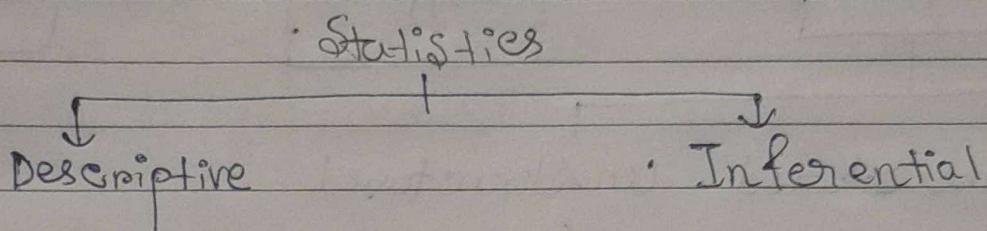


# Statistics

1. Definition: Statistics is a branch of applied maths that involves the collection, description, analysis, drawing and inference of conclusions from quantitative data.

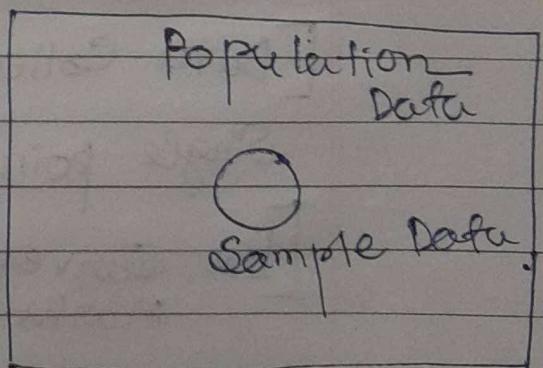


① Descriptive Statistics: All about summarizing and describing the data you have.

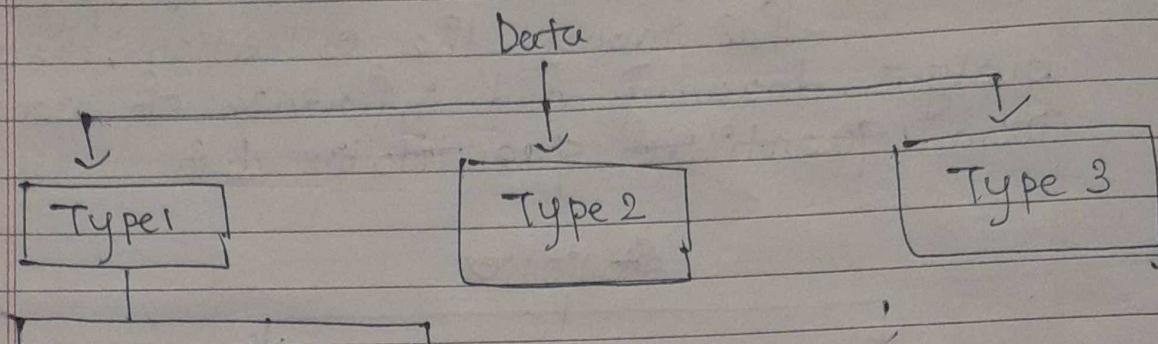
Example: Imagine you have a bunch of test scores from a class of students. Descriptive stat helps you to summarize the data and visualize the data.

② Inferential statistics: It is about making predictions or inferences about a larger group (population) based on the data you have. (sample)

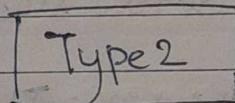
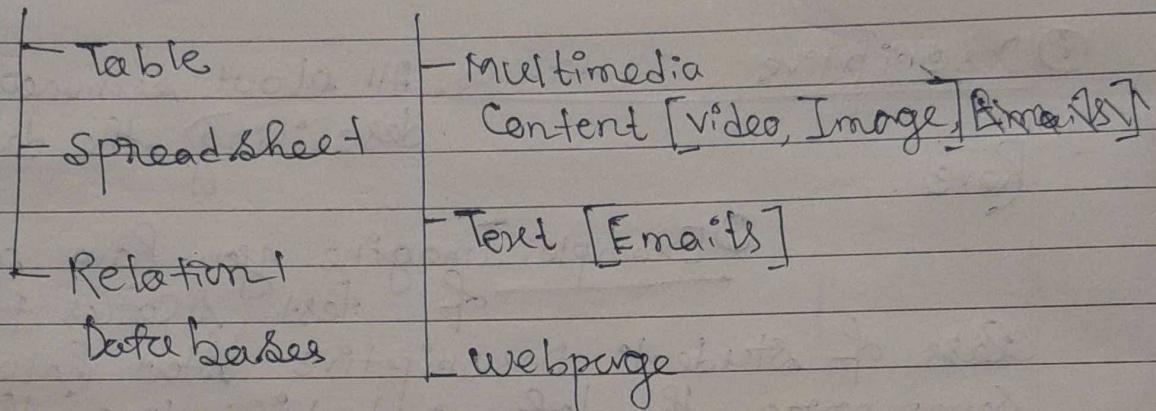
↳ Confidence Interval  
Estimation  
Hypothesis Testing



## 2. Types of Data:



Structured      Unstructured

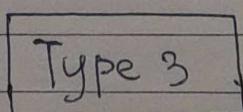


Cross Section

[Data Collected from at a Single point of time]  
[Ex: Survey, Market]

Time Series

[Data Collected over a sequence of time interval]  
[Ex: Daily stock price, Monthly Sales data]



Univariable

[Single variable]

Multivariable

[Two or more than 2 variables]

### 3. Types of Variables:

- (i) Nominal: Gender, Colours
- (ii) Ordinal: Rich-Poor, Education Levels, Customer Rating
- (iii) Numerical: Income, Age, Price
- (iv) Categorical: Types of cars, Product Categories
- (v) Interval: Temp, IQ Series
- (vi) Ratio: Height, weight

### 4. Population and Sample:

- (i) Population: It is the entire group of individuals

Ex: All people in India, All users of Netflix

- (ii) Sample: It is a subset of population.

### 5. Why Sample is needed?

- Ans: (i) To reduce the cost of data collection.  
(ii) When a full census of data can not be taken.

### 6. Sampling Techniques:

- (i) Random Sampling: Every member of the population has an equal chance of being selected.

- (ii) Stratified Sampling: The population is divided into subgroups (strata) and random samples are taken from each.

iii) Systematic Sampling: Every ~~n<sup>th</sup>~~<sup>nth</sup> member of the population is selected after a random starting point.

iv) Cluster Sampling: The population is divided into clusters, some clusters are randomly selected, and all members of chosen clusters are sampled.

## 7. Factors to choose the Sampling Techniques

- fns:
- i) Nature of population
  - ii) Research Objective.
  - iii) Available Resource.

## 8. Steps for Statistical Data Analysis

- fns:
- i) Define the problem or research question.
  - ii) Data collection.
  - iii) Data Cleaning.
  - iv) Exploratory Data Analysis [EDA]
  - v) Data Transformation.
  - vi) Hypothesis Formulation.
  - vii) Statistical Testing
  - viii) Interpretation of Results.
  - ix) Draw Conclusions.
  - x) Document the Analysis Process / Report Making.

## 9. Measures of Central Tendency:

- i) Mean: The average of data points. Use when data is evenly distributed without extreme outliers.  $\bar{x}$

⑪ Median: The middle value, when the data points are ordered. Use, when data has outliers or is skewed, as it is not affected by extreme values.

In case of even numbers, median is the average of the two middle values.

$$x = \{3, 1, 7, 5, 9\} \\ = \{1, 3, 5, 7, 9\} \rightarrow \text{odd}$$

median

$$x = \{1, 3, 5, 9, 11, 13\}$$

↓

$$\frac{5+9}{2} = \frac{14}{2} = 7$$

[Median]

⑫ Mode: The most frequently occurring value(s) in dataset.

→ Use, when you need to identify the most common value, especially in categorical data.

$$x = \{2, 4, 3, 4, 2, 6, 7, 2\}$$

Mode = 2

⑬ Summary:

① Distribution of data (Depend)

② Numerical Variable without outliers (mean)

③ Numerical Variable with outliers. (Median)

④ Categorical Variable. (Mode)

## b. Measures of Dispersion:

It describes the spread/variability of data points within a dataset.

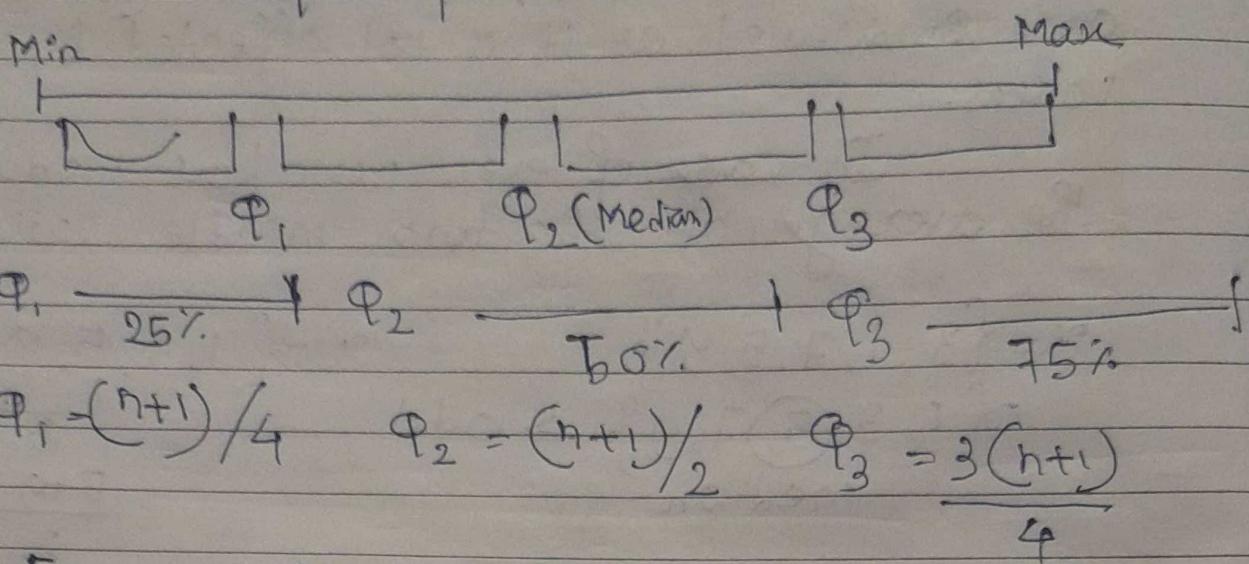
① Range: Max - Min -

② Sensitive to extreme values.

③ It does not consider the distribution.

④ Suitable for small dataset

② Quartile: Value that divide a dataset in 4 equal parts.



Example:  $\{4, 6, 7, 8, 10, 23, 34\}$ ,  $n=7$

$$Q_1 = \frac{(7+1)}{4} = \frac{8}{4} = 2 \quad Q_1 = 6$$

$$Q_2 = \frac{8}{2} = 4 \quad Q_2 = 8$$

$$Q_3 = \frac{3(n+1)}{4} = \frac{3(7+1)}{4} = \frac{24}{4} = 6$$

$$Q_3 = 23$$

③ Percentile:

$$\text{Formula: } P = \frac{P}{100} \times (n+1)$$

A value, below which a given percentage of data points fall.

Use to understand the relative standing of a data point within a dataset.

④ Interquartile Range (IQR): The difference between first ( $Q_1$ ) and third ( $Q_3$ ) quartiles.

Formula:  $Q_3 - Q_1 = 50\%$  of the data.

To see the middle portion of the data.  
Less sensitive to extreme values.

Use to measure variability and to identify outliers in skewed distributions.

⑤ Variance:

Formula:  $(V) = \frac{\sum_{i=1}^N (x_i - \bar{u})^2}{N}$

$N$  = Total number of data points

dataset = {1, 2, 3, 4, 5, 6, 7, 8, 9, 10}

$x_i$  = Individual data points  
 $\bar{u}$  = mean

$$\bar{u}(\text{mean}) = \frac{(5+6)}{2} = 5.5$$

$$(x_i - \bar{u})^2 = (1-5.5)^2$$

$$= (2-5.5)^2$$

$$= (3-5.5)^2$$

$$\vdots$$

$$\Rightarrow (10-5.5)^2$$

Sum / N

The average of the squared differences from the mean.

Use to measure the overall variability within a dataset.

⑥ Standard deviation (std): The square root of Variance, indicating

the average distance from the mean. Use when you need to measure of spread, that is in the same units as the data and to understand data dispersion.

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

ii) Frequency: No of times, a value of data occurs  
Use in Categorical Variables.

iii) Relative Frequency: Percentage or portion of the data value present in a complete dataset.

Formula:  $\frac{\text{Frequency}}{\text{Total no of Observations}} \times 100 \text{ if Percentage}$

iv) Cumulative Frequency: It is the running total of frequencies upto a certain point in a dataset.

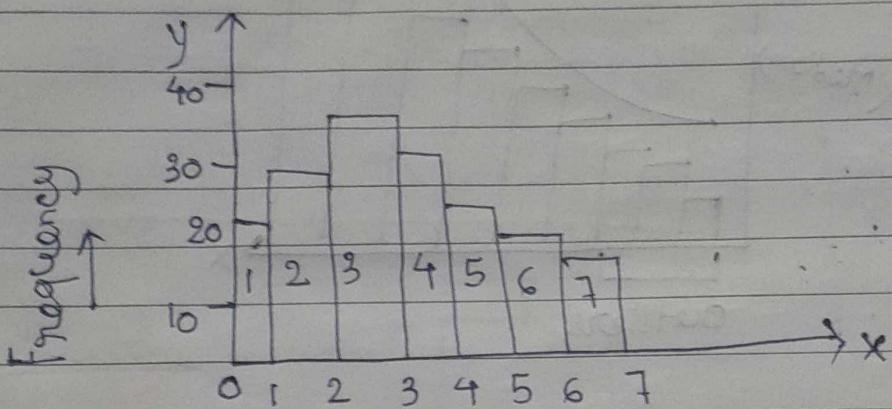
Ex: If the frequency of test scores are 2, 3, 5, then the cumulative frequency are  $\Rightarrow 2$ ,  $5 = [2+3]$ ,  $10 = [2+3+5]$

\* Dataset = {1, 2, 1, 2, 3, 4, 2, 3, 4, 1, 1, 2, 3, 4, 4, 2, 1, 3, 4, 1, 2, 2, 3, 4}

values	Frequency	Cumulative Frequency
1	6	6/24
2	7	7/24
3	5	5/24
4	6	6/24

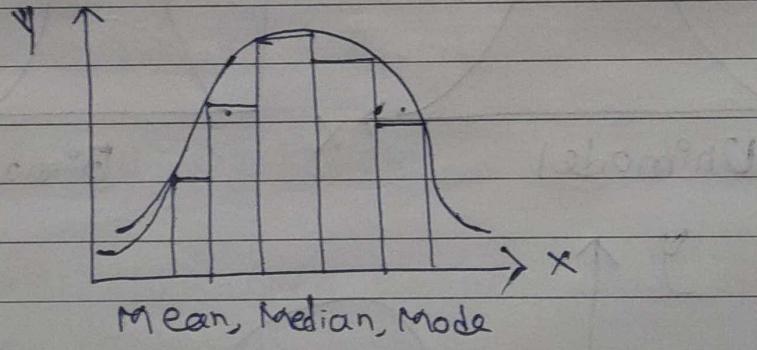
## 12. Graphical Representation:

① Histogram: A bar graph, telling the frequency distribution of a dataset.

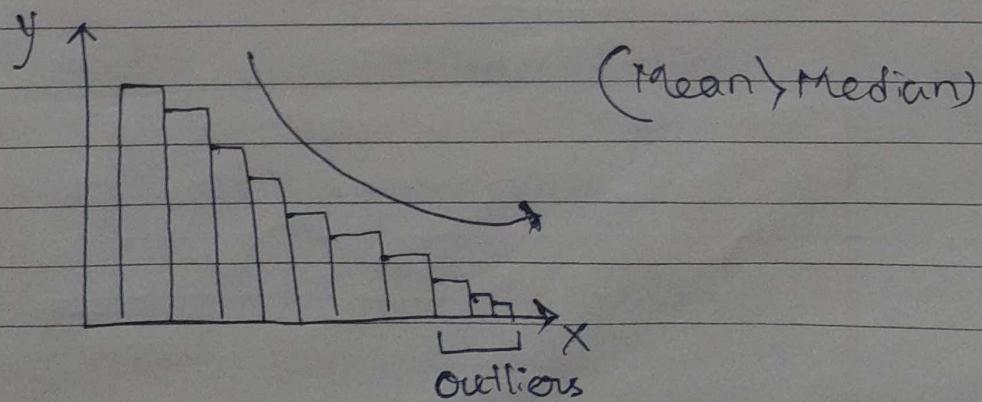


② Types of Skewed Histograms:

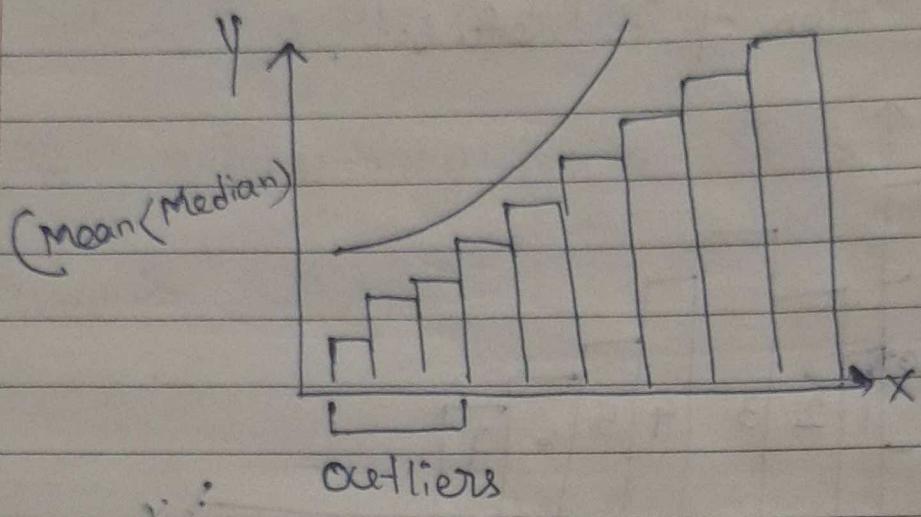
③ Symmetric (Normal Distribution): A histogram, where left and right sides are approximately mirror images.



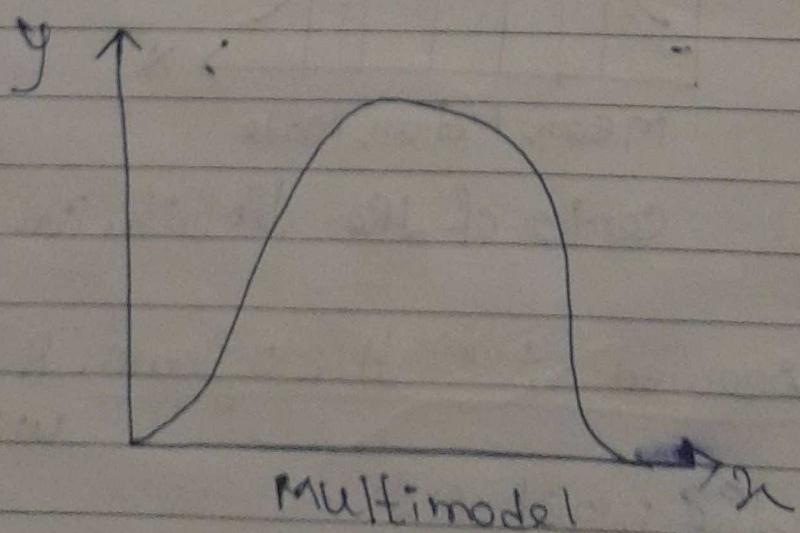
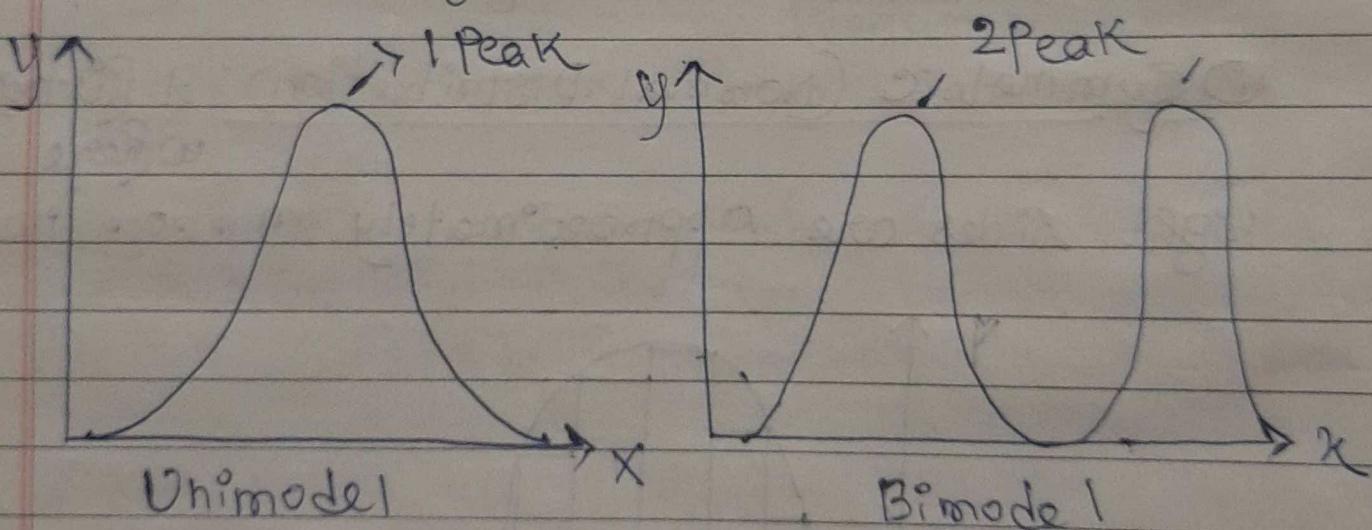
④ Right Skewed (+ve) Histogram: A histogram with the longer tail on the right side.



③ Left skewed (-ve) Histogram: A histogram with the longer tail on the left side.

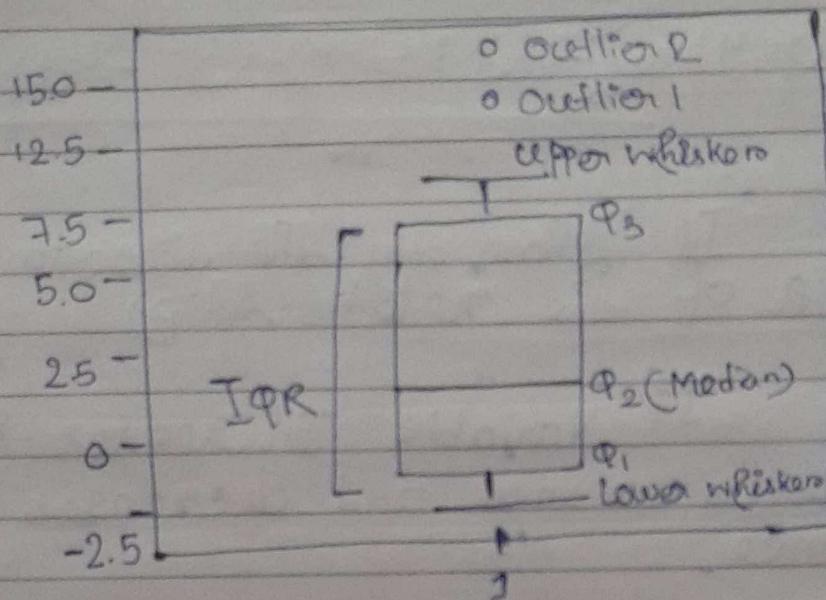


④ Types of Histogram based on number of modes:



(ii) Box plots: A graphical representation showing the distribution of the data and identify outliers. Also known as whisker plot.

$$\text{Box - IQR} = \Phi_3 - \Phi_1$$

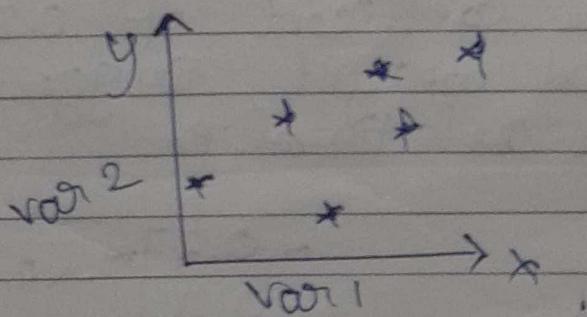


$$\text{Upper whisker} = \Phi_3 + (1.5 \times \text{IQR})$$

$$\text{Lower whisker} = \Phi_1 - (1.5 \times \text{IQR})$$

## 5 Summary Statistics

(iii) Scatter plots: Relationship between two continuous Variable.



Helps identifying patterns, trends, correlation, useful for identifying outliers & understanding strength and direction of relation.

13. (i) Covariance: a measure of how two variables change together.

Ex: Positive covariance between study time and test scores indicates they tend to increase together.

$$\text{cov}(x, y) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})$$

It shows only the direction (not the strength of the relationship)

If variable increases it is positive covariance and decreases negative covariance.

(ii) Correlation: A standardized measure of the strength and direction of the relationship between two variables. It shows linear relation.  $(-1 < r < 1)$

-ve    0    +ve  
No Correlation

Ex: Correlation coefficient of 0.8 between height and weight suggests a strong positive relationship.

$-0.5 < r < 0.5$

No Correlation

or weak Correlation

Pearson's Correlation Coefficient

$$r(x, y) = \frac{\text{cov}(x, y)}{\sqrt{\text{var}(x)} \sqrt{\text{var}(y)}}$$

$r = 0.7$  to  $0.9$   $\rightarrow$  +ve correlation

Strong

Correlation

$$\text{Cov}(x,y) = \sum_{i=1}^N \frac{(x_i - \bar{x})(y_i - \bar{y})}{N}$$

Correlation: A relationship where one variable directly affects another. It is direct connection in which one variable influences the other.

Limitations:

i) Outliers

ii) It does not measure non linear relationship.

14. Removing Outliers:

i) Z score

ii) IQR =  $Q_3 - Q_1$

$Q_3 + 1.5 \times \text{IQR} < \text{outlier}$

$Q_1 - 1.5 \times \text{IQR} > \text{outlier}$

$$\text{Z Score} = \frac{x - \mu}{\sigma}$$