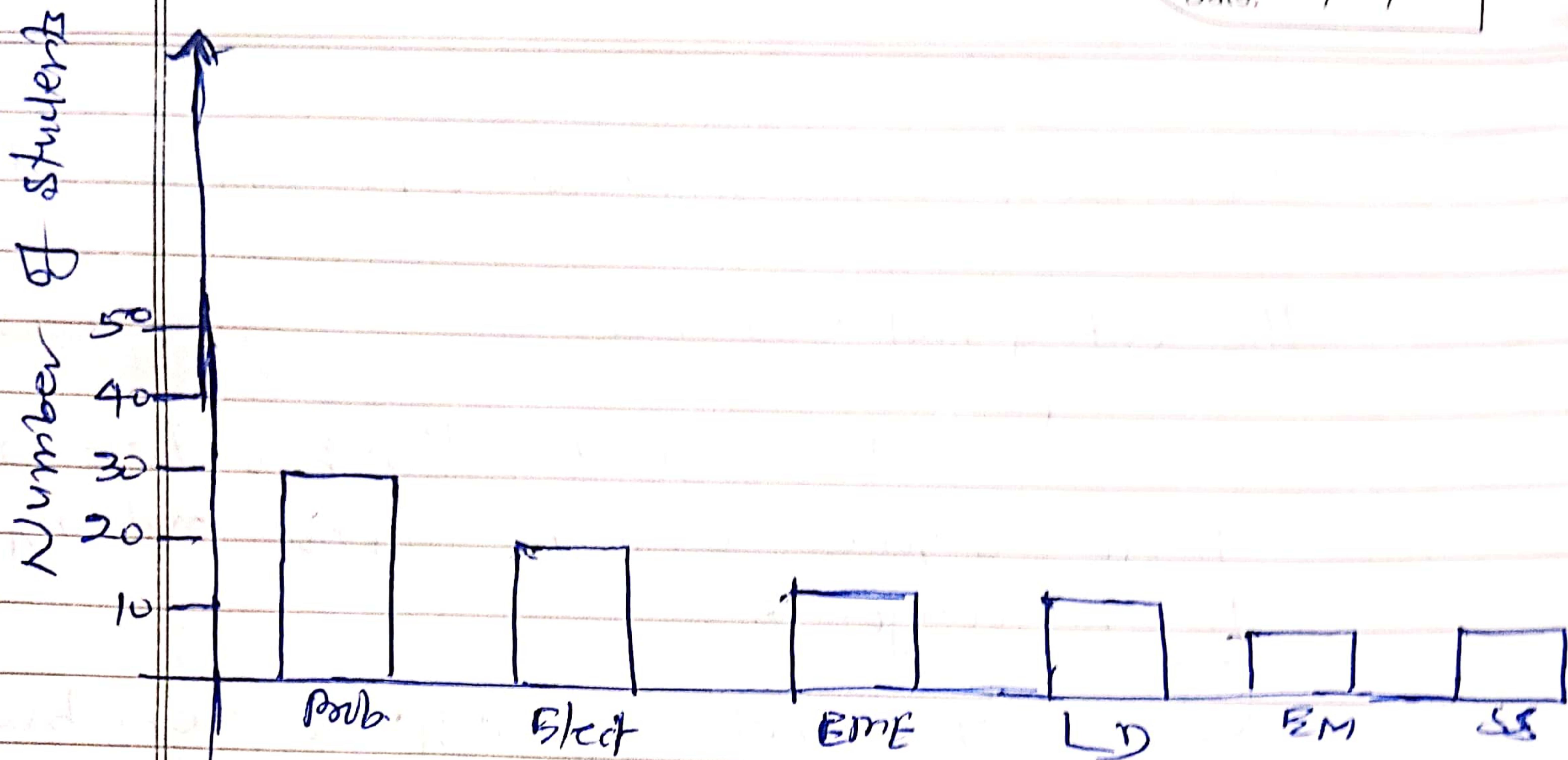
(Elect) Electronics : 20

(EM) Electromechanics : 15

(LD) Logic Design : 15

(Em) Electromagnetics : 10

(SS) Signal and System : 10

Note:

- (i) Gaps are included between the bars of each category.
- (ii) Bars can be arranged in any order without affecting the start, after, older.

Pie chart

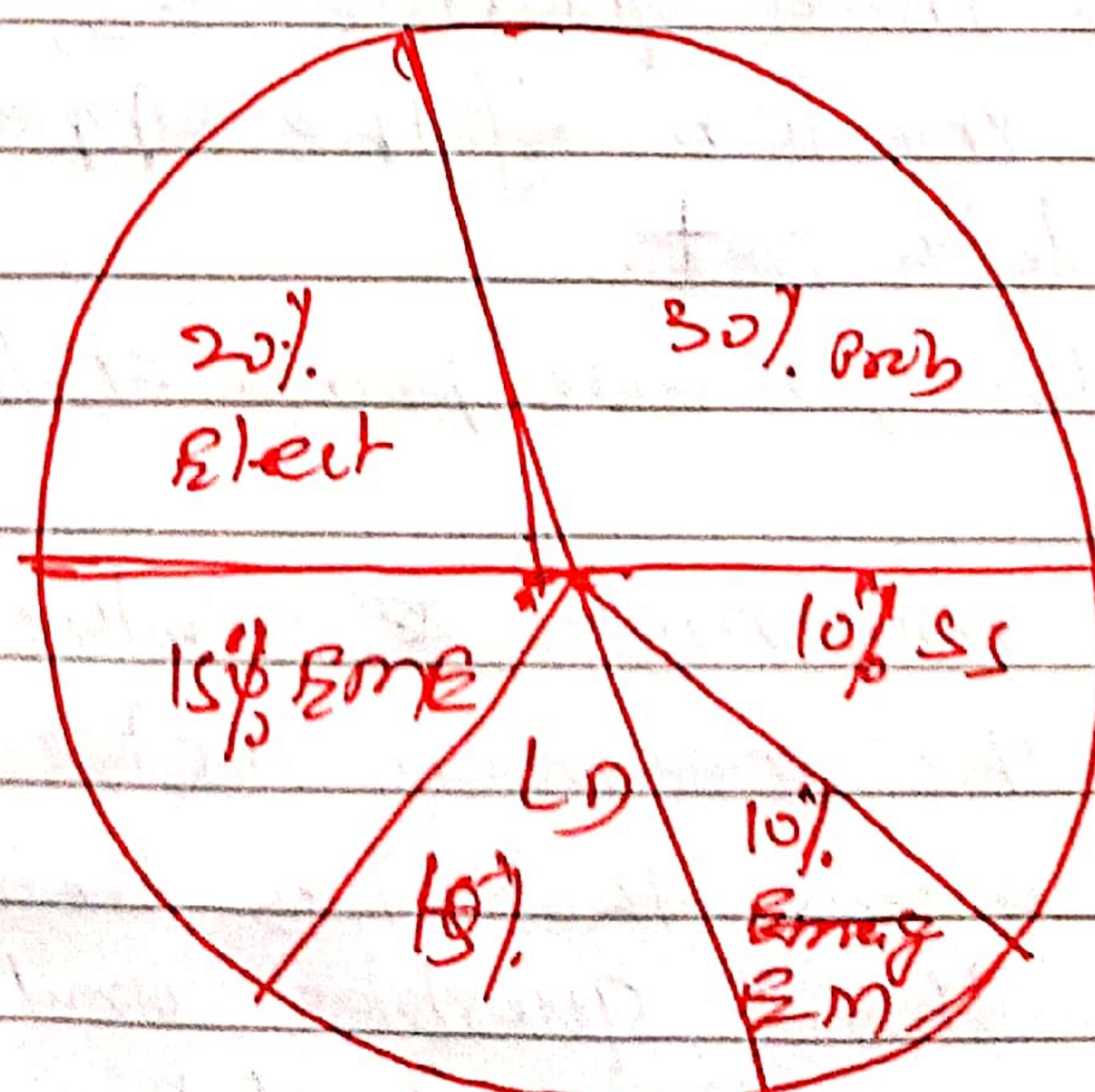
- "Pie slices" are used to show relative sizes of data.
- The size of each slice is proportional to the probability of the event that the slice represents.

Example the survey of ECE student to find out their favorite subjects.

i. Probability (Prob) : 30% ii. Electronics & 20% (Elect)

iii. Electromechanics (EM): 15%; iv. Logic Design 15%. v. Electromagnetics (EM): 10%.

vi. Signal and Systems: 10%

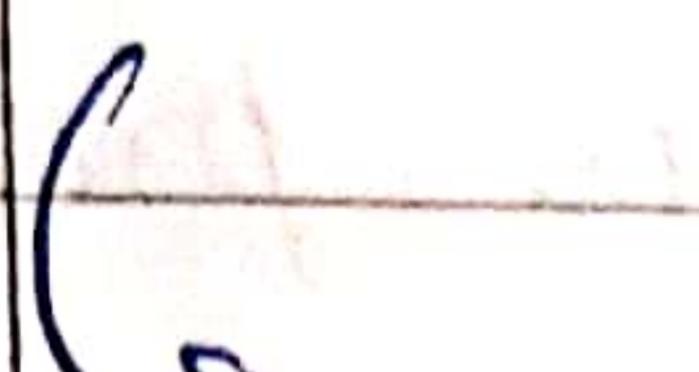


Box and Whiskers Plot



Way to ~~isolate~~ visually organize data

into fourths or ~~quart~~ quartiles



make a box including first and third quartiles and straight lines (whiskers) extending from the ends of

the box to the minimum and maximum

data values.

Procedure :-



1. Arrange the data in increasing order

2. Find the median

3. Find the first quartile Q_1 , which is the median of the lower half of the data set; and third quartile, Q_3 , which is the median of the upper half of the data set.

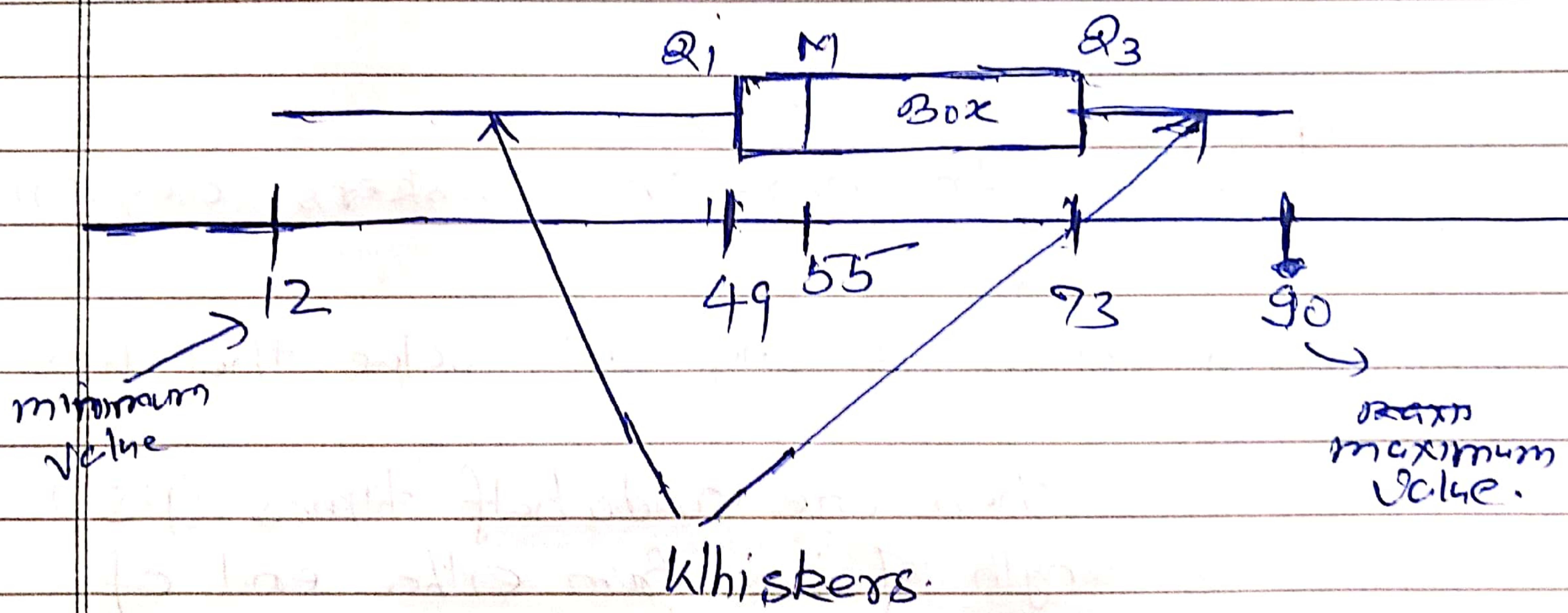
4. on a line, mark points at the median, Q_1 , Q_3 , the minimum ~~value~~ value of the data set and the maximum value of the data set

5. Draw a box that lies between the first and third quartiles and thus represents the middle 50% of the data.

6. Draw a line from the first quartile to

the minimum data value, and another line from the third quartile to the maximum data value. These lines are whiskers of the plot.

e.g.: 12, 34, 49, 50, 50, 54, 56, 65, 66, 80, 80, 80, 90



Outliers — Values much higher or much lower than all of the other values.

such values are known as outliers

These outliers are usually excluded from the whisker portion of the ~~the~~ box and whiskers diagram. They are plotted individually and labeled as outliers.

Interquartile range (IQR) and outliers

$IQR = Q_3 - Q_1$ = width of box in whiskers plot diagram

$IQR \Rightarrow$ is measure of dispersion



Statistics assumes that data values are

clustered around some central value.

Outliers

In box-and-whiskers diagram,

An outlier is any data value that lies

more than one and a half times (1.5)

the length of box from either end of

box.

i.e. data point below $Q_1 - 1.5 \times IQR$

or above $Q_3 + 1.5 \times IQR$ is ~~viewed~~

as being too far from the central values

to be ~~too~~ reasonable.



Thus $Q_1 - 1.5 \times IQR$ and $Q_3 + 1.5 \times IQR$

are the "fences" that marks of the

1) ~~"reasonable"~~ values form the outlier values.

Lower fence : $Q_1 - 1.5 \times IQR$

Upper fence : $Q_3 + 1.5 \times IQR$

In previous example:-

$$IQR = Q_3 - Q_1 = 73 - 49$$

$$IQR = 24 \quad IQR \times 1.5 = 36$$

$$\text{Lower fence} = Q_1 - IQR \times 1.5 = 49 - 36 = 13$$

$$\text{Lower fence} = 13$$

$$\begin{aligned}\text{Upper fence} &= Q_3 + IQR \times 1.5 \\ &= 73 + 36 = 109\end{aligned}$$

$$\text{Upper fence} = 109.$$

The only data value outside the ~~fence~~ fences

is 12, and all other values are within two fences

∴ Thus 12 is only outlier in the data set

Example 1 — Find outlier with the help of

Box plot — Whiskers plot

10, 12, 8, 1, 10, 13, 24, 15, 15, 24

Rank this in increasing order:

$$\begin{array}{ccccccccc} \text{Position} & \rightarrow & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \\ \hline \text{Value} & \rightarrow & 1 & 8 & 10 & 10 & 12 & 13 & 15 & 15 & 24 & 24 \end{array}$$

Total No. of data (N) = 10 = even

→ median is average of 5th and 6th position

Element in sorted sequence

$$\begin{aligned} M &= \frac{12+13}{2} \\ M &= 12.5 \end{aligned}$$

Lower half data set: 1, 8, 10, 10, 12

whose median $Q_1 = 10$

Upper half data set: 13, 15, 15, 24, 24

whose median $Q_3 = 15$

Thus $IQR = Q_3 - Q_1 = 15 - 10 = 5$

Lower fence = $Q_1 - 1.5 \times IQR$

$$= 10 - 1.5 \times 5$$

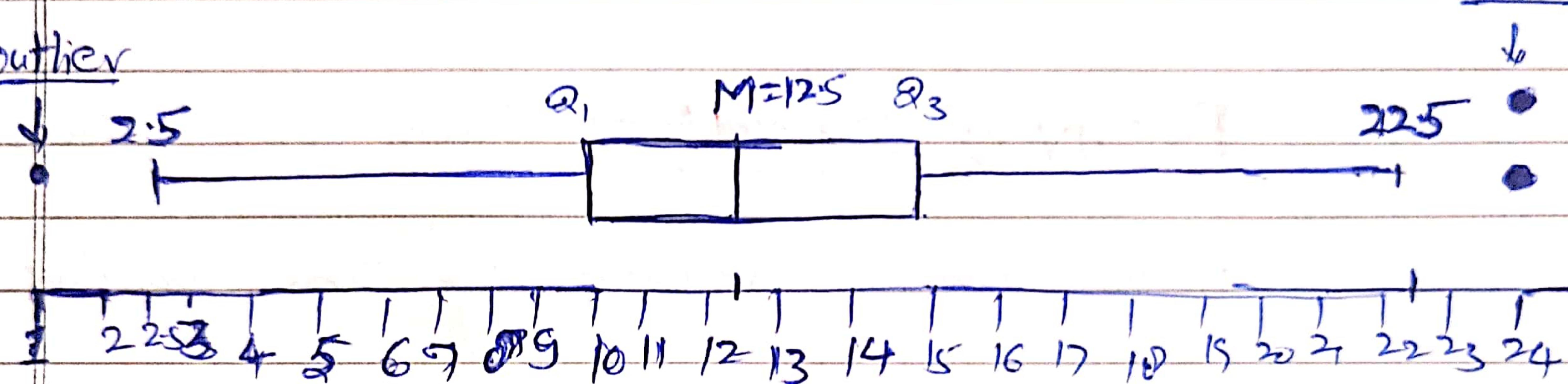
$$= 10 - 7.5$$

Lower fence = 2.5

$$\text{Upper fence} = Q_3 + IQR \times 1.5$$

$$\text{Upper fence} = 15 + 7.5$$

$$\text{Upper fence} = 22.5$$

outlieroutlier

↓

•

•

22.5

•

•

→ Shape of Frequency Distribution

~~Skew = $\frac{M_3 - 3M_1}{2M_2}$~~

(SKEWNESS)

Any distribution / curve which is not symmetrical about longitude is said to be skewed.

Positively skewed — lower data values

have higher frequencies. The positively

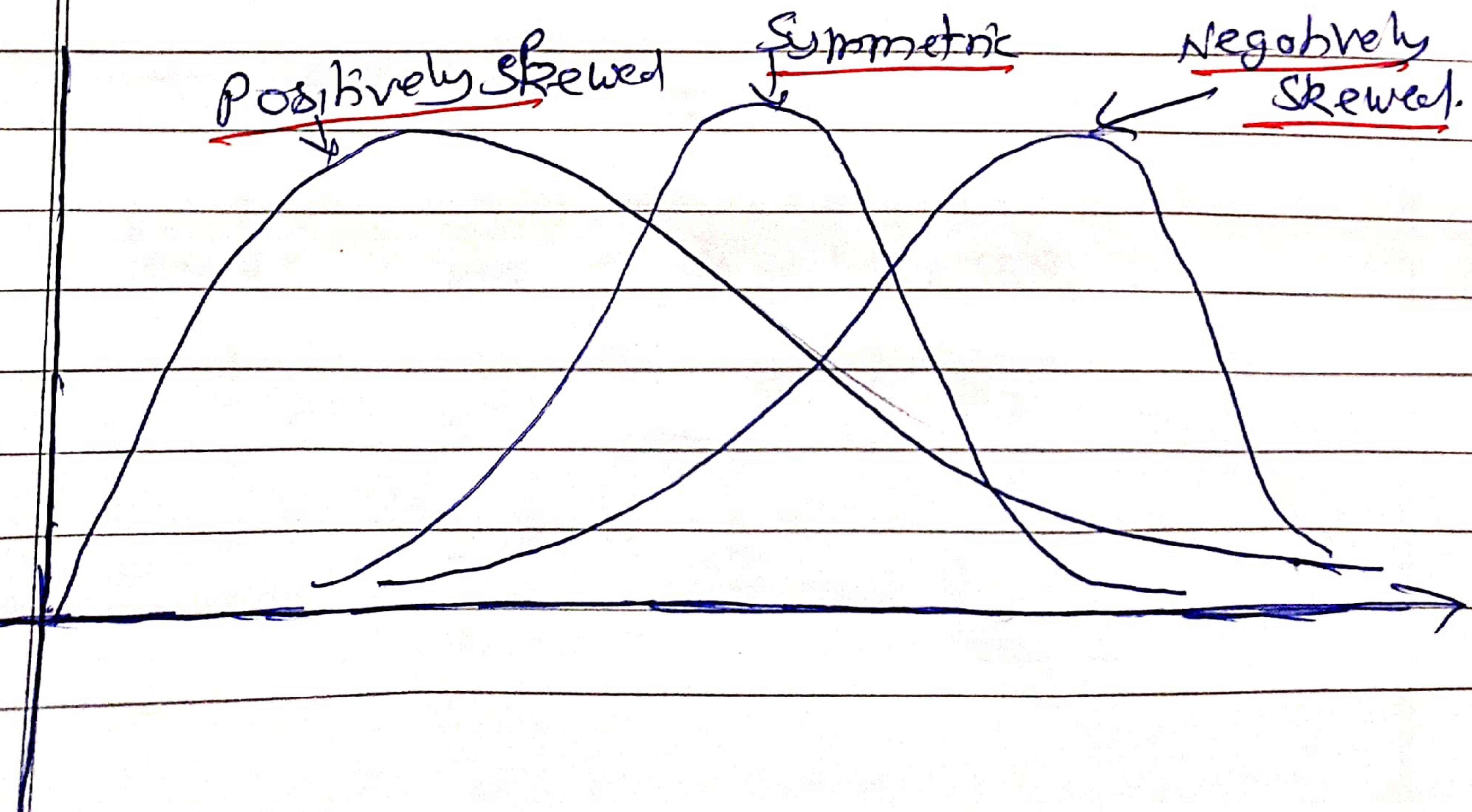
skewed distribution is right-tailed distribution

Negatively skewed — Higher data

values have higher frequencies, the distribution

is negatively skewed. The negatively skewed

distribution is a left-tailed distribution.



1

2

Quantitative measure of Skewness :-

Quantitatively, the skewness of a random variable

X is denoted by γ_1 and defined as

$$\gamma_1 = E \left[\left(\frac{X - \mu_X}{\sigma_X} \right)^3 \right]$$

$$\gamma_1 = \frac{1}{\sigma_X^3} E[(X - \mu_X)^3] \quad \text{--- (1)}$$

where μ_X : mean

σ_X : Standard deviation.

(1) Skewness computation based on population:-

Let dataset size (N) (size of population)

Then skewness is calculated as follows

$$\gamma_1 = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{\left[\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{3/2}} \quad \text{--- (2)}$$

(2) Skewness Computation Based on Sample Size (n):

The skewness is given by

$$G_1 = \frac{\sqrt{n(n-1)}}{n-2} \gamma_1 \quad \text{--- (3)}$$

{ Skewness Measure Based on Pearson's }
 Coefficient of skewness }

Ist Version

x_{med} : median of a data set

x_{mode} mode

The Pearson's mode / first skewness coefficient is defined by:

$$\gamma_1 = \frac{M_x - x_{\text{mode}}}{\sigma_x} \quad \text{--- (4)}$$

IInd Version

The Pearson's median / second skewness coefficient is used when the mode of a

distribution is not known and is given by:

$$\gamma_1 = \frac{3(M_x - x_{\text{med}})}{\sigma_x} \quad \text{--- (5)}$$

Note: The factor 3 in the equation (5) is due to fact that empirical results indicate that

$$\text{Mean} - \text{Mode} \approx 3(\text{Mean} - \text{Median})$$

$$[M_x - x_{\text{mode}} \approx 3(M_x - x_{\text{med}})]$$

~~Note~~ Note: ① Pearson Median Coefficient is more popularly used.

① ②

③ Positively skewed distribution ($\gamma_1 > 0$):

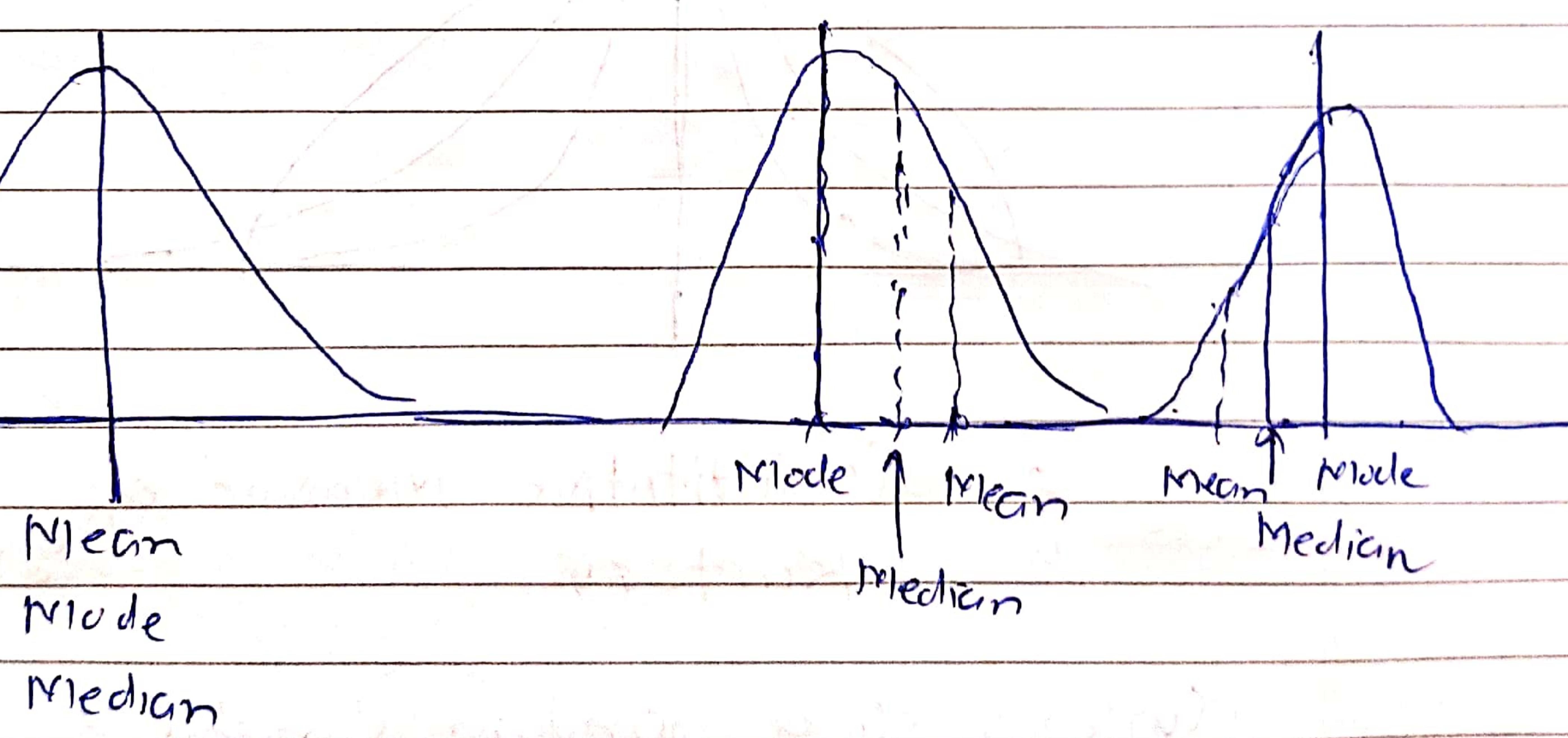
Mean > Median > Mode

④ For negatively skewed distribution ($\gamma_1 < 0$):

Mean < Median < Mode

⑤ For symmetric distribution ($\gamma_1 = 0$):

Mean = Median = Mode.



Shape of frequency distribution (Peakedness)

Kurtosis (કર્ટોઝિસ) = Refers to their degree of peakedness.

⇒ Reference is normal distribution

(i) more peaked than the normal distribution

Lepto = slender = skinny

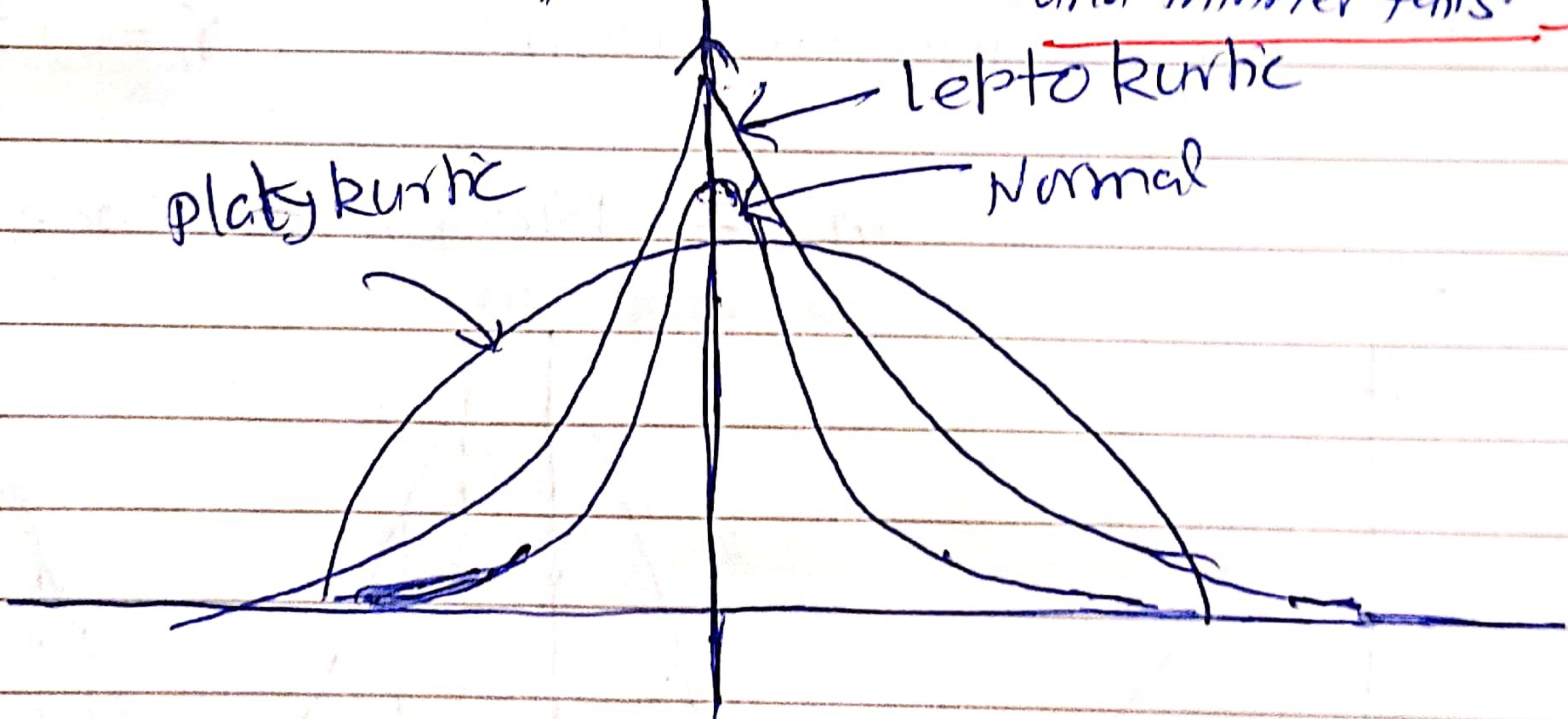


It is said to be leptokurtic; more acute peak around the mean and fatter tails.

(ii) less peaked than the normal distribution

Platy = broad "

platykurtic; wider peak around mean, and thinner tails.



Quantitative Measure of Kurtosis

Kurtosis of a random variable X is

denoted by β_2 and defined by:

$$\beta_2 = E\left[\left(\frac{X-\mu_x}{\sigma_x}\right)^4\right] = \frac{E[(X-\mu_x)^4]}{[\sigma_x]^4}$$

$$\beta_2 = \frac{1}{(\sigma_x^2)^2} E[(x - \mu_x)^4] \quad \text{--- (6)}$$

Excess Kurtosis — Excess Kurtosis, γ_2 , of the

frequency distribution is given by:

$$\gamma_2 = \beta_2 - 3$$

(L) Three (3) is subtracted from β_2 is used as

a correction to make the kurtosis of the ~~normal~~

normal distribution equal to zero.

Leptokurtic \Rightarrow distribution: Have positive excess

Kurtosis

Platykurtic \Rightarrow distribution: Have negative

excess kurtosis.

Normal distribution has excess kurtosis of zero.

Mesokurtic distribution: The distribution

that is peaked the same way as the normal distribution is called mesokurtic.

Computation of Kurtosis

Data \in Population of size N , the Kurtosis is given by

$$b_2 = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^4}{s^4} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^4}{\left[\frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \right]^2}$$

— (6)

As in the case of the skewness parameter, this

computation applies to a population of N members.

Data \in Sample of size n , the excess Kurtosis is

given by:

$$g_2 = \frac{n-1}{(n-2)(n-3)} \left\{ (n+1)b_2 + 6 \right\} \quad — (7)$$