

Statistics \Rightarrow Method of collecting, organizing and analysing data in such a way that meaningful conclusion can be drawn.

Why?

data \rightarrow fact | information that can be stored and measured.

ML

Learn Patterns

Example \rightarrow Scores made by Virat Kohli in last five matches

$\checkmark [92, 80, 40, 70, 60]$

Ex.2. Height of students in the classroom

$\checkmark [150, 160, 170\text{cm}, 180, 190\text{cm}]$

Motivation behind Statistics

- ① Weather forecast \longrightarrow 45°C
- ② Sports analysis \longrightarrow SR - VK.
- ③ Election campaign \rightarrow Exit Polls \rightarrow BJP
- ④ PMS | e-commerce \longrightarrow D MART
- ⑤ medical | genetics

Types of Statistics

Descriptive

Describe

It consists of

summarising and organising data.

Inferential

\rightarrow

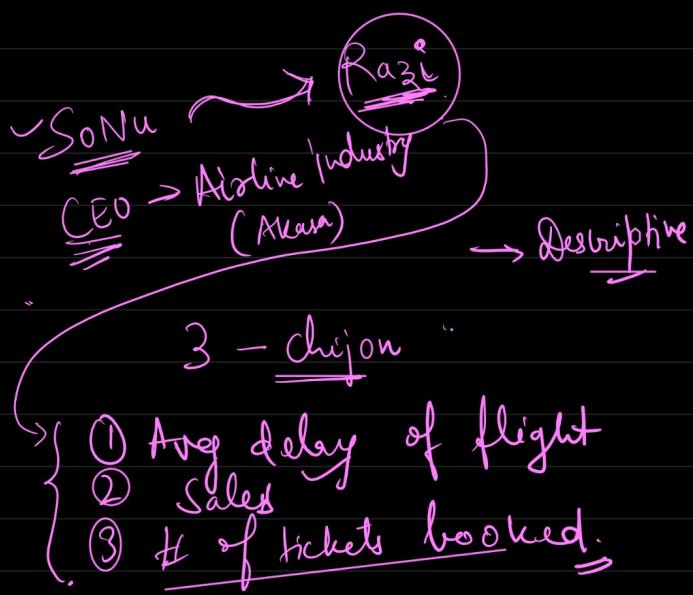
Inference

Using sample data from a population

→ Complete population

and calculate statistics

but can also be
calculated for sample.



Sakshi

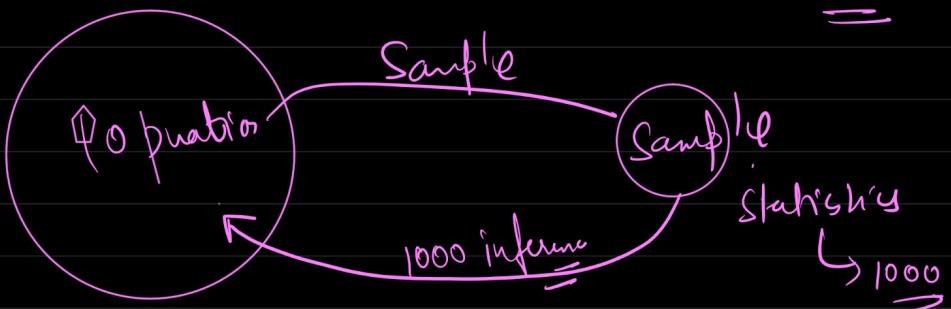
T.

Traveling → Kaziranga
National park

Sample

No of trees

$$\text{Inferential } 1000 \text{ m}^2 \rightarrow \underline{\underline{1000}}$$
$$100000 \text{ m}^2 - 100000 \text{ trees} =$$



100m

Dipam

Khalid = 80kg
180kg

Wt.

inferential stats

* Class:

Description - Avg wt?

- most frequent wt?

- Strike rate of VK?

Descriptive

① Measure
of CT

② Measures
of Symmetry

measure of
dispersion
=

* Inferential

Avg wt of people of India



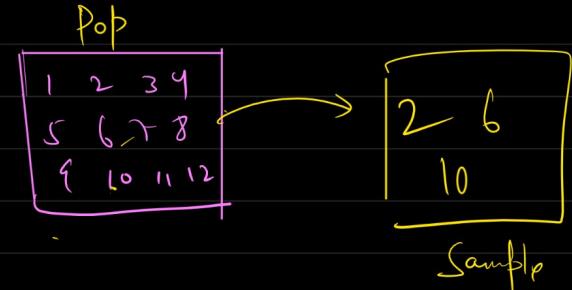
NE

{ Sample
wt → India
represent?

Types of Sampling

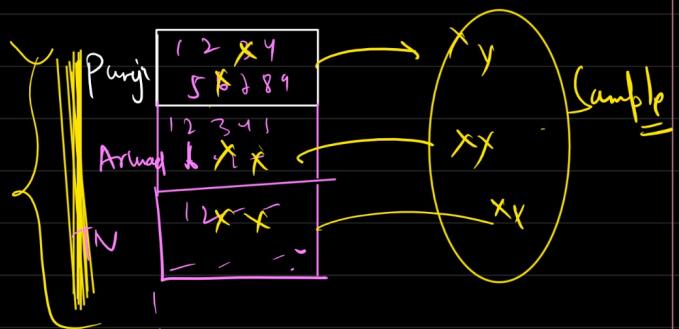
 ① Simple random Sampling.

 Each member of pop. has an equal chance of being selected $\frac{1}{n}$.



Q Stratified Sampling

Strata | layer | groups



③ Cluster Sampling

→ divides the population into groups / clusters. Some of these clusters are randomly selected.

Covid days

CovidIndia.org



④ Systematic Sampling

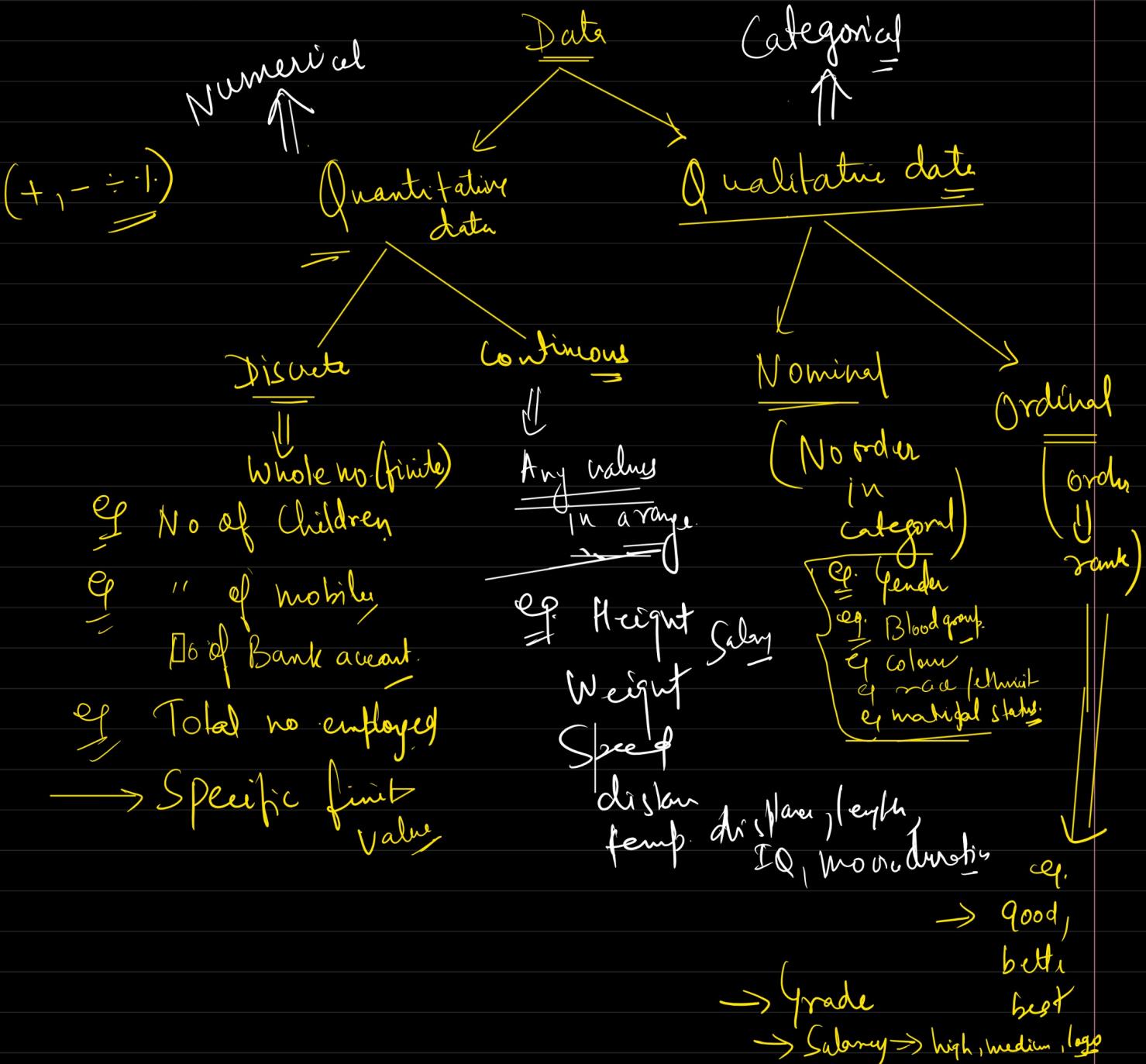
Every n^{th} element.

	1	2	3	4
5 -	8	7	8	9
10	11	12	13	

- * Quota Sampling
 - * Min Max Sampling
 - * Convenience Sampling
 - * Accidental Sampling

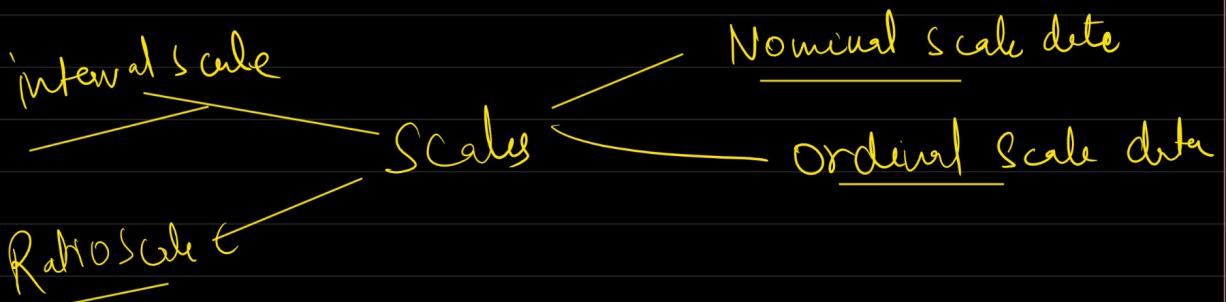
* Types of data

Calculus



Scales of measurement

→ Military ranks



① Nominal scale data

→ Qualitative / Categorical data

→ Gender, ethnicity, colour, location.

→ No order in the data

Example
Employees

M
F
M
F
M
M
M
M

M - 6
F - 2

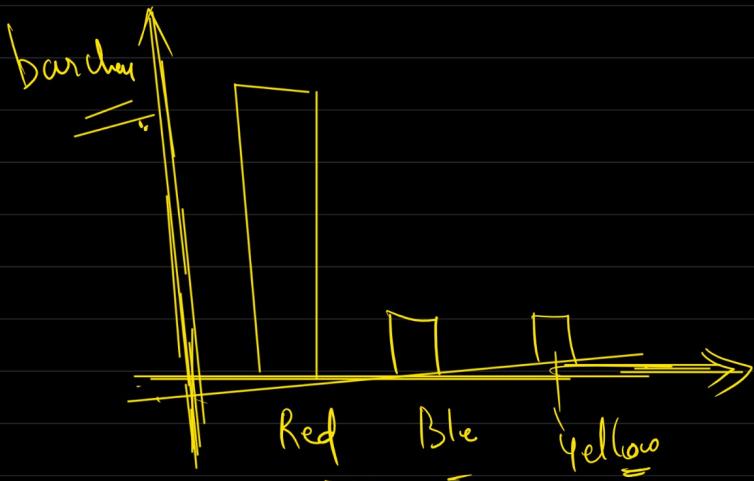
* Count =

Red - 5 → 50%

Blue - 3 → 30%

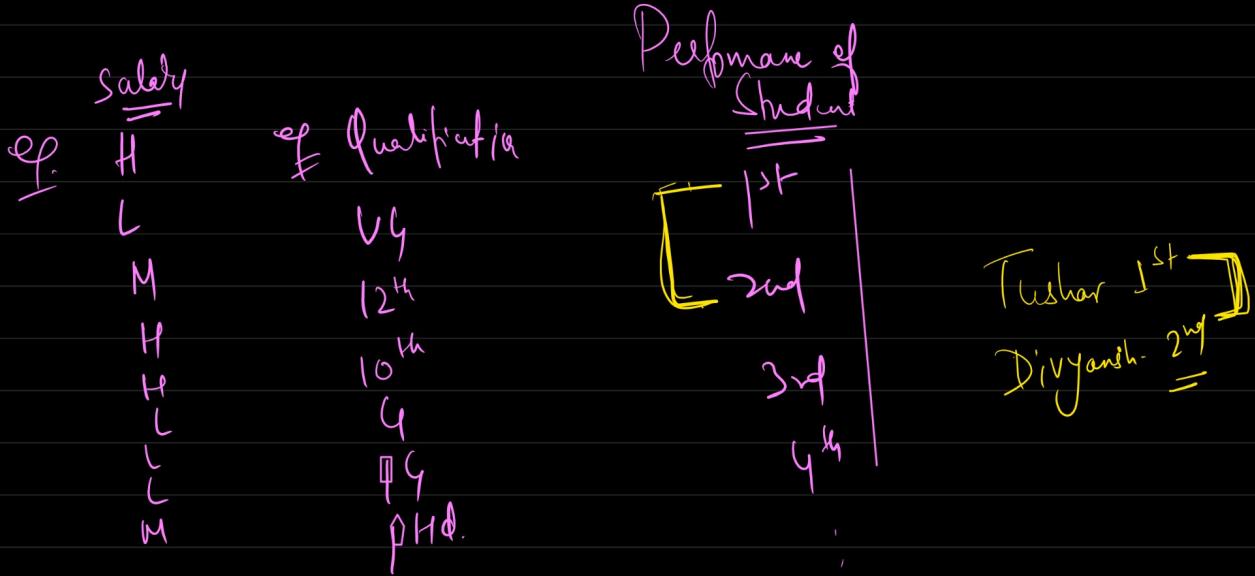
Orange - 2 → 20%

Pie-chart



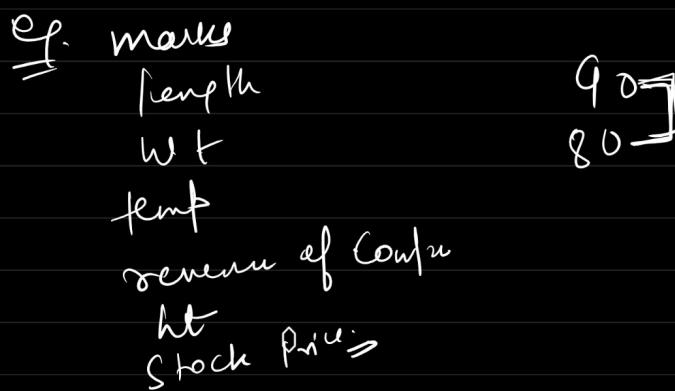
② Ordinal Scale data

- Order / rank matters.
- Difference can not be measured.



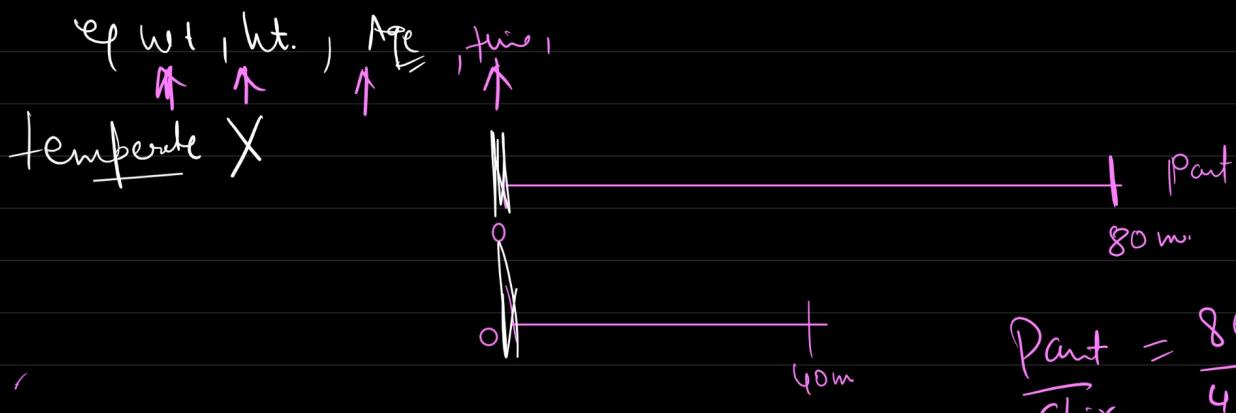
③ Interval Scale data

- rank & order has a meaning
- Difference can be measured (excluding ratio)
- It doesn't have a starting value.



④ Ratio Scaled data

- Order / rank has meaning
- Difference and ratio is measurable
- It does have a starting point.



Temp.

-5, 0, 5°
inside 30°C
outside - 60°C
-8, -10, -20

$$\frac{30}{60} = \frac{1}{2}$$

Temp can be negative

Data	Nominal	Ordinal	Interval	Ratio
------	---------	---------	----------	-------

labelled ✓ ✓ ✓ ✓

order ✗ ✓ ✓ ✓

differs ✗ ✗ ✓ ✓

true zero
start at 1st
fail ✗ ✗ ✗ ✓
lit

Example → gender, salary, post office code, location, satisfied, rating, grade, rank, military rank, score, lit, wt, ht, time, age.

* Descriptive Statistics (Summarization of data)

- ① Measure of central tendency
- ② Measure of dispersion
- ③ " " Symmetry.

* MCT

Central \Rightarrow What is that one value around which all the data is revolving.



CT \rightarrow represents the center point of a data.

- ① mean
- ② Median
- ③ Mode

} EDA and feature engineering.

① Mean

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

$$\frac{1+2+3+4+5}{5}$$

$$x = 1, 2, 3, 4, 5 \\ \sum \sim 1+2+3+4+5$$

$$(M) \sum x_i$$

$$\bar{x} = \frac{\sum x_i}{n}$$

$$x = 1, 2, 3, 4, 5$$

$$\bar{x} = \frac{\sum x_i}{n}$$

$$\frac{1+2+3+4+5}{5}$$

$$x = (8, 9, 7, 6, 1)$$

$$\sum x_i = 8+5+6+7$$

$$\text{mean} =$$

$$\frac{\sum x_i}{n}$$

$$1, 2, 3, 4, 5$$

$$\frac{x_1 + (1) + (1) + (1) + (1)}{5}$$

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

② Median (Physical mid point)

4, 5, 2, 3, 1, 2

* Sort the data $\rightarrow 1, 2, 2, 3, 4, 5$.

Count the no of element = 6

if count is even \rightarrow avg of two middle element.

$$\begin{array}{c} 1 \\ | \\ 2 \\ | \\ 2 \\ | \\ 3 \\ | \\ 4 \\ | \\ 5 \end{array}$$

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

$$\left| \begin{array}{l} 4, 5, 2, 3, 1 \\ \text{Sort} = 1, 2, 3, 4, 5 \\ \text{Count} = 5 \\ \text{Odd} \\ \text{median} = \frac{1+3+4+5}{4} \\ \underline{\underline{3}} \end{array} \right.$$

.mean(), .median()

$$\begin{array}{c} \text{Sort} \\ \Rightarrow \\ 1 \\ | \\ 2 \\ | \\ 3 \\ | \\ 4 \\ | \\ 5 \\ | \\ 6 \\ | \\ 7 \\ | \\ 8 \\ | \\ 9 \\ | \\ 9 \end{array}$$

$\text{avg} = \frac{5+6}{2} = 5.5$

$$\begin{array}{c} \text{or} \\ | \\ 1 \\ | \\ 2 \\ | \\ 3 \\ | \\ 3 \\ | \\ 4 \\ | \\ 5 \\ | \\ 6 \\ | \\ 6 \\ | \\ 7 \\ | \\ 8 \\ | \\ 9 \end{array}$$

S is mean

Scenario-1

1, 2, 3, 4, 5

mean = 3

→ 1, 2, 3, 4, 1000 outlier

$$\text{mean} = \frac{1+2+3+4+1000}{5} = \frac{1010}{5} = 202$$

Scenario-2

1, 2, 3, 4, 5

median = 3

Sorting → 1, 2, 3, 4, 1000
median = 3

* The mean is affected by outliers whereas median is not affected
by outliers

③ mode ⇒ Categorical data

maximum frequency.

{1, 1, 1, 2, 2, 3, 3, 1, 1, 2, 3}

col

good - 1

better - 2

best - 3

mode → 1

* Use of Central tendency

Numerical → mean, median /
Categorical → mode / if outlier present

⇒ Age
25
26

Gender
M-1
male

Weight
80
70

Scary(k)
—

median

23

Mode
||

70

50

Age

continuous num
dt

mean

$$\frac{25+26+23+25}{4}$$

4

23

M-1
Highest
freq

30

60

25

M-1
med

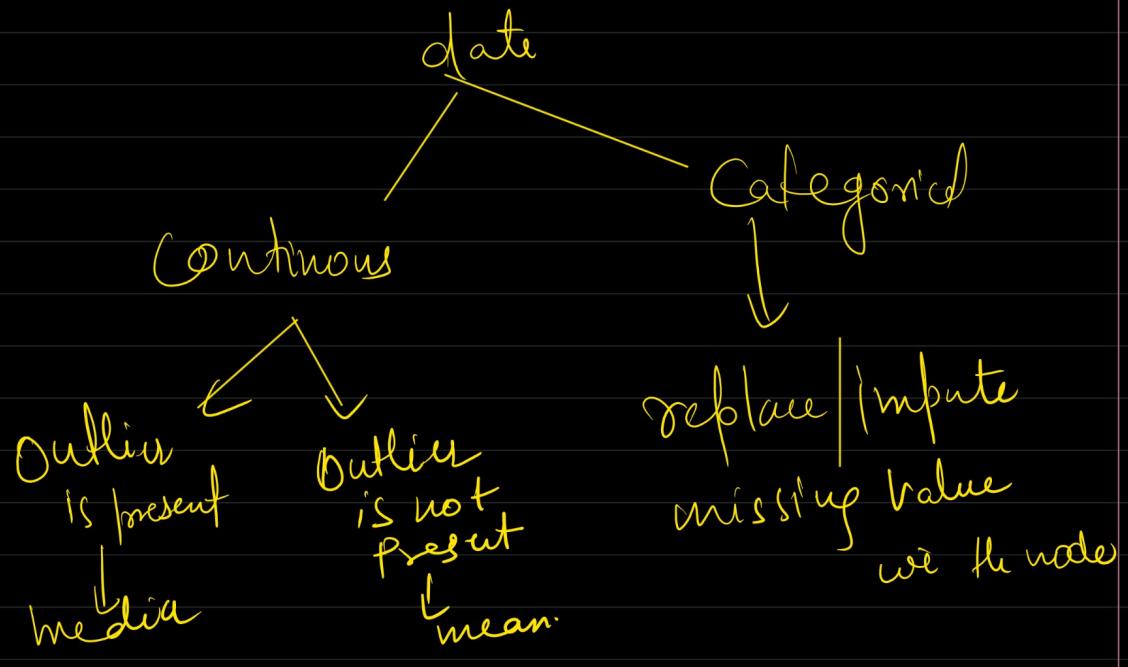
—

70

Fo

60

45



* Measures of dispersion / spread

$$\underline{S_1} = 1, 2, 3, 4, 5 \quad , \quad \underline{S_2} = 3, 3, 3, 3$$

mean | median, 3, 3.

3, 3.



* Measure of spread?

- Range
- Percentage & percentile
- Quartiles (Box plot)
- Variance
- Std dev^n

① Range \rightarrow difference b/w max and min value

$$\begin{array}{c} \{1, 2, 3, 4, 5\} \\ \hline \text{Range} = 5 - 1 = 4 \end{array}$$

Range is affected by Outlier

$$\boxed{\{1, 2, 3, 4, 5\}} \quad | \quad \textcircled{1000} \\ 1000 - | = 999$$

Inter Quartile Range

* Percentage

{ 1, 2, 3, 4, 5 }

What is the percentage of Nos that are odd.

$$\frac{3}{5} \times 100^{\circ} = \underline{\underline{60^{\circ}}}$$

* Percentile: A percentile is a value below which a certain percentage of observations lie.

$$Abir = \frac{99}{100} \rightarrow$$

$$Abir = 99\% \text{ et } \frac{99}{100}$$

Score

$\frac{99}{100}$

Percentile rank of a no = No of values below than no $\times \frac{100}{n}$

$$\frac{2}{10} \times 10^6 = 20 \underline{\underline{J}} \quad \text{Total nos (n)}$$

20. If the nos are below 3 =

* What values exist at 75th percentile?

$$\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)$$

$$= \frac{75}{100} \times \underline{\underline{(10+1)}}$$

$$= \frac{3}{4} \times 11 = \underline{\underline{8.25}}$$

$$\begin{aligned} & 3 \\ & 5 \\ & \frac{75+10}{100} = \frac{15}{2} \\ & +_2 = 7.5 \\ & \frac{7m}{2} \\ & \underline{\underline{= 7.5}} \\ & 8^{th} \text{ no.} + 9^{th} \text{ no.} \\ & \underline{\underline{= 7.5}} \end{aligned}$$

$$75^{\text{th}} \text{ percentile} = \underline{\underline{7}}$$

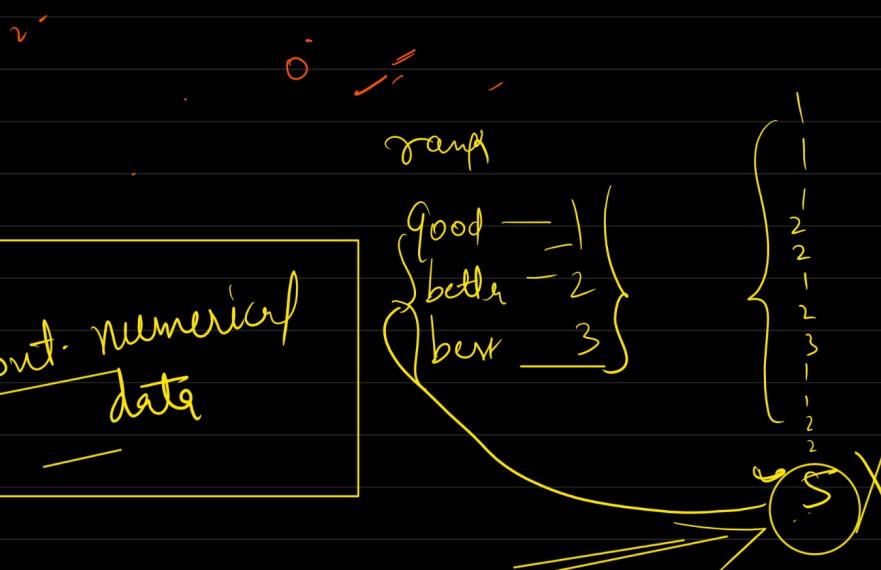
* mode - highest frequency

$\downarrow \sim \text{male}$

Col.	1	Male
1	1	Male
1	1	Male
0	0	Fem
0	0	Fem
		Mal



Outlier \rightarrow cont. numerical data



UP-200

J

2

Nagaland

50 Laks

J

2

=

S.Ru