

3. MODE

unimode

no mode

multiple mode

SURVEY → Favourite FRUIT among

APPLE, BANANA, ORANGES :-

∴ MODE = FRUIT mentioned ↑ Frequently

- MOST frequent value.
- Applicable Best FOR CATEGORICAL VARIABLES.
- In CASE OF NUMERICAL VARIABLES → Unique Count SHOULD be less.

$$\text{E.g. } \textcircled{2} \text{ } X = \{2, 4, 3, 4, 2, 6, 7, 2\}$$

⇒ $\textcircled{2} \equiv \text{MODE}$.

II MEASURES OF DISPERSION

→ also known as SPREAD or VARIABILITY.

- It refers to statistics that describes how data points in a dataset are DISTRIBUTED OR SPREADOUT.

→ Measures of Dispersion like

- RANGE
- INTERQUARTILE RANGE (IQR)
- VARIANCE
- STANDARD DEVIATION, help us understand how much data values diverge or spread from central values like mean or median.

∴ ↑ SPREAD ⇒ ↑ VARIABILITY ⇒ ↓ CONSISTENCY OR PREDICTIBILITY

- Higher variability can be useful in some scenarios and less desirable in others.

1. RANGE

- Range = $\{ \text{maximum value of dataset} - \text{minimum value of dataset} \}$

↳ SENSITIVE to EXTREME VALUES / OUTLIERS.

$$\begin{aligned} \text{E.g. } & \{1, 2, 3, 4\} \\ \text{Range} &= 4 - 1 \\ &= 3 \end{aligned}$$

$$\begin{aligned} \text{E.g. } & \{1, 2, 3, 4, 50\} \\ \Rightarrow \text{Range} &= 50 - 1 \\ &= 49 \end{aligned}$$

↳ Does not give any information on how data is distributed around the mean.

↳ Suitable for small datasets.

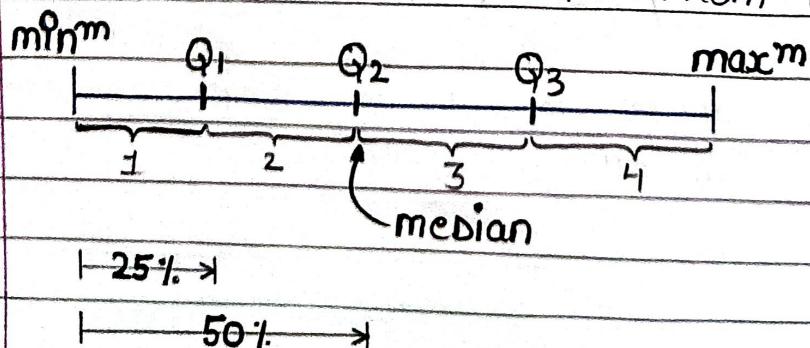
2. QUARTILES [Q₁, Q₂ & Q₃]

- Quartiles are values that divide a dataset into 4 equal parts each containing a quarter of the data. They help describe the distribution.

↳ FIRST QUARTILE, Q₁ - Value that separates lowest 25% of data

↳ SECOND QUARTILE, Q₂ - MEDIAN or the midpoint of the data.

↳ THIRD QUARTILE, Q₃ - Value that separates the lowest 75% of data from highest 25%.



- For finding quartiles try to arrange your dataset in increasing order.

Page No.:

Date:

$$\rightarrow Q_1 = \frac{(n+1)}{4}^{\text{th}} \text{ term}$$

$$\text{E.g. } x = (4, 6, 7, 8, 10, 23, 34)$$

$$\rightarrow Q_2 = \frac{(n+1)}{2}^{\text{th}} \text{ term}$$

$$\bullet Q_1 = (7+1)/4 = (8/4) \equiv Q_1 = 6$$

$$\rightarrow Q_3 = \frac{3(n+1)}{4}^{\text{th}} \text{ term}$$

$$\bullet Q_2 = (7+1)/2 = 4^{\text{th}} \text{ Term} \equiv Q_2 = 8$$

$$\bullet Q_3 = (3(7+1)/4) = 6^{\text{th}} \text{ Term} \equiv Q_3 = 23$$

3. PERCENTILES :

→ Percentiles divide a dataset into 100 equal parts.

For example the 25th percentile is the value below which 25% of the data falls.

→ QUARTILES = 4 Equal parts

PERCENTILES = 100 Equal parts.

POSITION, $P = \frac{p}{100} \times (n+1)$

• p = percentile

100

• n = no. of observations.

E.g. $x = 4, 6, 7, 8, 10, 23, 34$

find 5 percentile of the given dataset.

$$\text{We know: Position, } P = \frac{p}{100} \times (n+1) = \frac{5}{100} \times (7+1) = \frac{1 \times 8}{20} \\ = 0.4$$

This indicates a position just before the 1st data point. Therefore, the 5% percentile is closest to the lowest value which is 4.

NOTE :-

To find the value of a datapoint that falls between two datapoints, we would interpolate between them.

For Example :- 2.5 position value means the percentile value is halfway between 2nd and 3rd data point, so take average of these two datapoints (value) to get the answer.

For Example :- If the position value comes out to be 2.7? This means we're at 70% of the way between the second and third data point.
Here's how to calculate percentile :-

Let say,

Value at position 2 = V_2
Value at position 3 = V_3

Then the interpolated value at position 2.7 is =

$$V = V_2 + 0.7(V_3 - V_2)$$

$$\Rightarrow V = 6 + 0.7(10 - 6)$$

$$= 6 + 0.7 \times 4$$

$$= 8.8$$

∴ The percentile value would be around 8.8.

4. IQR (Interquartile Range)

$IQR = Q_3 - Q_1 = 50\% \text{ of data from our dataset}$

→ IQR is less sensitive to extreme values as it finds the central 50% of data leaving first and last 25% called as extreme values.

5. VARIANCE

→ Variance measures how spread out the values in a dataset are from the mean.

→ Variances are the average of the squared differences from the mean.

→ FORMULA :-
For a dataset with value $x_1, x_2, x_3, \dots, x_n$ and mean μ .

- POPULATION VARIANCE ;

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

- SAMPLE VARIANCE

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

where,

- $\mu \equiv$ Population mean

- $\bar{x} \equiv$ Sample mean

- $n \equiv$ no. of data points.

x_i

* Example : Dataset :- [4, 7, 10]

- Mean, $\mu = \frac{4+7+10}{3} = 7$

- Squared Difference $\Rightarrow (4-7)^2 = 9$

- $(7-7)^2 = 0$

- $(10-7)^2 = 9$

$$\therefore \text{Population Variance} = \sigma^2 = \frac{(9+0+9)}{3} = \frac{18}{3} = 6 \text{ (Ans)}$$

6. STANDARD DEVIATION (Std)

- $S.D. = \sqrt{\text{Variance}}$

- The unit of standard deviation is same as of the dataset while the unit of variance is squared of dataset. That's why most of the time we prefer calculating S.D over variance.

FREQUENCY

↳ no. of times a value of data occurs.

↳ We prefer to calculate FREQUENCY in case of CATEGORICAL VARIABLES.

↳ We create FREQUENCY DISTRIBUTION TABLE.

* Ex :- {1, 3, 3, 2, 4, 1, 2, 2, 1, 2, 3, 5, 4, 1, 2, 1, 3, 1, 4, 1}

DATA VALUE	FREQUENCY	RELATIVE FREQUENCY
1	7	7/20
2	5	5/20
3	4	4/20
4	3	3/20
5	1	1/20
	21	21/20

RELATIVE FREQUENCY :-

→ Percentage or position of data values present in complete dataset.

→ $\frac{\text{FREQUENCY}}{\text{Total number of observations}} \times 100$

CUMULATIVE FREQUENCY :-

• Cumulative frequency is the running total of frequencies in a dataset.

Example: Customer complaints at a company —

Imagine working in a customer service department of an e-commerce company and you've collected data on types of complaints received over a month.

You have categorized the complain into main 4-types :-

Shipping Delays - 15 complaints

Product Quality - 10 complaints

Billing Issues - 8 complaints

Returns - 12 complaints

∴ CUMULATIVE FREQUENCY TABLE :-

COMPLAINT TYPE	COMPLAINT FREQUENCY	CUMULATIVE FREQUENCY
SHIPPING DELAYS	15	15
PRODUCT QUALITY	10	$15+10 = 25$
BILLING ISSUES	8	$25+8 = 33$
RETURNS	12	$33+12 = 45$