

X Day - 4 X

X ML with Python X

X Statistics X

Types

- X Descriptive
- X Inferential

example - Screw manufacturing

1M

500 → quality (sample)

which parts will win the election

500 to 1000 → Sample

1M → Population

X Descriptive → 1m row or 1M → 20

X Inferential sample → except infer about population  
population  
Population

X Descriptive →

- ① measures of central tendency
- ② measures of dispersion
- ③ covariance correlation

3M's (mean, median, mode)

Range, quartiles,  
SD, variance,

## \* Mean \*

- ① 5 4 2 1 3  $15/5 = 3$   
 ② 5 4 2 5 3 1  $515/6 = 85.83$  *So outlier*  
 ③ 1 2 3  $6/3 = 2$  *(extreme)*  
 M.F. M.F.  $\Rightarrow ?$  mean

Data  $\rightarrow$  missing values Numerical  $\Rightarrow$  missing values  $\rightarrow$  mean  
 (Blank, ?, 0) categorical  $\rightarrow$  if there are no outliers

## \* Median \*

- rules  $\Rightarrow$  ① Sort it  
 ②  $n$  (no. of items)  
 ③  $n$  is odd  $\Rightarrow$  median =  $\frac{1}{2}(n+1)$ th term  
 ④  $n$  is even  $\Rightarrow$  median =  $\frac{1}{2}(\frac{n}{2} \text{th} + (\frac{n}{2} + 1) \text{th})$  terms

5 3 1 4 2

$$\begin{array}{c} 1 2 3 4 5 \\ \hline n = 5 \text{ (odd)} \\ \frac{n+1}{2} = 3 \text{rd term} \end{array}$$

$$\text{median} = 3$$

5 3 1 4 2 1 3 5

$$n = 6 \text{ (even)}$$

$$\frac{1}{2} \left[ \left( \frac{6}{2} \right) \text{th} + \left( \frac{6}{2} + 1 \right) \text{th} \right]$$

$$\frac{1}{2} (3^{\text{th}} + 4^{\text{th}})$$

= 1 2 3 4 5 6

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

## \* Mode \*

6 0 0 R R 0 G G G mode = 9  
 1 2 3 3 3 2 3 mode = 3

## \* Measures of Dispersion

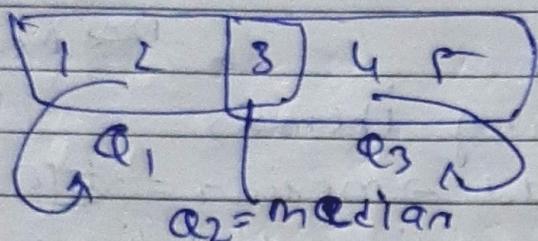
\* Range  $\rightarrow$  Max - Min

$$\text{Ages} = 1, 2, 3, \dots, 20$$

$$20 - 1 = 19$$

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

Quantiles



$$Q_1 = 2$$

$$Q_2 = 3$$

$$Q_3 = 4$$

① Short if

② Median

③  $Q_2 = \text{median}$

④ Exclude

\* Outliers  $\Rightarrow$

$$IQR = Q_3 - Q_1$$

$$< Q_1 - 1.5 * IQR$$

$$> Q_3 + 1.5 * IQR$$

$$IQR = Q_3 - Q_1 = 4 - 2 = 2$$

$$< 2 - 1.5 * 2 = -1$$

$$> 4 + 1.5 * 2 = 7$$

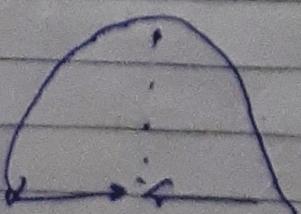
< -1

> 7

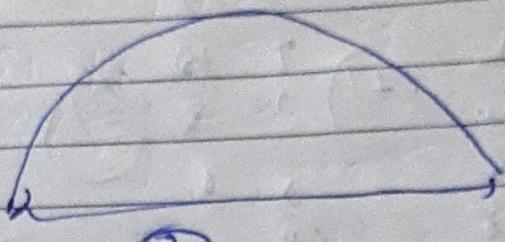
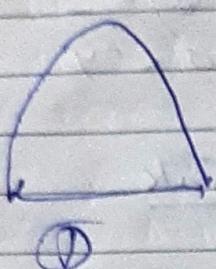
\* General Tendency  $\Rightarrow$  Mean, Median, Mode } Center

Normally distribution.

Mean = Median = Mode



measure of spread



(2) - more spread  
more Std deviation  
Variance

### Covariance & correlation

$X$	$Y$	$(x-\bar{x})(y-\bar{y})$	Covariance = $\frac{1}{n} \sum (x-\bar{x})(y-\bar{y})$
1	2	(1-3)(2-6)	
2	4	(2-3)(4-6)	
3	6	(3-3)(6-6)	
4	8	(4-3)(8-6)	
5	10	(5-3)(10-6)	

Correlation =  $\frac{\text{Covariance}}{\sigma_x \sigma_y}$

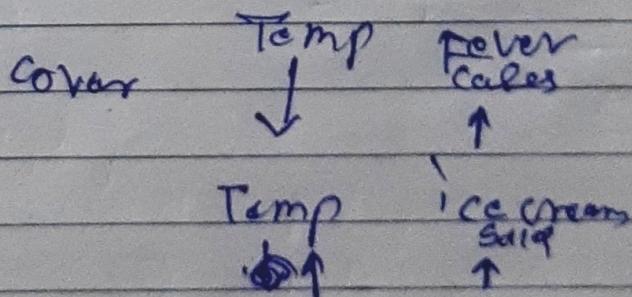
$\Rightarrow$  Std deviation

$$\text{mean} = \overline{y}$$

$$= 4$$

$X$	$Y$
1	10
2	8
3	1
4	9
5	2

$$\text{mean} = 4$$



## Correlation

X	Y	$\bar{x}$	$\bar{y}$	$\text{cov}$
1	2	4	10	$\sigma_{x-y}$
2	4	7	8	(7)
3	6	5	6	<u>Salary Sank</u>
4	8	4	4	31
8	10	8	2	

Price of Hour      Age of human      Age of buyer      sqft of floor

$$\text{cov} = 0.02$$

$$-0.08$$