

## Project - 1

### Part A : Theory & Definitions

(Q.1) Define the following with examples from the dataset

(A) Types of data : Numerical & Categorical

Ans. According to the dataset Numerical data's are as follows.

(1). Numerical Data : Measurable quantities which are countable in numbers.

There are 4 types of Numerical Data :

- Discrete  $\rightarrow$  Countable Numbers
- Continuous  $\rightarrow$  Measurable quantities with decimals.

$\rightarrow$  According to my dataset Numerical data's are as follows:

(1). Age of household Head  $\rightarrow$  Discrete

(2). Household Income  $\rightarrow$  Discrete

(3). Family Size  $\rightarrow$  Discrete

(4). Height  $\rightarrow$  Continuous

(2). Categorical Data  $\rightarrow$  Descriptive labels or names

There are 2 types of Categorical Data.

- Nominal  $\rightarrow$  No specific order
- Ordinal  $\rightarrow$  Ordered Categories.

As per my dataset some example of Categorical data are: Household, Education level, Owns House, Urban Rural

- (1). Household
- (2). Education level
- (3). Owns House
- (4). Urban Rural

(B). Types of Statistics: Descriptive vs. Inferential

(1). Descriptive Statistics → It describes about the data, it summarises the data for eg: mean, median, mode → Measure of Central Tendency

Measure of Dispersion.

(2). Inferential Statistics → It provides us the conclusions from the data and we can use it for making predictions.

Eg. of Inferential are: - Linear Regression, logistic regression etc.

(C). What is Descriptive Statistics?

→ The process of presenting the data in a summarising and presenting raw data in a meaningful ways is known as descriptive statistics.

(A) Explain the difference between it.

(A) Mean, Median, Mode

Mean = The mean is the sum of all values divided by the number of values.

Median = The median formula =  $\frac{\text{Sum of all values}}{\text{Numbers of values}}$

Median → The middle value of a sorted list

Mode → The mode is the value that appears most frequently.

(B) Range, Variance, Standard Deviation

Range → The difference between the maximum and minimum values is known as range. (Range = Maximum - Minimum).

Variance → It shows how far each value in the dataset is from the mean.

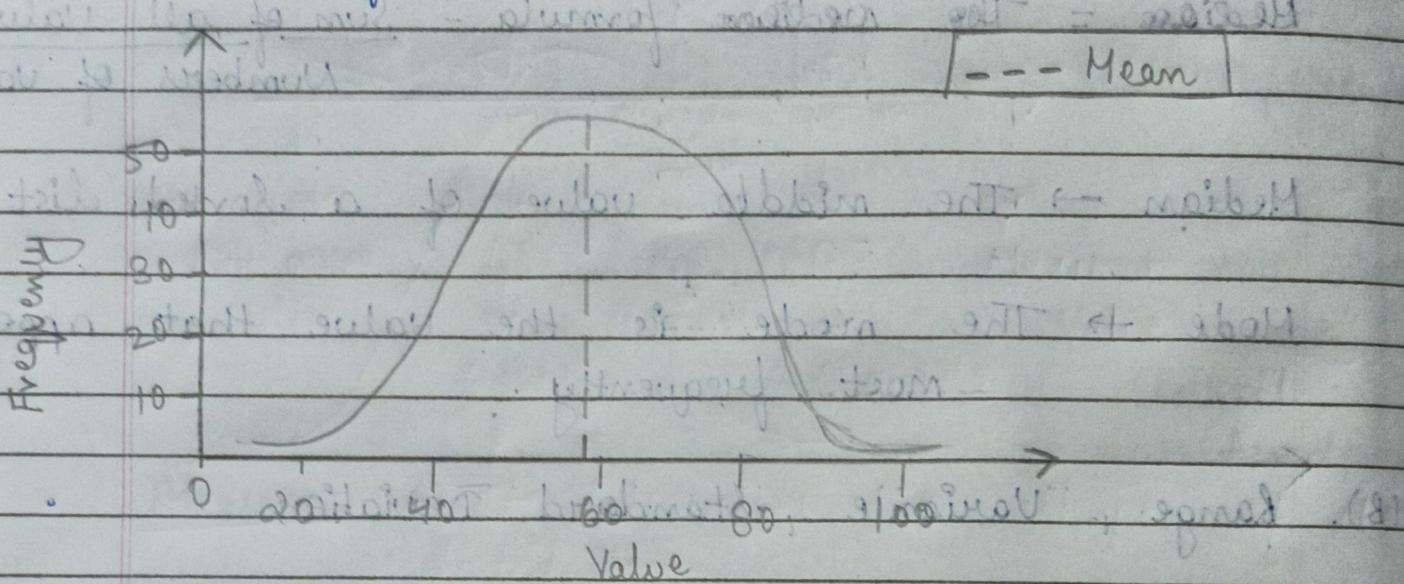
formula → Variance =  $\frac{\sum (x_i - \bar{x})^2}{n}$

Standard Deviation → It shows how much on average the values deviate from the mean.

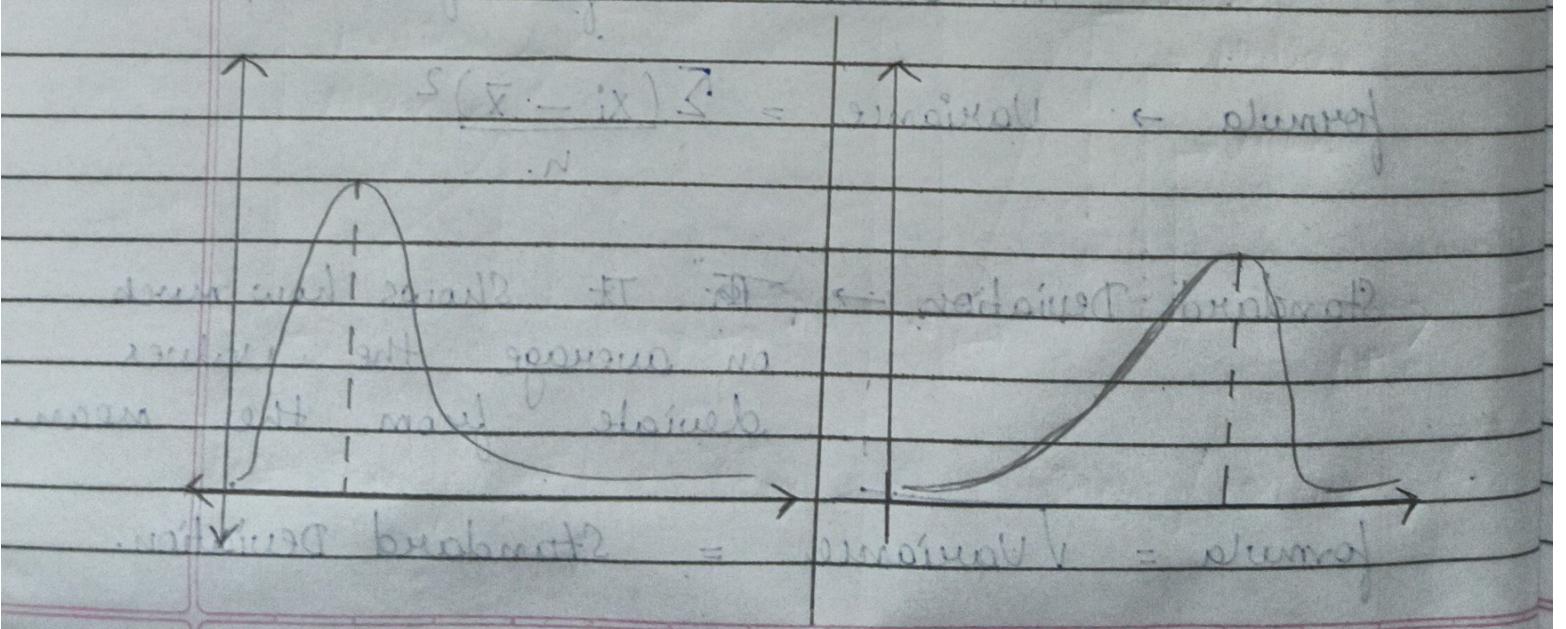
formula =  $\sqrt{\text{Variance}} = \text{Standard Deviation}$ .

(Q.3). Explain the following term with neat and clean diagram along with its formula:

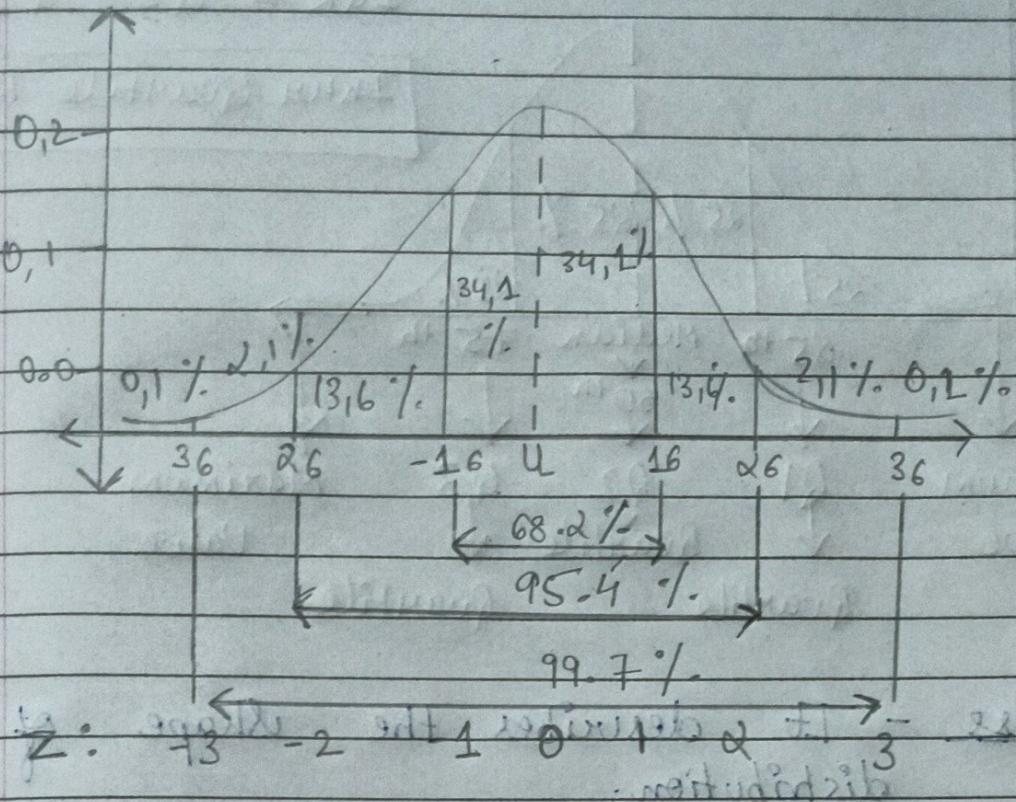
- Gaussian Distribution - Gaussian Distribution is a bell shaped curve that shows how data is distributed for eg. Height of people Marks of students, Blood pressure.



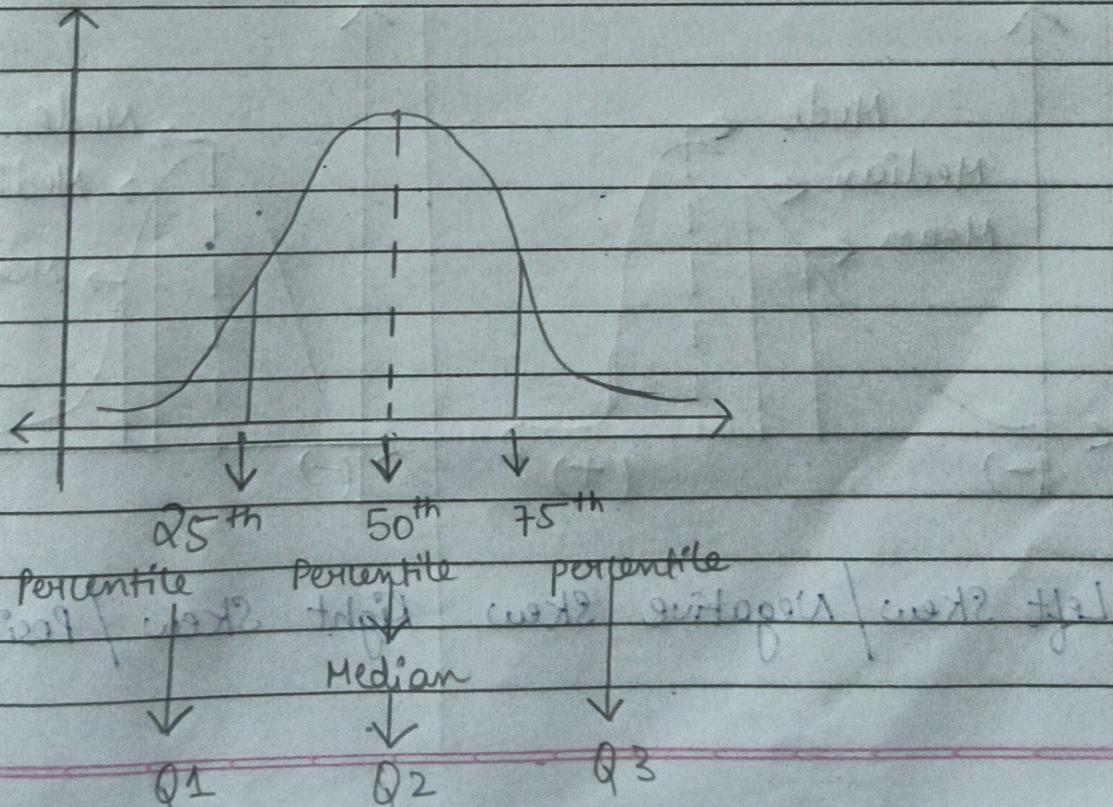
- Log Normal Distribution - An unevenly distributed data is known as log normal data. It may be right or left skewed.



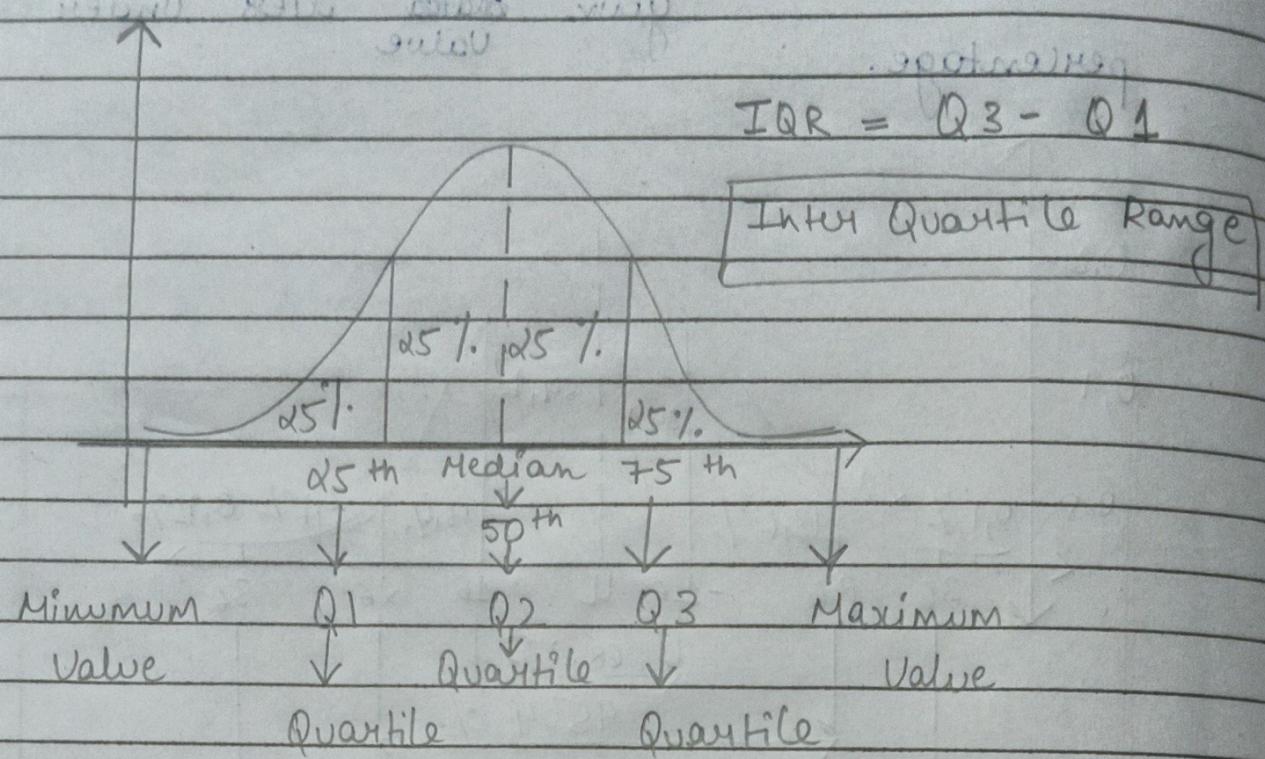
- 3 Sigma Rule or Empirical Rule: It shows where your data lies under the given percentage.



- Percentiles: A percentile indicates the value below which a given observation lies/fall.

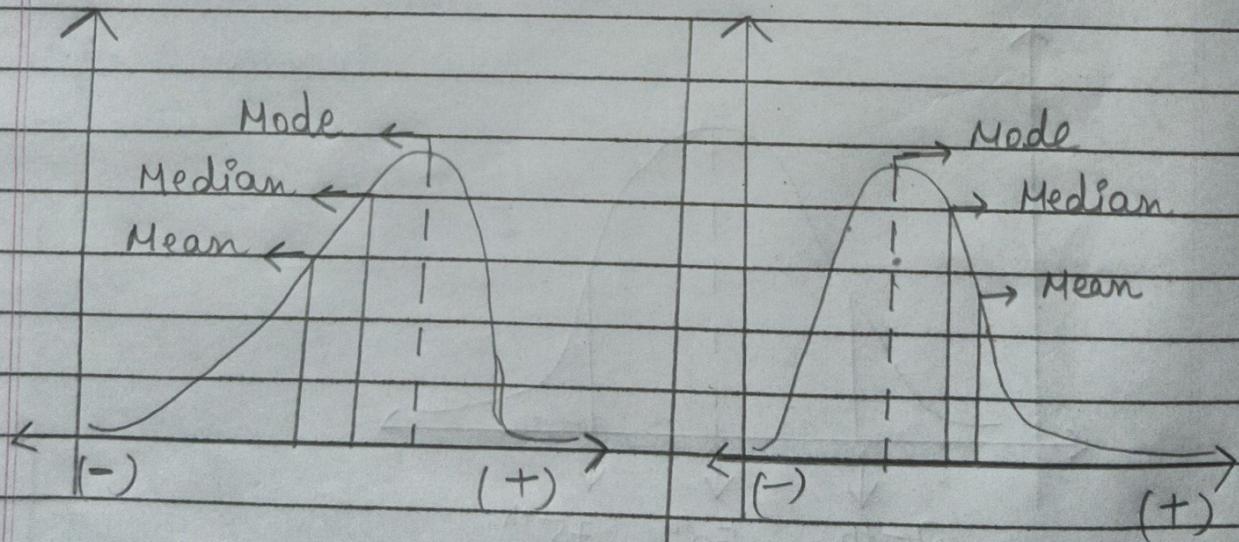


• Quantiles & Five number summary



- Skewness - It describes the shape of the distribution.

~~There are 2 types of skewness.~~



Left Skew / Negative Skew

Right Skew / Positive Skew

## Kurtosis

