

Statistics

Def: statistics is the science of collecting, organizing, interpreting, presenting and analyzing data.

Data: "facts or pieces of information that can be measured".

Ex: Heights of students in a classroom.

2) IQ of students.

3) Daily Activities.

4) Weight of people, Age.

Types of statistics:

① Descriptive statistics

Def: It consists of organizing and summarizing data.

① Measure of central tendency

{ Mean, Median, Mode }

② Measure of dispersion.

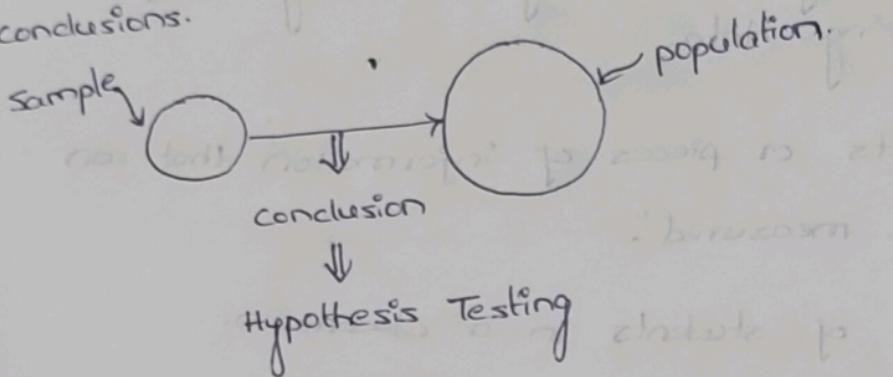
{ Variance, standard deviation }

③ Different types of distribution of data.

Ex: Histogram, pdf, pmf, cdf

② Inferential statistics

Def: It consists of data you have measured from conclusions.



- ① z-test
- ② t-test
- ③ chi-square Test
- ④ ANOVA
- ⑤ F test

⇒ conclusion of sample
on population.

Ex: Let say there are 20 classes in your college.
and you have collected the heights of student
in the class.

Heights are recorded as [175 cm, 180 cm, 140 cm, 135 cm, 160 cm]

Descriptive:

"What is the average height of the students in the classroom".

$$\text{Mean} = \frac{175 + 180 + 140 + 135 + 160}{5}$$

Inferential:

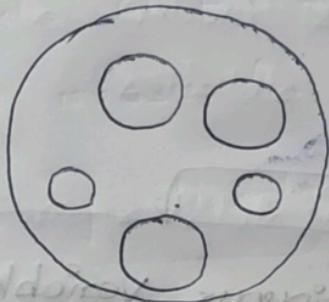
"Are the weight of the students in the classroom
similar to what you expect in the college?"

sample.

population.

population(n) and sample Data (n)

Exit poll



Population :- The entire group you want to study.

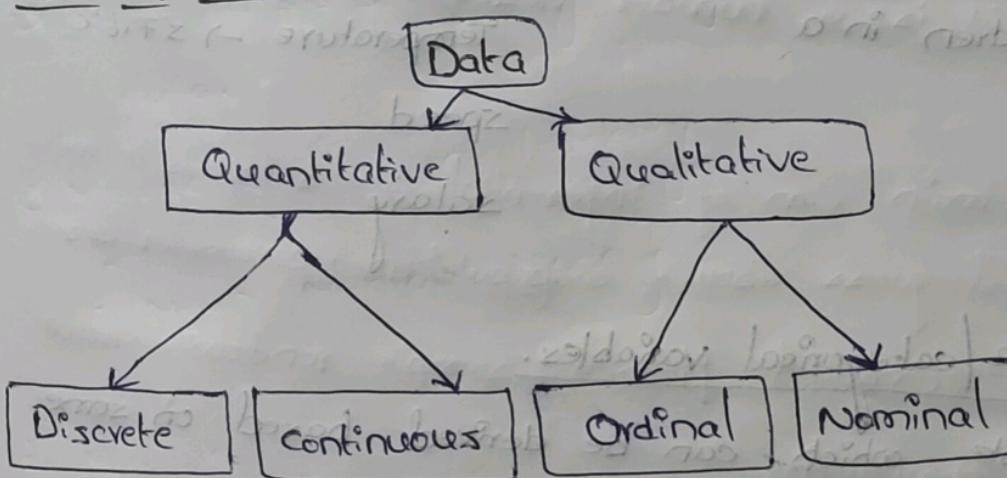
Sample :- The portion of the population you select to study.

Sample Size = 500

Population = 10000

Media people will go and ask people in different areas whom they have voted for. Based on this sample of data, they will make conclusions for population, who will going to win in that election.

Types of data:



① Quantitative Variables: Measured numerically
we can perform operations → { Add, subtract, multiply, divide }

Ex: Age

Weight

Height.

(i) discrete (ii) continuous

Discrete variable

Continuous Variable

whole numbers with some range

Ex: No. of bank accounts of people.
2, 3, 4, 5

→ It can't be 2.5, 3.5, 4.5

Ex: No. of children in a family.

Ex: Weight → 100 kg, 99.5, 99.75
Height → 172.5, 162.4 cm

Age → 22.5 years

Temperature → 37.6°C

speed

salary

② Qualitative (categorical variables).

→ The variable which can be derived based on some categorical variable.

Ex: Gender

Ex: Blood Group

A⁺ve, A⁻ve

Ex: IQ

0-10
↓
low IQ

10-50
↓
Medium IQ

50-100
↓
Good IQ

Ex: T-shirt size

L, XL, M, S

Qualitative variable

Ordinal

Rank and order of data matters, but value doesn't.

Ex:- Ranks

Good - 3

Better - 2

Best - 1

Ex:- Students (Marks)

100

96

57

85

44

1

2

4

3

5

Ordinal data

Nominal

→ is nothing but categorical data {Grouped into different classes}

→ Here nominal means there is no such specific rank.

Like Male is greater than female, one colour is greater than other colour.

Ex:- Gender

M, F

Blood Group → A⁺, A⁻

→ color of hair → Black, Green

→ pincode

population	sample
1) whole group	1) part of the group
2) Group we want to know about.	2) Group we do know about.
3) characteristics are called parameters.	3) characteristics are called statistics.
4) parameters are generally unknown.	4) statistics are always known
5) parameters are fixed.	5) statistics change with the sample.

parameter: It is a number describing a whole population

Ex: population mean.

statistic: It is a number describing a sample

Note: The main aim of inferential statistics is to make an educated guess about a population parameter based on statistic computed from a sample randomly drawn from that population.

Scales of Measurement:

① Nominal scale data

② Ordinal scale data.

③ Interval scale data.

④ Ratio scale data.

① Nominal scale data.

1) Qualitative / Categorical Data. (This type of data can be measured)

Ex: Gender, colors, Labels

2) Order does not matter.

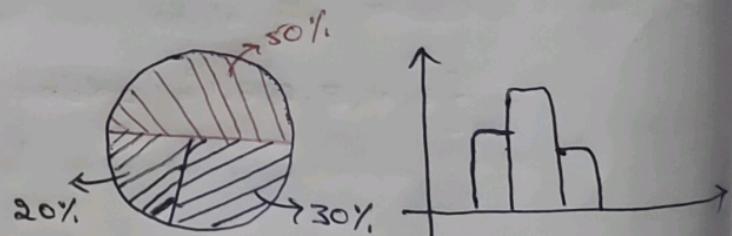
Ex:- Favourite colour

Red \rightarrow 5 \rightarrow 50%

Blue \rightarrow 3 \rightarrow 30%

Orange \rightarrow 2 \rightarrow 20%

This data can be represented through bar chart, pie chart.



② Ordinal scale data:

① Categorical data.

② Ranking and order matters.

③ Difference cannot be measured.

Ex:- { Best \rightarrow 1 }

Good \rightarrow 2

Bad \rightarrow 3

Instead focusing on these values we will mostly focus on these ordinal values while dealing with

Rank	Value
1 st	Rank \rightarrow 90
2 nd	Rank \rightarrow 70
3 rd	Rank \rightarrow 40

③ Interval data scale data :-

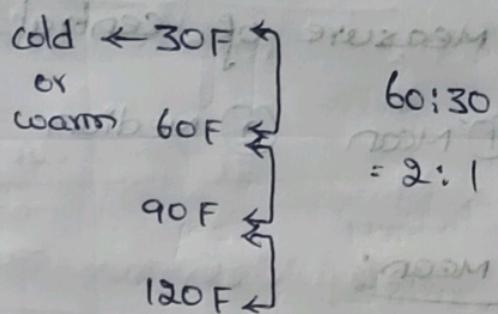
① The order matters.

② Difference can be measured.

③ Ratio cannot be measured.

④ NO "0" starting points.

Ex: Temperature.



This means that

→ suppose room 1 with 30°F

room 2 with 60°F

this doesn't mean that we will feel twice the

temperature in room 2 compared to room 1. It

will completely depends on external factors. So, Ratio cannot be measured for this type of data.

④ Ratio scale data :-

① The order matters

② Differences are measurable including Ratios.

③ Contain a "0" starting point.

Ex: Student marks in class.

↓ ↓
0, 30, 45, 60, 90, 95, 99
S1 S2 S3 S4 S5 S6 S7

Ex:
2, 4, 6, 8, 10

① Marital status [Nominal scale
nominal data : E.g. data]

② Favourite food based on Gender? [Nominal]

③ IQ measurements [Ratio scale]

Descriptive statistics

① Measure of central tendency.

- ① Mean ② Median ③ Mode.

i) Mean's

population (N)

$$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$$

$$\text{population mean}(\bar{\mu}) = \frac{\sum_{i=1}^n x_i}{N}$$

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

$$\bar{\mu} = \frac{1+1+2+2+3+3+4+5+5+6}{10} = 3.2$$

ii) Median's

$$X = \{4, 5, 2, 3, 2, 1\}$$

steps :-

① sort the random variables $\{1, 2, 2, 3, 4, 5\}$

② no. of elements

③ if count == even

$$1, 2, \boxed{2, 3}, 4, 5$$

$$\text{median} = \frac{2+3}{2} = 2.5 \text{ median}$$

④ if count == odd

$$1, 2, 2, \boxed{3}, 4, 5, 6$$

3 median

why median? → means are affected by outliers.

$$x = \{1, 2, 3, 4, 5\} \quad x = \{1, 2, 3, 4, 5, 100\}$$

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

Median = 3

$$\bar{x} = \frac{1+2+3+4+5+100}{6} = \frac{115}{6} \approx 19$$

Incase of median:

$$x = \{1, 2, 3, 4, 5, 100\}$$

$$\text{median} = \frac{3+4}{2} = 3.5$$

Note: Median is used to find the central tendency and also when outlier is present.

Outliers → can be completely different numbers compared to other distribution [due to these outliers, there is more difference or movement in the central distribution of data].

iii) Mode: maximum frequency occurring element.

$$\{2, 1, 1, 1, 4, 5, 7, 8, 9, 9, 10\}$$

$$\text{Mode} = 1$$

→ Mean, Median, Mode concepts used for EDA and Feature Engineering.

EDA and feature Engineering

	Age	Weight	Salary	Gender	Mode	categorical fea.
24	70	40K	M	B.E		
25	80	70K	F		B.E	
27	95	45K	F			
outliers	24	50	M	PHD		
Median or Mode	32	60	60K	M	B.E	
	60	-		M	Master	
	65	55K		M	B.Sc	
	40	77	-	M	B.E	

→ Missing values

→ In case of missing values of we run in a particular row, if we drop that row means there will be a loss of data.

→ Let's say I want to handle the missing values in age, weight, salary.

for that → we need to check whether there are any outliers or not.

→ If there are outliers means we will replace missing values with median.

→ If there are no outliers means we will go with mean.

→ But what categorical features in order to handle the missing values we can specifically use the mode.

Q) Measure of Dispersion?

i) Variance (σ^2)

ii) Standard deviation (σ)

(i) Variance : It measures the spread of the data with emphasis on extreme low and high values.

Population variance

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

x_i = Data points

μ = population mean

N = population size

Sample variance

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

x_i ⇒ data points

\bar{x} ⇒ sample mean

n ⇒ sample size

Q) Why we divide sample variance by $n-1$?

$$Ex: \{1, 2, 3, 4, 5\}$$

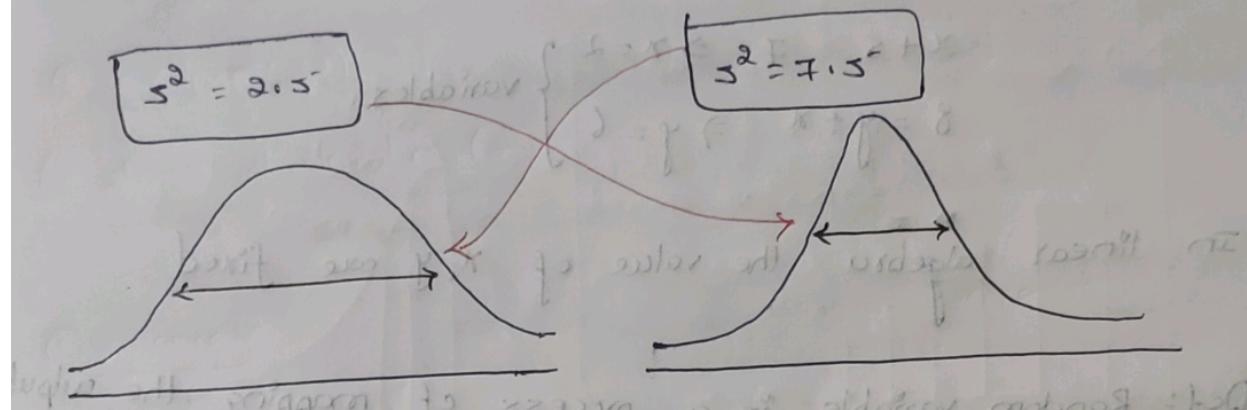
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

x_i	\bar{x}	$(x_i - \bar{x})^2$
1	3	4
2	3	1
3	3	0
4	3	1
5	3	2
<hr/>		$\sum (x_i - \bar{x})^2 = 10$
$\bar{n} = 3$		

$$s^2 = \frac{10}{4}$$

$$x = \{ \quad \}$$

$$y = \{ \quad \}$$



ii) Standard deviation: It measures the dispersion of a dataset relative to its mean and is calculated as the square root of the variance.

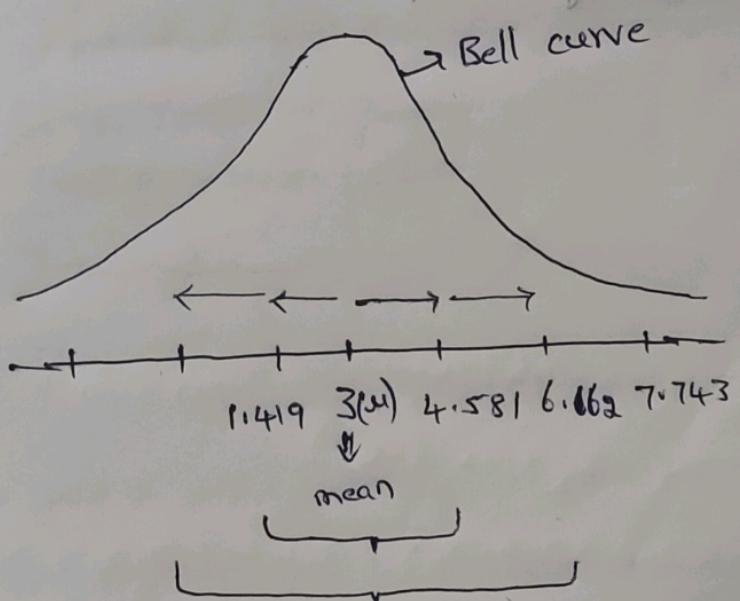
$$\text{population standard deviation } (\sigma) = \sqrt{\text{variance}} = \sqrt{\sigma^2}$$

$$\text{sample standard deviation } (s) = \sqrt{s^2}$$

$$\text{midland } \mu = \{1, 2, 3, 4, 5\}$$

$$\bar{x} = 3$$

$$\sigma = 1.581$$



Random Variable :-

$$\begin{aligned} x+5 = 7 &\Rightarrow x=2 \\ 8 = y+x &\Rightarrow y=6 \end{aligned} \quad \left. \begin{array}{l} x=2 \\ y=6 \end{array} \right\} \text{variables.}$$

In linear algebra \rightarrow the value of x, y are fixed.

Def: Random variable is a process of mapping the outcome of a random process or experiment to a number.

Ex: Tossing a coin $\{ \text{Head, Tail} \} \rightleftharpoons$ process \rightarrow to

$$X = \begin{cases} 0 & \text{if Head} \\ 1 & \text{if Tail} \end{cases}$$

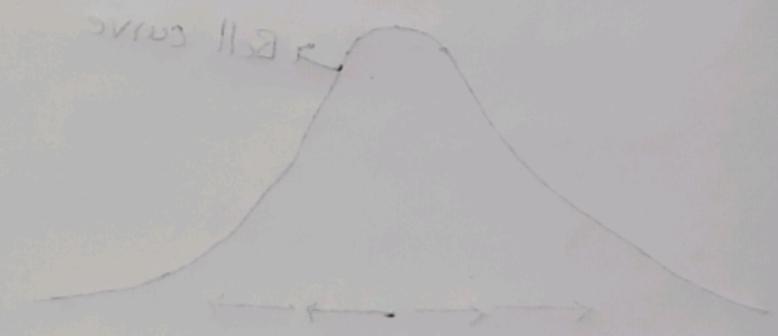
Ex: Rolling a dice $\{1, 2, 3, 4, 5, 6\}$.

$$Y = \{\text{sum of rolling a dice 7 times}\}$$

\rightarrow Random process will help us to solve few questions like $\rightarrow \Pr(Y \geq 15)$?

$$\rightarrow \Pr(Y < 10) ?$$

$$\rightarrow \Pr(10 \leq Y \leq 15) ?$$

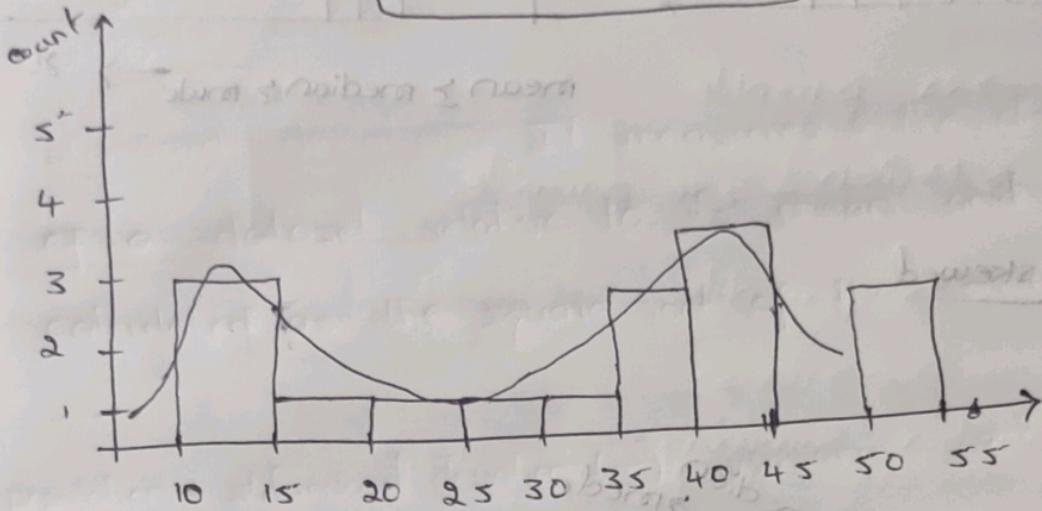
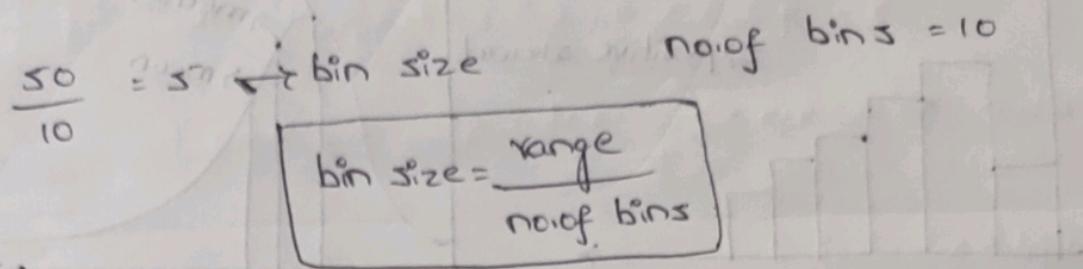


Ex: $\Pr(18 \leq Y \leq 19)$

Answer

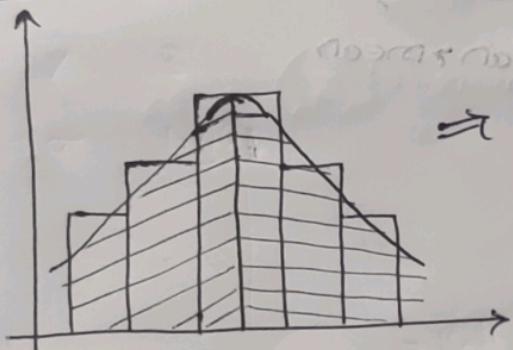
Histogram

Ages = { 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51 }



Skewness:

① NO skewness



⇒ Normal distribution

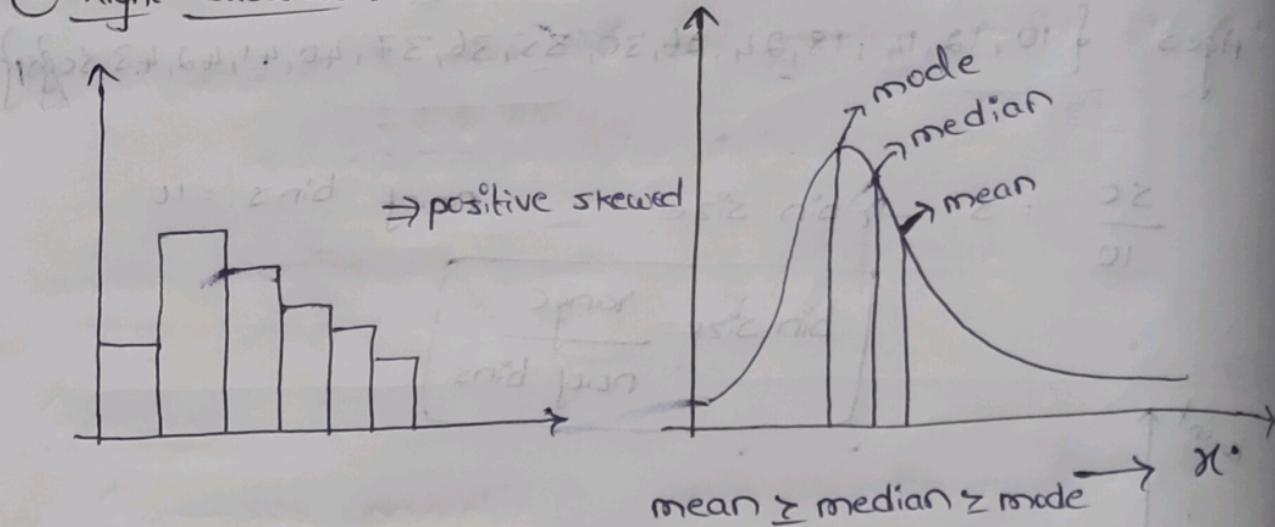
⇒ symmetrical distribution



NO skewness

Median = Mean = Mode

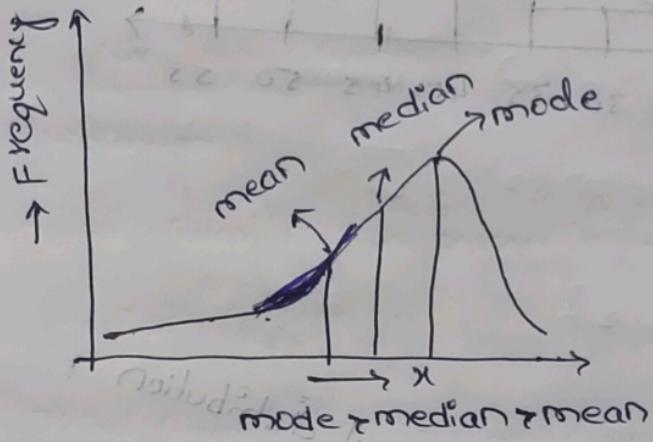
② Right skewed:



Ex: Length of youtube comments

income data set (most people having average incomes, smaller no. of people having higher incomes.)

③ Left skewed



Ex: Test scores of a group of college students who took a relatively simple exam.

(Most of the students have high test scores, and a smaller number of people would have low scores)

that skew the curve toward the left of the group)