

# 1. Measure of Central Tendencies :-

Here we discuss about mainly 3 main topics :-

① Mean

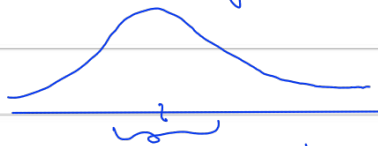
② Median

③ Mode

} # EDA (Exploratory Data Analysis) and Feature Engineering will be using these concepts

Q1 What does measuring of Central Tendency mean?

Ans → Suppose below is the given data distribution :-



So, measure of central tendency talks about this central region where maximum amount of data is present

## 1.1 Mean :-

Suppose we have a population of 10 ppl and sample is taken from 5 ppl the calculating both population and sample mean :-

Population ( $N=10$ )

Sample ( $n=5$ )

$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$

$\therefore$

$$\text{Population mean } (\mu) = \frac{\sum_{i=1}^N x_i}{N}$$

$\therefore$

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

$$\Rightarrow \mu = \frac{32}{10} = \boxed{3.2} = \mu_N$$

## 1.2 Median :-

It is the central Element of the sorted list of data

eg: data = 4, 5, 2, 3, 2, 1

↓

sorting data :- 1, 2, 2, 3, 4, 5

↓

Median: Count the no. of Elements:-

if

a) Even Count:- then median is the mean of the two middle nos.

ex: 1, 2, 2, 3, 4, 5.

$$\text{median} = \frac{2+3}{2} = 2.5$$

b) Odd Count:- median is  $(n+1)^{\text{th}}$  element

ex: - for 1, 2, 2, 3, 4

$$\text{median} = 2$$

☆5

Q Why do we need to calculate median?

Soln → Suppose we have a sample data {1, 2, 3, 4, 5}

$$\text{mean} = \frac{1+2+3+4+5}{5} = 3, \quad \text{median} = 3$$

Now,

Suppose an outlier (100) is added to the above data set such that the dataset becomes

{1, 2, 3, 4, 5, 100}

# an outlier is a value that does not belong/fits into the dataset

∴

$$\text{median} = \frac{1+2+3+4+5+100}{6} = 19.16, \quad \text{mode} = \frac{3+4}{2} = 3.5$$

Observation:- Even after adding a large outlier the value of median have not shifted that much but there is a significant shift in the value of mean

So, whenever we have an outlier it is a better strategy to use median coz change in median value is not significant as compared to the mean.

### 1.34 Mode :-

Element with the maximum frequency will have the mode value

Ex:  $\rightarrow$  for  $\{2, 1, 4, 5, 1, 7, 8, 1, 9, 1, 10\}$   
mode = 1

Q) Where mode is used?

$\Rightarrow$  Suppose dataset :- Type of Flowers

Lilly

Age

10

Rose

3

Rose

5

Sunflower

Median on mean

Rose

8

Replaced by  
'mode'

If outlier is present  
the replace with  
Median or Mean

**EDA on Feature Engineering**