

# Introduction to statistics

Def: Statistics is the science of collecting, organizing and analyzing the data.

Eg: Height of student in the class

## # Types of Statistics

### 1. Descriptive stats

Def: It consists of organizing and summarizing the data.

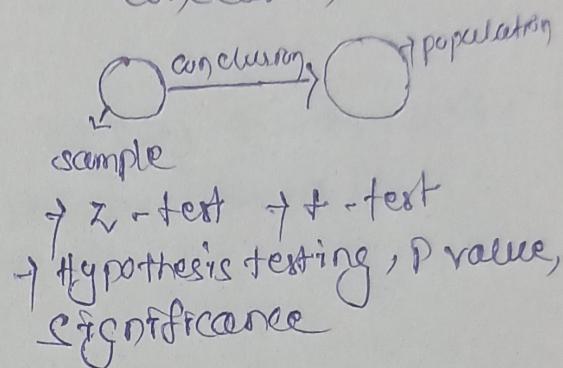
→ Measure of central Tendency  
[mean, median, mode]

→ Measure of dispersion  
[variance, std]

→ Histograms, Bar chart, pie chart

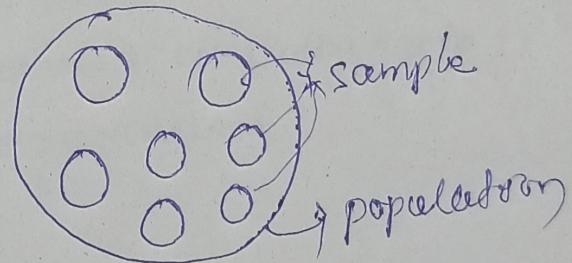
### 2. Inferential stats

Def: It consists of using data you have measured to form conclusion.

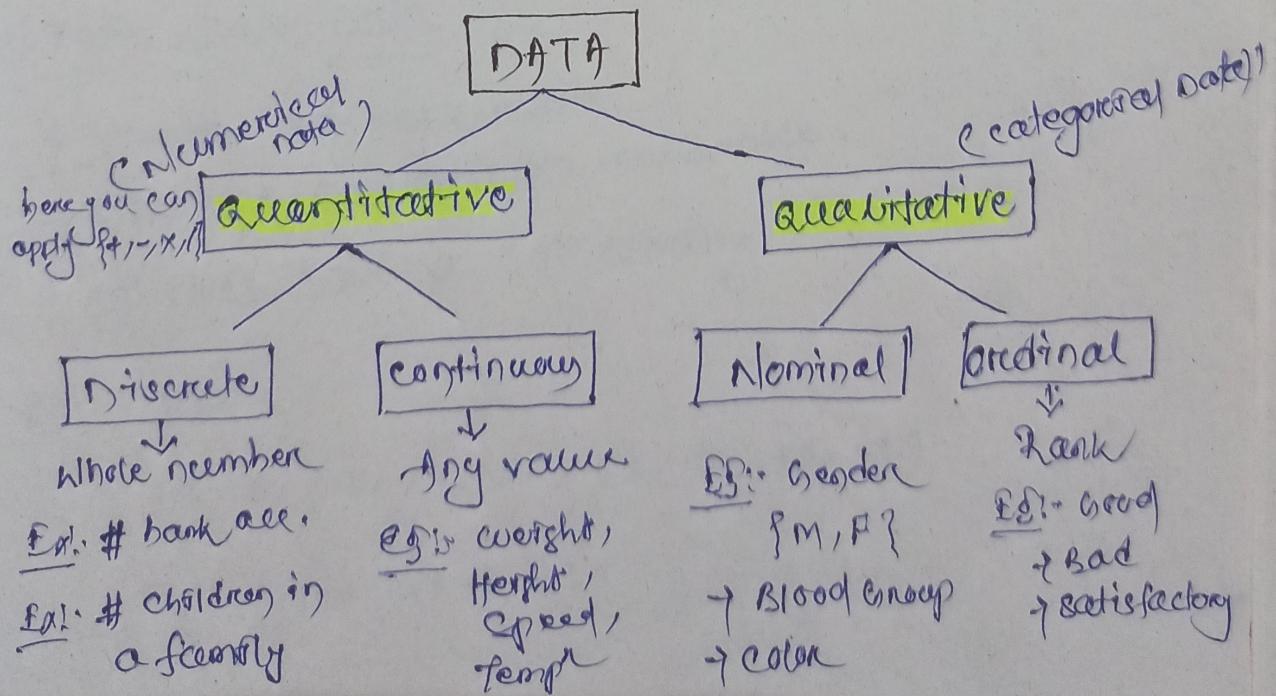


## # Sample Data and population Data

Eg: Exit poll



## # Types of Data



## # Scale of measurement of data

### 1) Nominal scale data:

→ Qualitative / categorical data.  
Eg.: Gender, colors, blood group

→ Order does not matter

### 2) Ordinal scale data:

→ Ranking and order matter.  
→ Difference can't be measured.

Eg.: Qualification

P.H.d	1 <sup>st</sup>
M.S	2 <sup>nd</sup>
B.Tech	3 <sup>rd</sup>
B.com	4 <sup>th</sup>

∴ here Ranking & order possible  
but difference can't be  
measured.

### 3) Interval scale data

→ The rank and order matters.

→ Difference can be measured (excluding ratio).

→ Doesn't have "0" starting value.

Eg.: Temp

### 4) Ratio scale data

→ Order and rank matter.

→ Difference and ratio are measurable.

→ It does have a "0" starting

Eg.: mark

[ 0, 10, 20, 30, 40, 50, 60, 70, 80 ]

$$100 : 50 \Rightarrow 2 : 1$$

# # Measures of central Tendency

## 1. Mean ( $\mu$ ) [Average]

Population ( $\mu$ )

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

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

$$\mu = \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= \frac{32}{10} = 3.2$$

Sample ( $s$ )

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

## 2. Median

Ex: 4, 5, 2, 3, 2, 1

here # element = even

$$\begin{array}{c} \text{sort} \rightarrow 1, 2, \boxed{2, 3}, 4, 5 \\ \downarrow \\ \boxed{(2+3)/2 = 2.5} \\ \boxed{\text{median: } 2.5} \end{array}$$

Ex: 1, 2, 2, 3, 4, 5, 7

# element = odd

median = 3

## 3. Mode [maximum frequency]

$$\{2, 1, 1, 1, 4, 5, 7, 8, 9, 10\}$$

mode = 1

## # measures of central Tendency using Python

# mean

age = [12, 13, 14, 21, 24]

import numpy as np

np.mean(age)

Output

16.8

# median

np.median(age)

Output  
14

# mode

height = [4.2, 5.0, 5.0, 5.2, 6.0]

np.mode(height)

Output

\* module (numpy) has no attribute 'mode'

from scipy import stats  
stats.mode(height)

Output  
5.0

## # Measure of Dispersion

1. Variance : spread of the data.

Population variance

Sample variance

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (\alpha_i - \mu)^2$$

where  $\alpha_i$  = data points

$\mu$  = population mean

$N$  = population size

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (\alpha_i - \bar{\alpha})^2$$

[Why does the sample variance have  $n-1$  in the denominator?

Ans: to create unbiased estimator of the population variance.]  $\rightarrow$  Bias & consistency

where  $\alpha_i$  = data point

$\bar{\alpha}$  = sample mean

$n$  = sample size

Ex:  $\{1, 2, 3, 4, 5\} \rightarrow$  sample

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (\alpha_i - \bar{\alpha})^2$$

$$\bar{\alpha} = \text{sample mean} = \frac{1+2+3+4+5}{5} = 3$$

$\alpha_i$	$\bar{\alpha}$	$(\alpha_i - \bar{\alpha})^2$
1	3	4
2	3	1
3	3	0
4	3	1
5	3	4

$$\text{Now } s^2 = \frac{4+1+0+1+4}{5-1}$$

$$= \frac{10}{4} = 2.5$$

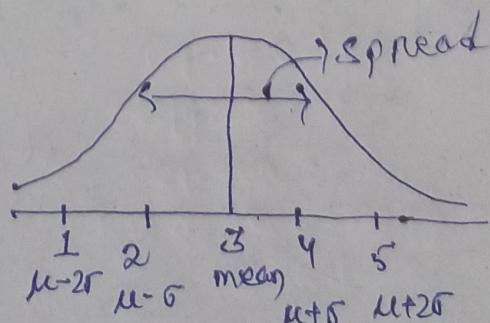
$\sqrt{s^2} = s = \text{sample std deviation}$

Consider:

$$\{1, 2, 3, 4, 5\}$$

$$\mu = 3$$

$$\sigma = 1$$



# measures of dispersion using python

Variance & Standard deviation

ages = [23, 43, 23, 56, 74, 32, 68, 98, 45, 32]

import numpy as np

mean = np.mean(ages)

mean :-  $\bar{x} = 49.4$

var = np.var(ages)

var :-  $s^2 = 541.64$

std = np.std(ages)

std :-  $s = 23.278169507331188$

OR

import pandas as pd

data = [[10, 12, 13], [34, 23, 45], [32, 34, 21]]

data :-  $\bar{x} = \frac{1}{3} [10 + 12 + 13 + 34 + 23 + 45 + 32 + 34 + 21] = 29.67$

df = pd.DataFrame(data, columns=['A', 'B', 'C'])

df :-  $\bar{x} = \begin{array}{ccc} A & B & C \\ 10 & 12 & 13 \\ 34 & 23 & 45 \\ 32 & 34 & 21 \end{array}$

# column wise

df.var() :-  $\bar{x} = \begin{array}{ccc} A & B & C \\ 10 & 12 & 13 \\ 34 & 23 & 45 \\ 32 & 34 & 21 \end{array}$

# Row wise

df.var(axis=1) :-  $\bar{x} = \begin{array}{ccc} 0 & 2.3333 \\ 1 & 121.0000 \\ 2 & 49.0000 \end{array}$

# Random Variable: Random variable is a process of mapping the output of a random process or experiment to a specific number.

Eg:- Tossing a coin  
Rolling a dice

# Covariance And Correlation: covariance indicates the relationship of two variables whenever one variable changes, if an increase in one variable results in an increase in the other variable, both variance variables are said to have a positive covariance. Decrease in one variable also cause a decrease in the other, both variable said to have a +ve covariance. But in case of increase in one variable other is decrease then both said to have negative covariance.

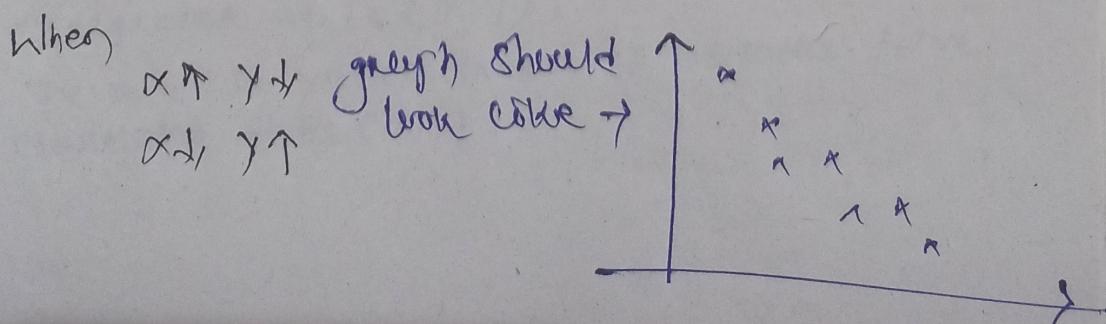
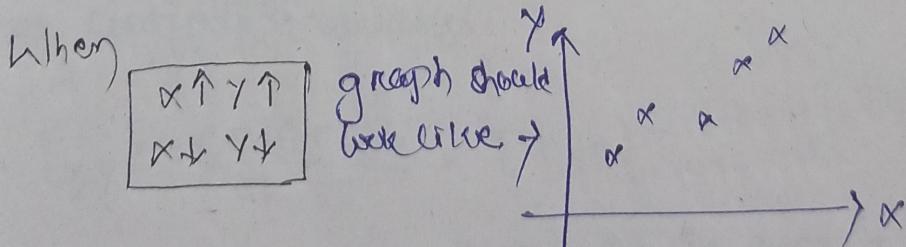
Eg:-

X	Y
2	3
4	5
6	7
8	9

 ∵ here X & Y are two variable  
 { Possible Relationship betw' X and Y can made? }

		$\left\{ \begin{array}{l} X \uparrow \\ Y \uparrow \end{array} \right.$	$\left\{ \begin{array}{l} X \downarrow \\ Y \downarrow \end{array} \right.$	When X & Y both are increasing or decreasing then both variable said to have +ve covariance.
		$\left\{ \begin{array}{l} X \uparrow \\ Y \downarrow \end{array} \right.$	$\left\{ \begin{array}{l} X \downarrow \\ Y \uparrow \end{array} \right.$	
		$\left\{ \begin{array}{l} X \downarrow \\ Y \uparrow \end{array} \right.$	$\left\{ \begin{array}{l} X \uparrow \\ Y \downarrow \end{array} \right.$	

When among  
X & Y variable  
one of them increase  
other is decrease  
then both variable  
said to have -ve  
covariance.



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

where  $x_i$  = data point

$\bar{x}$  = sample mean

$n$  = sample size

Also we know that,

$$\text{var}(x) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \Rightarrow \text{cov}(x,y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

[spread of data]

Ex:-

	$x$	$y$
2	3	
4	5	
6	7	
$\bar{x} = 4$	$\bar{y} = 5$	

$$\begin{aligned} \text{cov}(x,y) &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{2} \\ &= \frac{4+0+4}{2} = 4 \text{ (+ve)} \end{aligned}$$

# Advantage

# Disadvantage

→ relationship between  $x$  &  $y$   
+ve or -ve value. → covariance does not specify  
cimed value.

# Pearson correlation coefficient [-1 to 1]

The Pearson coefficient is a type of correlation coefficient that represents the relationship between two variables that are measured on the same interval or ratio scale. The Pearson coefficient is a measure of the strength of the association between two continuous variables.

$$r_{xy} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} = [-1 \text{ to } 1]$$

→ The more the value toward +1 the more +ve correlated it is ( $x|y$ ).

→ The more the value toward -1 the more -ve correlated it is ( $x|y$ ).

## # Spearman Rank Correlation

$$r_s = \frac{\text{cov}(R(X), R(Y))}{\sqrt{R(X)} \times \sqrt{R(Y)}}$$

Ex:

X	Y	R(X)	R(Y)
1	2	4	5
3	4	3	4
5	6	2	3
7	8	1	1
0	7	5	2

$\overrightarrow{R(X)} = 3 \quad \overrightarrow{R(Y)} = 3$

## # Covariance And Correlation with python

```
import seaborn as sns
```

```
df = sns.load_dataset("xyz")
```

### # Covariance

```
df.cov()
```

### # OR using numpy

```
import numpy as np
```

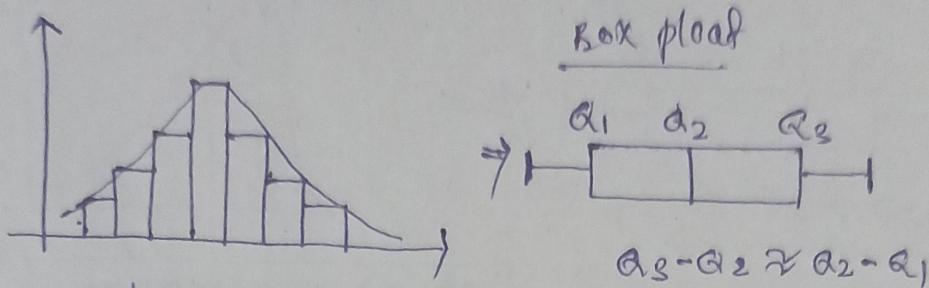
```
np.cov()
```

### # Correlation

```
df.corr(method='spearman')
```

```
df.corr(method='pearson')
```

# Skewness: skewness is a measure of how asymmetrical a distribution is about its mean.  
 → Symmetrical Distribution



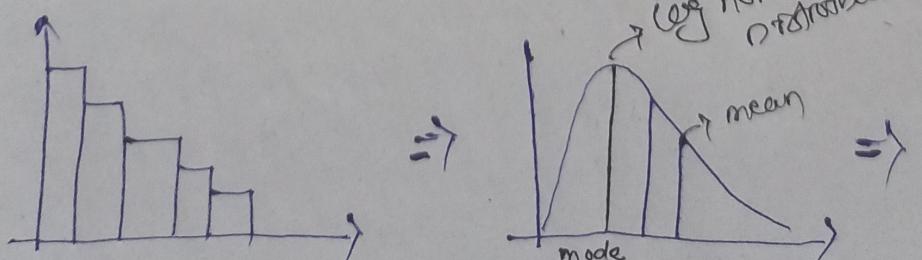
- 
- Normal / Gaussian Distribution
  - Symmetric Distribution
  - No skewness
- ↳ The mean, median, mode all are perfectly at center.  
 [mean = median = mode]

Def' of Skewness: A distribution is said to be skewed when the mean and the median fall at different points in the distribution, and the balance is shifted to one side or the other to left or right side.

measures of skewness tell us the distribution and the extent of skewness. In symmetric distribution the mean, median and mode are identical. The more the mean moves away from the mode, the larger the asymmetry or skewness.

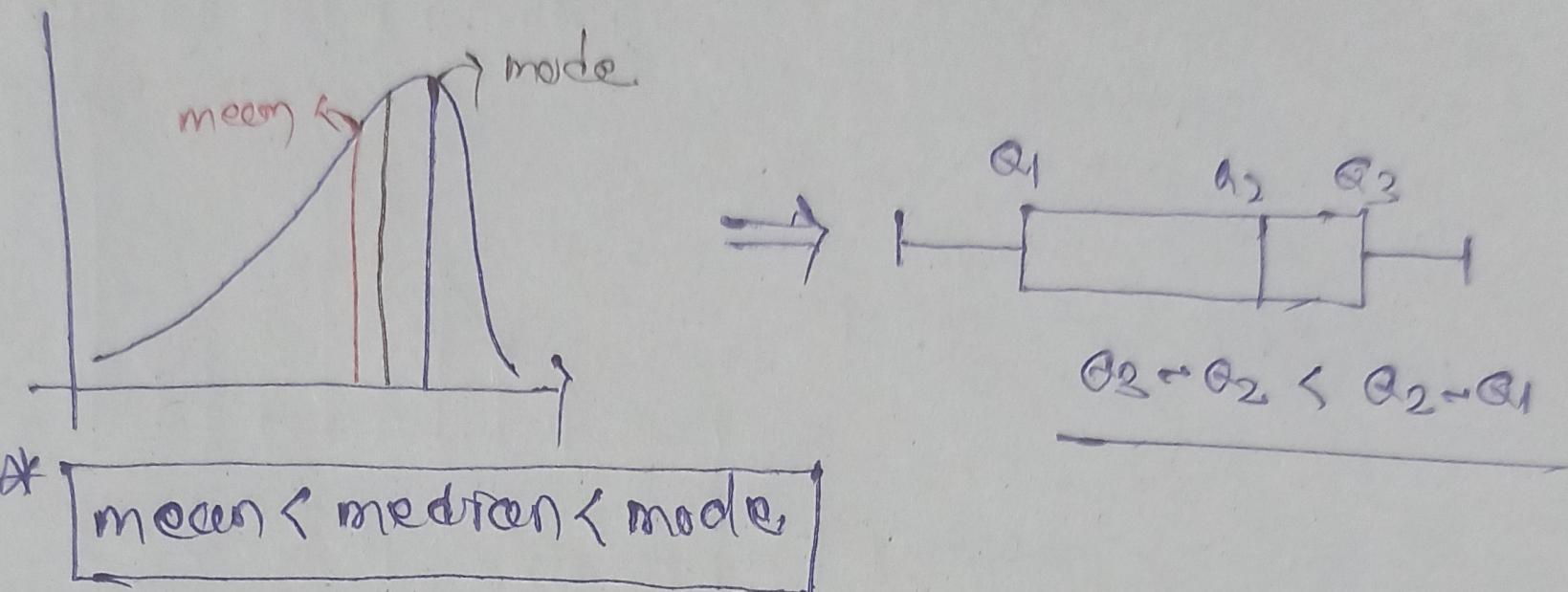
When data points on a bell curve are not distributed symmetrically to the left and right sides of the median, the bell curve is skewed. distributions can be positive and right skewed, or negative and left skewed.

# Right Skewed Distribution [Positive Skewed]



$[Q_3 - Q_2 > Q_2 - Q_1] \Rightarrow \text{mean} > \text{median} > \text{mode}$

## # Left skewed Distributions



## Histogram [Frequency]

$$\text{ages} = \{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 48, 51\}$$

Bins = 10

$$\text{Bin size} = \frac{50 - 0}{10} = 5$$