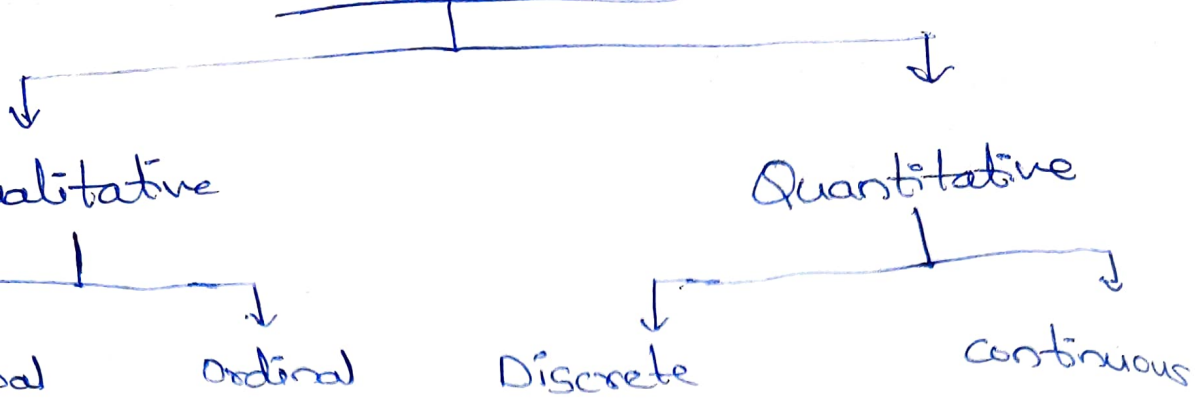


# Types of Data



| Qualitative            | Quantitative |
|------------------------|--------------|
| Category data          | Numbers      |
| Men, Women             |              |
| Blue, Black            |              |
| American, Non American |              |

Nominal! Men, Women  
Blue, Black  
American, Non American

Ordinal! Poor, Rich, Middle (Hierarchy defined)  
Poor, Average, Outstanding

Discrete! 10, 20 (Not having float values)  
11, 13  
Ex! No. of people in Room

Continuous: Weight of people (can be int, float)

## Descriptive Statistics

Measures of Central Tendency:-  
a) Mean      b) Median      c) Mode

1.1) Mean: Average of the data

Ex 1: 10, 11, 13, 17, 19, 23

$$\Rightarrow \frac{10+11+13+17+19+23}{6} = 15.5$$

~~Ex 2~~

Ex 2: 10, 11, 13, 17, 19, 23, 24, 500 <sup>↑ outlier</sup>

$$\Rightarrow \frac{10+11+13+17+19+23+24+500}{8} = 77.125$$

Outliers: Extremely high or low value

\* Problem with mean is that it is affected hugely in the case of outliers present in the data.

## 1.2 Median - Exact middle of a sorted data

7, 9, 13, 11, 10, 5, 4

\*a) Sort the data

b) Take the middle value

$$\left(\frac{n+1}{2}\right)$$

4, 5, 7, (9), 10, 11, 13 → odd data  
↓  
Median

Ex 2: 4, 5, 7, (9, 10), 11, 13, 14

→ even data

$$\frac{9+10}{2} = 9.5$$

$$\frac{\left(\frac{n}{2}\right) + \left(\frac{n}{2} + 1\right)}{2}$$

### Advantage of Median

→ It is robust in the case of outlier present in the data.

10 11 13 (17) 19 23 24

10 11 13 (17 19) 23 24 500

$$\frac{17+19}{2} = 18$$

1-3 Mode :- Most frequently occurring value in a data.

P<sub>1</sub> Yes

P<sub>2</sub> Yes

P<sub>3</sub> Yes

P<sub>4</sub> No

P<sub>5</sub> Yes

P<sub>6</sub> No

Central value

Mode = Yes

Two mode  
Bimodal  
Three mode  
Multimodal

Measure of Spread/Dispersion

- a) Range    b) Variance    c) Standard Deviation
- d) Median absolute Deviation    e) Percentile
- f) Quartiles    g) Inter Quartile Range

a) Range → Max-Min

|    |     |
|----|-----|
| 10 | 5   |
| 11 | 50  |
| 14 | 55  |
| 16 | 105 |
| 17 | 300 |
| 19 | 550 |

①  $19 - 10 = 9$

②  $550 - 5 = 545$

Outlier 10, 11, 12, 13, 14, 15, 16, 17, 110

$110 - 10 = 100$

\* Effected by the presence of outlier in the data.

## 2.2 Variance $\rightarrow$ Spread of data around mean

$$\text{Var} = \sum \text{datapoint} - \text{mean}$$

10, 11, 13, 17, 19, 23

① Mean  $\rightarrow 15.5$

② datapoint - mean  $\Rightarrow$

$$\begin{aligned} 10 - 15.5 &= -5.5 + \\ 11 - 15.5 &= -4.5 + \\ 13 - 15.5 &= -2.5 + \\ 17 - 15.5 &= 2.5 + \\ 19 - 15.5 &= 4.5 + \\ 23 - 15.5 &= 7.5 = 2 \end{aligned}$$

$$= \sum (\text{datapoint} - \text{mean})^2$$

$$\begin{aligned} \Rightarrow & (-5.5)^2 + (-4.5)^2 + (-2.5)^2 + (2.5)^2 + (4.5)^2 \\ & + (7.5)^2 \\ = & \end{aligned}$$

Because of squaring each value, there will be explosion of value



$$\text{var} = \frac{\sum (\text{datapoint} - \text{Mean})^2}{\text{No. of points}}$$

①.3 Standard Deviation  $\rightarrow$

$$\text{SD} = \sqrt{\text{Variance}}$$

Q) Are variance & SD affected by outliers?

A) Yes, because of presence of mean in calculation of var & SD.

①.4 Median Absolution deviation  $\rightarrow$

$$\frac{\sum (\text{datapoint} - \text{median})}{N}$$

\* It is not effected by outliers

①.5 Percentile

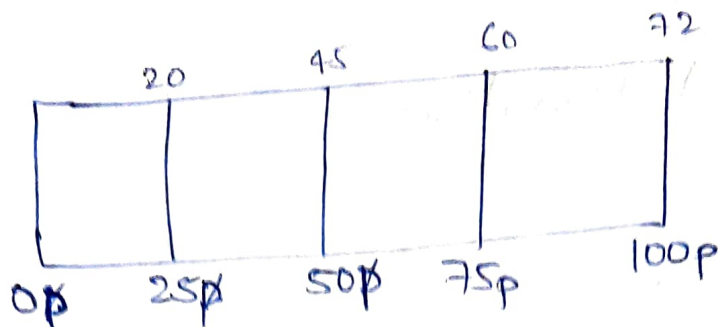
$$\text{Stud} = \frac{61}{100} = 61\% \rightarrow 99 \text{ percentile}$$

1% students who score  $> 61\%$



## 1.6) Quartiles

$$25p = 20$$



Step 1  $\rightarrow$  Sort the data

25% person score less than 20  
75% " " more than 75

100% is scored by 72

## 1.7) Inter Quartile Range ||

|       |          |       |     |
|-------|----------|-------|-----|
| 25%   | 25%      | 25%   | 25% |
| 25%   | 50%      | 75%   |     |
| $Q_1$ | $Q_2$    | $Q_3$ |     |
|       | (median) |       |     |

$$IQR = Q_3 - Q_1$$

\* IQR helps us remove outliers from the

data.

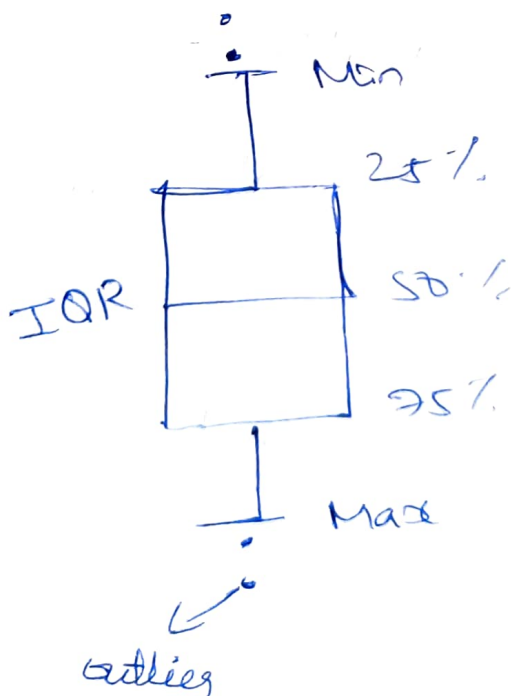
1, 2, 3, 50, 52, 55, 61, 64, 67, 510, 515, 520

① Sort the data.

$$IQR = Q_3 - Q_1$$

|       |     |                   |
|-------|-----|-------------------|
|       | 1   |                   |
|       | 2   | → lower outliers  |
| $P_1$ | 3   |                   |
|       | 50  | $Q_1$             |
| $P_2$ | 52  |                   |
|       | 55  | $Q_2$             |
| $P_3$ | 61  |                   |
|       | 64  |                   |
|       | 67  | $Q_3$             |
| $P_4$ | 510 |                   |
|       | 515 |                   |
|       | 520 | → higher outliers |

Boxplot → Five point summary of the data





Lower Outliers  $\leq Q_1 - 1.5 * IQR$

Upper Outliers  $\geq Q_3 + 1.5 * IQR$

Boxplot helps us detect outliers & also helps remove outliers.

Sometimes it may remove genuine data

## Bivariate Analysis

① Covariance

② Correlation

Covariance: Measures the direction of relationship between two variables

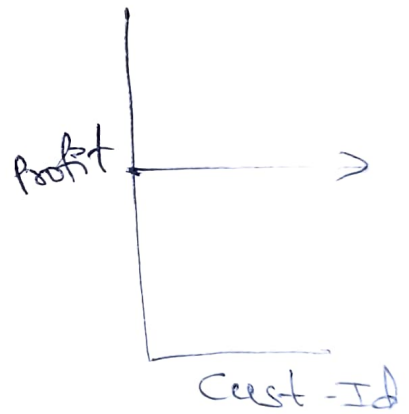
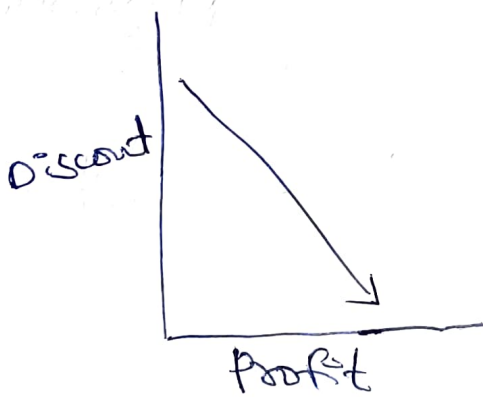
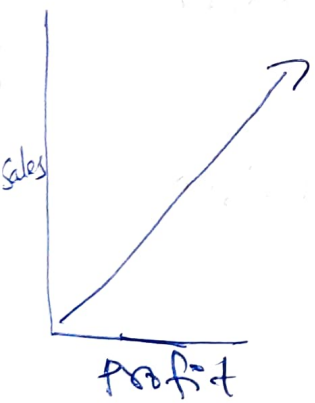
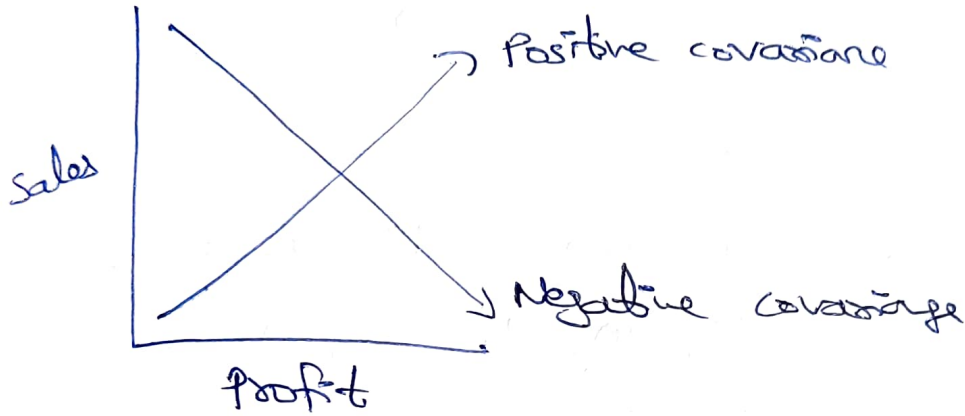
\* Positive Covariance: Means that both variables tend to be high or low at the same time.

\* Negative Covariance: Means that when one variable is high, the other tends to be low.

## Covariance Formula :-

For Populations

$$\text{Cov}(X, Y) = \frac{\sum (x_i - \bar{X}) * (y_i - \bar{Y})}{N}$$



~~Cov~~

$$\text{Cov} = \frac{\sum_{i=1}^n (x_i - \bar{X}) (y_i - \bar{Y})}{N}$$

$\bar{X}$  = mean of X data

$\bar{Y}$  = mean of Y data

Sales (X)

Profit (Y)

150

10

170

25

190

35

210

60

240

80

$X_i = 1, 2, 3, 4, \dots$

$x_1 = 150, x_2 = 170, x_3 = 190, \dots$

$y_1 = 10, y_2 = 25, y_3 = 35, \dots$

①  $(x_i - \bar{x})$

$(y_i - \bar{y})$

-

+

} → Negative Relationship

②

+

-

③

+

+

} → Positive Relation

④

-

-

the  
Positive

the  
Negative

If Sales ↑

Profit ↑

→ positive relationship

Sales ↓

Profit ↓

→ " "

Sales ↑

Profit ↓

→ Negative

Sales ↓

Profit ↑

→ Negative

$$\Rightarrow \bar{X} = 192, \bar{Y} = 42$$

$$\begin{aligned} \Rightarrow & (150 - 192) * (10 - 42) = -32 + \\ & 170 - 192 = -22 * 25 - 42 = -17 + \\ & 190 - 192 = -2 * 35 - 42 = -7 + \\ & 210 - 192 = 18 * 60 - 42 = 18 + \\ & 240 - 192 = 48 * 80 - 42 = 38 + \\ & = \end{aligned}$$

$$\begin{aligned} & (-42 * -32) + (-22 * -17) + (-2 * -7) + (18 * 18) + (48 * 38) \\ & = 3880 \end{aligned}$$

$$\Rightarrow \frac{3880}{5} = 776 \text{ (positive Relationship)}$$

Correlation (Pearson Correlation Coefficient):

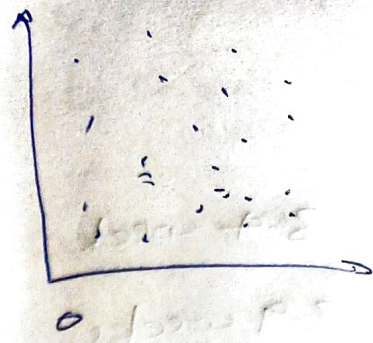
Measure the direction and strength and  
Direction of relationship b/w two variables

- Positive Correlation
- Negative Correlation
- No correlation.



$$r_{xy} = \frac{S_{xy}}{S_x S_y} \Rightarrow$$

$S_{xy}$  is covariance  
 $S_x$  &  $S_y$  are the standard deviations

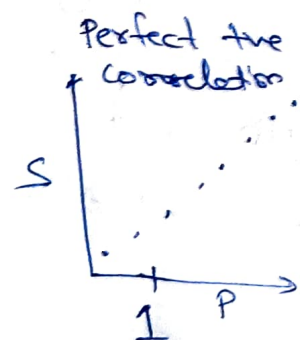


Correlation coefficient  $r$  is number between -1 to +1 and tells us how well a regression line fits the data

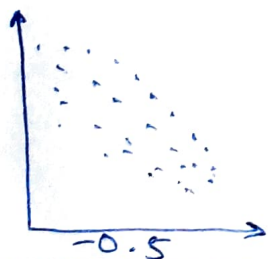
$$\text{Corr} = \frac{\text{Cov}(X, Y)}{SD(X) * SD(Y)}$$

-1 to +1  
 (-1 to +1)

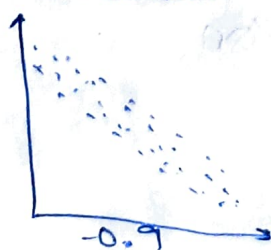
If a number is close to +1  $\rightarrow$  +ve relationship  
 close to -1  $\rightarrow$  -ve relationship  
 close to 0  $\rightarrow$  weak relationship



Low -ve correlation



High -ve correlation



Perfect -ve correlation

