(7)

# Measure of Central Tendency.

A measure of CT is a __single value__ that attempts to describe a set of data by identifying the __central position__.

The three main types : -

1) **Mean - (Average)** The mean is the Sum of all values divided by the no. of Values.

formula

$$Mean = \frac{Sum\ of\ all\ values}{No.\ of\ values}$$

eg ⇒ n = 70, 75, 80, 85, 90

$$= (70+75+80+85+90)/5$$
$$= \frac{400}{5} = 80$$

• **Mean w.r.t population (μ - mu) N**
  The entire group you're Studying.

formula

$$population\ Mean\ (\mu) = \sum_{i=1}^{N} \frac{x_i}{N}$$

eg ⇒ N = (24, 23, 2, 1, 28, 27)

$$\mu = \frac{24 + 23 + 2 + 1 + 28 + 27}{6}$$

$$\mu = 17.5$$

• **Mean w.r.t Sample (n) x̄**
  A part of population you actually collect data from.

formula

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

eg ⇒ n = (23, 2, 28, 27)

$$\bar{x} = \frac{23 + 2 + 28 + 27}{4}$$

$$\bar{x} = \frac{80}{4} = 20$$

Mean tells you the "typical" or "central" value in a data set.

eg) # Data Set (for mean eg.)

| Age | Salary | family size |
|-----|--------|-------------|
| 24 | | |
| 26 | | |
| NAN | | |
| 21 | | |
| 20 | | |
| 18 | | |

☆ NAN → not a number → Empty
How to handle NAN / Empty
1) Delete → loss of info.
2) ignore
3) fill it (Avg.
4) Drop if too may nans.

In this example we are fill it with Avg.

$$Avg = \frac{24 + 26 + 21 + 20 + 18}{5} = 21.8$$

Now the NAN values replaced by 21.8

• We use Mean where there is no outlier in dataset

2> Median - (Middle value) The median is the middle value of a dataset when it's arranged in order

Range = (5-1) = 4

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

$$\iff$$

Range = (100-1) = 99

$$en - \{1, 2, 3, 4, 5, 100\}$$ — Outlier

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

$$M = \frac{1+2+3+4+5+100}{6}$$

$$= \frac{115}{6} = 19.16$$

Outlier => It is a number that is complete different than the entire distribution.

+ Steps to find the median.    $\{1, 2, 3, 5, 4, 100\}$

1) Sort the no.    $(1, 2, 3, 4, 5, 100)$

2) find the central no.    $(1, 2, \boxed{3, 4} 5, 100)$
   if the no. of elements are even we find the avg of central no.s,  $\frac{3+4}{2} = \frac{7}{2} = 3.5$

   if the no. of elments are odd we find central no.  en => $(1, 2, \boxed{4} 5, 100)$
   => 4

we used Median when there is outlier in data.

(9)

3> Mode — Most frequent ~~app~~ appeard or occured element.

Example- Dataset ≠ Categorical variable

Types of flowers

Lily
Sunflower
Rose
Lily
Rose
(NAN)
Rose
Rose
Lily
Lily

So no mean and median used

$\Downarrow$

Mode => Lily or Rose
      4         4

$\Downarrow$

Replace with Something.

\* ~~we used Mode when the data have catogorical replacement,~~

——— ||| ———

# Measure of Dispersion.

It shows how much the data varies or spreads out from the center (like mean or median)

Why is it imp in DA => Because avg. alone can lie, you need dispersion to know if the data is stable or chaotic.

Ex=> $X = \{1,1,2,2,4\}$          $Y = \{0,2,2,3,3\}$
$\mu = \dfrac{1+1+2+2+4}{5}$  <=>  $\mu = 2.$
$= 2$

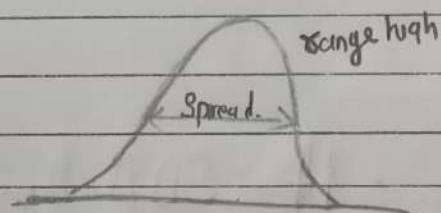these are same? No, but how it will
Know by the technique Variance.

Types of Measure of Dispersion.

1) Variance
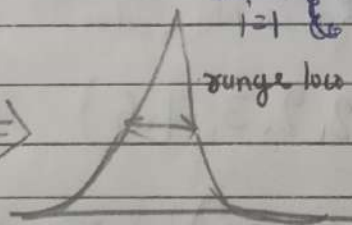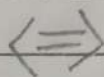
It measures the avg. of the squared difference from the mean. (spread of distribution).

formula —

population (N) variance $(\sigma^2) = \sum\limits_{i=1}^{N} \dfrac{(x_i - \mu)^2}{N}$

→ population mean $(\mu)$



range high

Spread

⟺

range low

Variance highest          Variance low

• if range is high variance will high and vice versa.

formula —

Sample (n) Variance $(S^2) = \sum\limits_{i=1}^{n} \dfrac{(x_i - \bar{x})^2}{n-1}$

→ Bessel's correction
degree of freedom.

Bessel's correction = Correcting the bias in estimating population variance from a sample.

Use n-1 because when you work with a Sample you are guessing — so you adjust for that guess with a tiny corrections,

2) Standard Deviation (SD)

It tells us how much data values differ from the mean — in the same units as the data.

It is literally Jus the square root of variance

formula -

for population (N)                    for Sample (n)

$$\sigma = \sqrt{\sum_{ni=1}^{N} \frac{(x_i - \mu)^2}{N}}$$     $$S = \sqrt{\sum_{ni=1}^{n} \frac{(x_i - \bar{x})^2}{n-1}}$$

Example :— for sample data

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

| X | $\bar{x}$ | $(x - \bar{x})$ | $(x - \bar{x})^2$ |
|---|---|---|---|
| 1 | 2.83 | -1.83 | 3.34 |
| 2 | \| | -0.83 | 0.6889 |
| 2 | = | -0.83 | 0.6889 |
| 3 | \| | 0.17 | 0.03 |
| 4 | \| | 1.17 | 1.37 |
| 5 | \| | 2.17 | 4.71 |

· if SD↑ then dispersion or spread↑

· if SD↓ spread↓

Small SD ⇒ data is close to mean (consistent)

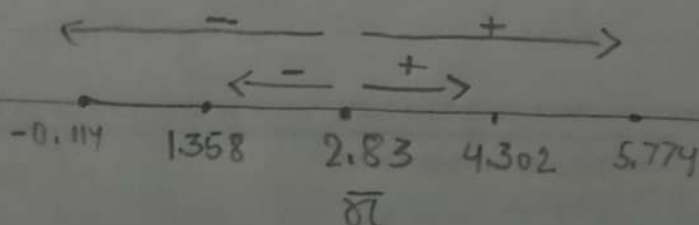Large SD ⇒ data is spread out (more varient)

$$S^2 = \frac{10.84}{5} = 2.168$$

$$S = \sqrt{2.168} = 1.472.$$

2.83 ($\bar{x}$)
+1.472 (S)
———
4.302
+1.472
———
5.774
~~+~~

2.83
-1.472
———
1.358
-1.472
———
-0.114



-0.114   1.358   2.83   4.302   5.774
                  $\bar{x}$

## Range :-

Range = Highest value - lowest value

It tells you the total spread of data

$$en - X = \{12, 18, 25, 30, 40\}$$
$$Man = 40 \quad Min = 12$$
$$Range = 40 - 12 = 28$$

So the data values stretch across 28 units

- Range is super sensitive with outliers. one extreme value can totally mess it up.

————||•||————

## Percentile And Quartile

Percentile - $\{1, 2, 3, 4, 5, 6, 7, 8\}$

percentage of even no. = $\dfrac{No. \; of \; even \; no.}{total \; no. \; of \; nos.} = \dfrac{4}{8} = \dfrac{1}{2}$
$$= 0.5 = 50\%.$$

percentage of odd no. = $\dfrac{4}{8} = \dfrac{1}{2} = 0.5 = 50\%$

Percentile - we have seen Gate, SAT CAT JEE etc gave
result in percentile.

A percentile is a value below which a certain percentage
of observation lie.

Example - if A person in 99 percentile, It means the
person has got better marks then 99% of the
entire Students.

- Data set Example (Sorted)

$$[2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, \underline{10, 11, 11}, 12]$$
                                                        Same

✱What is percentile ranking of 10?

percentage Rank of $x = \dfrac{\text{# of values below } n}{n} \times 100$

$$10 = \dfrac{\cancel{16}^{\,8}}{\cancel{20}} \times \cancel{100} = 80 \text{ percentile}$$

10 is greater than 80% of entire distribution.

for 11,     $11 = \dfrac{17}{\cancel{20}} \times \cancel{100}^{\,5} = 85 \text{ percentile}$

✱What is the value that exists at 25 percentile

formula Value $= \dfrac{\text{percentile}}{100} \times (n+1)$

$$= \dfrac{\cancel{25}^{\,1}}{\cancel{4}\cancel{100}} \times (20+1)$$

$$= \dfrac{20}{4} = 5 = \text{Index}$$
                    between $5^{th}$ and $6^{th}$ position

for 55 percentile

$$= \frac{55^{th}}{100\%} \times (20 \text{ +i}) = \frac{(20)}{20} = 11 \text{ } ^{th} \text{ positive}$$

for 40 percentile

$$= \frac{40^\circ}{500} \times 20 = 8^{th} \text{ position}$$

What if you index is in decimal value.

$5.5^{th}$ position

blw then $5^{th}$ and $6^{th}$ position i.e 5 and 5 in ques.

Avg. $\frac{5+5}{2} = \frac{10}{2} = 5.$

## Five Numbers Summary.

1) Minimum

2) Q1 first quantile (25th %)

3) Median (Q2 / 50th %)

4) Q3 third quantile (75th %)     } used to remove outliers

5) Maximum

(15)

Dataset

$\{\dot{1}, 2,2,2,3,3,4,5,5,5,6,6,6,6,7,88,9 \underset{\text{outlier}}{(27)}\}$

$> -3?$            $< 13$

Quartile - divide a sorted dataset into 4 equal parts.

Inter Quartile Range IQR - how spread out the 50% of data
                              • b/w Q3 and Q1

$$IQR = Q3 - Q1$$

Lower Bound $= Q1 - 1.5 \times (IQR)$

Higher Band $= Q3 + 1.5 \times (IQR)$

$$Q1 = (25\%) = \frac{25}{100} \times \overset{5}{\cancel{20}} = 5^{th} \text{ position} = 3$$

$$Q3 = (75\%) = \frac{\overset{15}{75}}{\cancel{5100}} \times \cancel{20} = 15^{th} \text{ position} = 7$$

$$IQR = Q3 - Q1 = 7 - 3 = \underline{\underline{4}}$$

$LB = 3 - 1.5 \times 4 = -3$
$HB = 7 + 1.5 \times 4 = 13$

After removing outlier, the remaining value will be

$\{1, 2,2,2,3,3,4,5,5,5,6,6,6,6,7,8,8,9\}$

1) Minimum = 1     2) Q1 = 3     3) Median/Q2 = 5

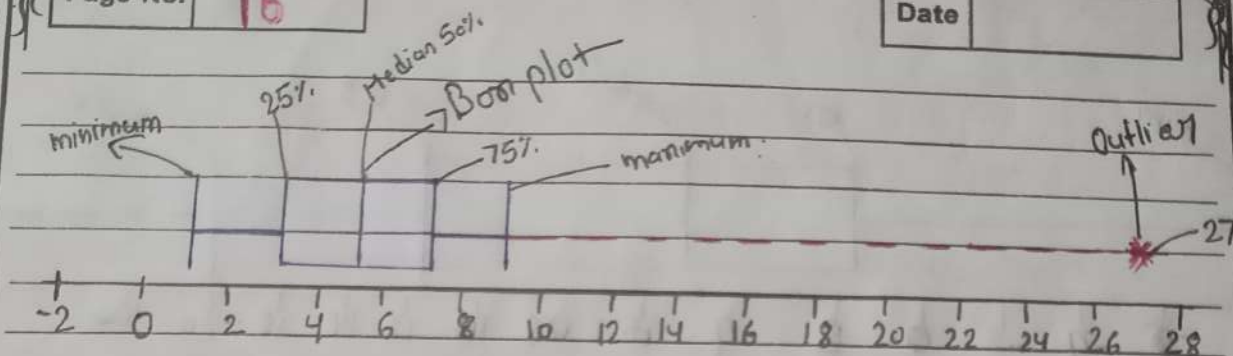4) Q3 = 7       5) Maximum = 9

fig Box plot (identifies outliers).

## Dataset

$$\{ -8, 1, 2, 4, 5, 6, 8, 15, 20, \boxed{120} \}$$

→13.5                                        ←26.5

$$Q_1 = (25\%) = \frac{25}{2 \text{to} 100} \times 10 = \frac{5}{2} = 2.5^{th} \text{ position}$$

$$= 2^{nd} + 3^{rd}$$

$$= \frac{1+2}{2} = \frac{3}{2} = 1.5$$

$$Q_3 = (75\%) = \frac{75}{100} \times 10 = 7.5^{th}$$

$$= 7^{th} + 8^{th} = \frac{8+15}{2} = \frac{23}{2} = 11.5$$

$$IQR = Q_3 - Q_1 = 11.5 - 1.5 = 10$$

$$LB = 1.5 - 1.5 \times 10 = -13.5$$
$$HB = 11.5 + 1.5 \times 10 = 26.5$$

After removing outliers

$$\{ -8, 1, 2, 4, 5, 6, 8, 15, 20 \}$$

(17)

1) Minimum = -8   2) Q1 = 1.5   3) Median = 5   4) Q3 = 11.5   5) Max = 20



mini   Q1   median   Q3   man

-8  -6  -4  -2  0  2  4  6  8  10  12  14  18  20  100        120

Hum