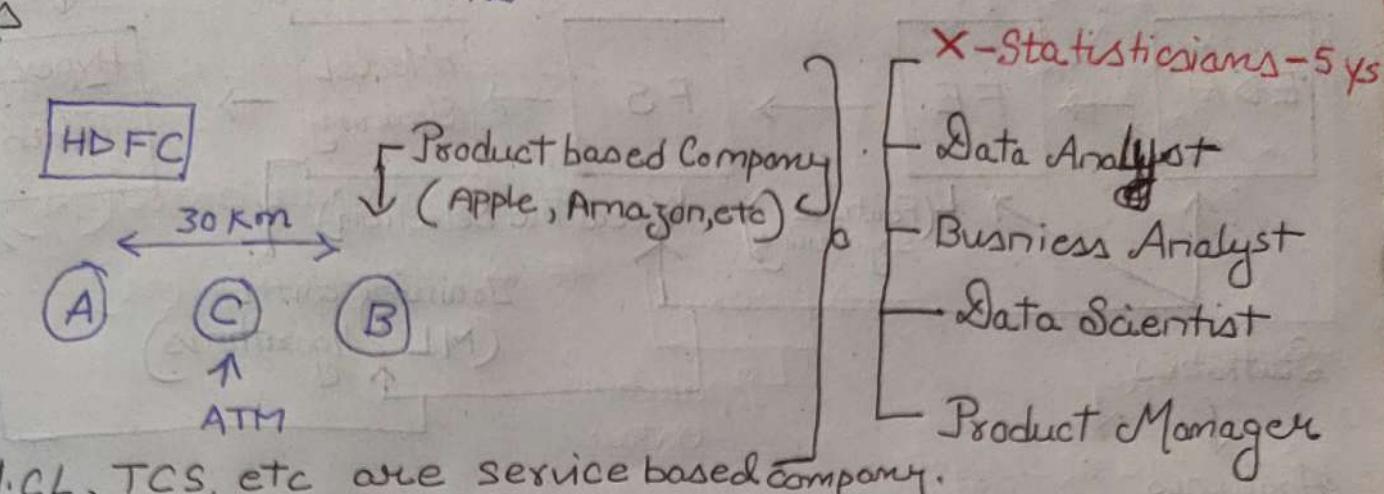


# Statistics

→ Use Cases: For banks  
examples

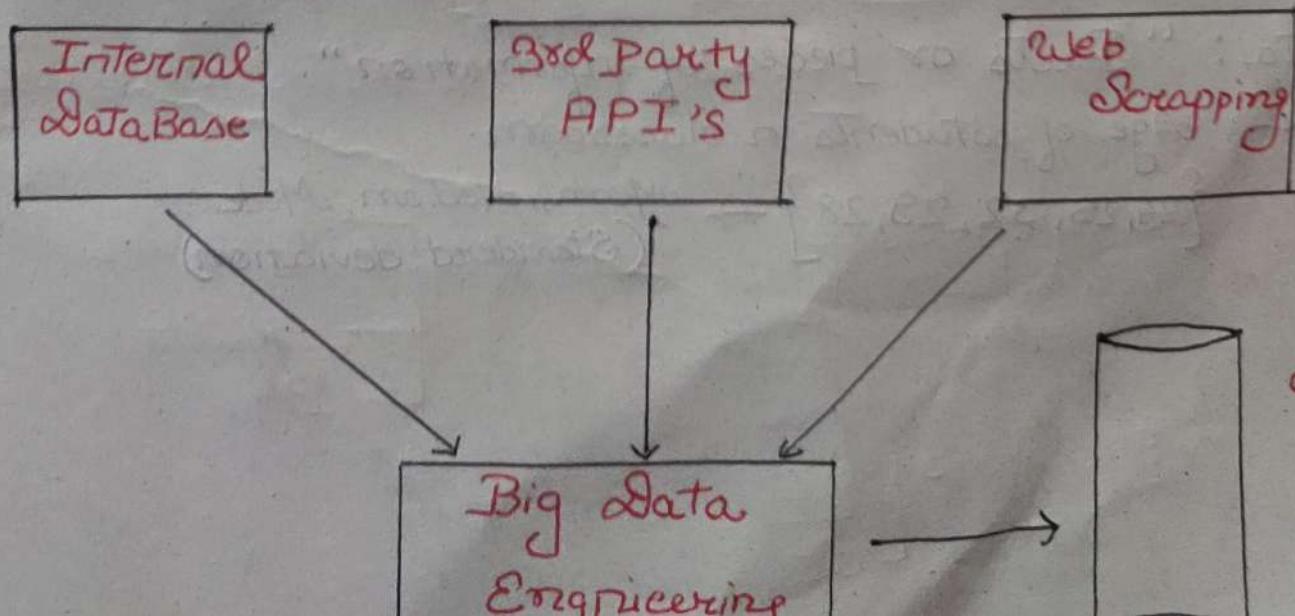
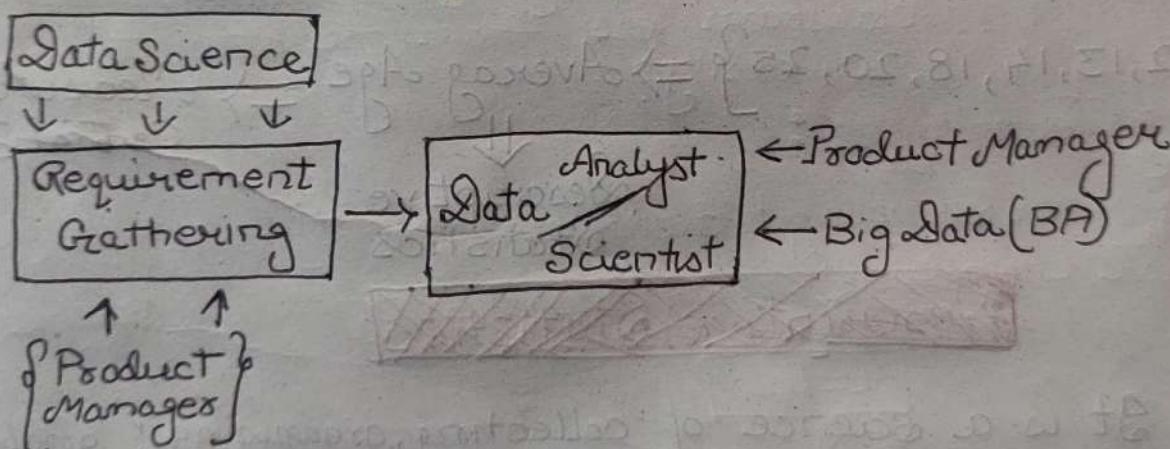
①



→ HCL, TCS, etc are service based company.

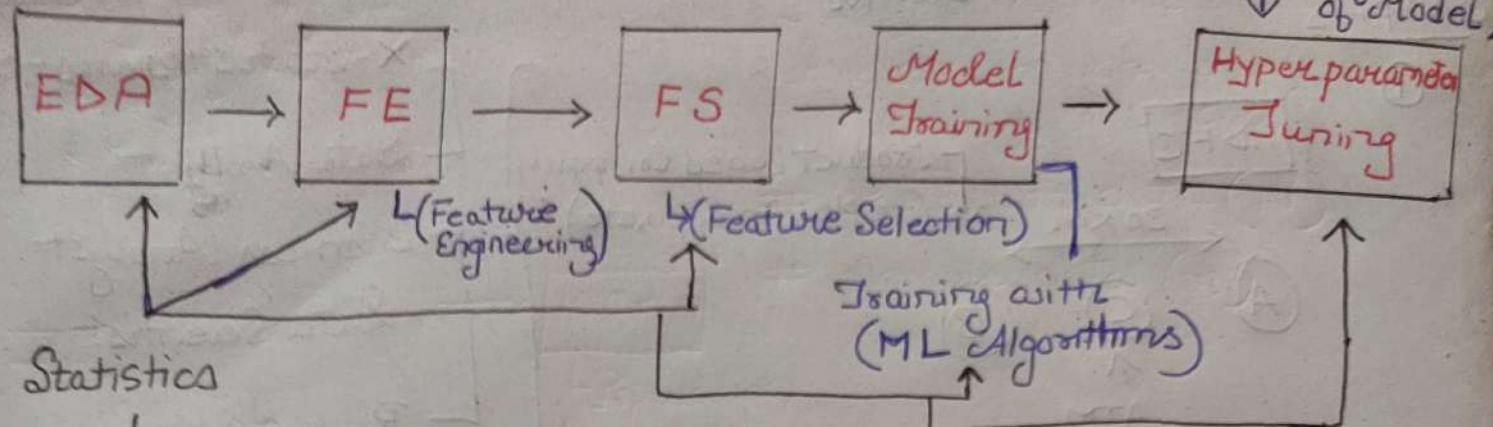
- ② Identifying the avg. osize of the chart in the world.
- ③ Amazon or Flipkart big-billion day sale, identifying the month of cselection.

## Definition in Life cycle of Data Science Project



# Life Cycle of DS Project

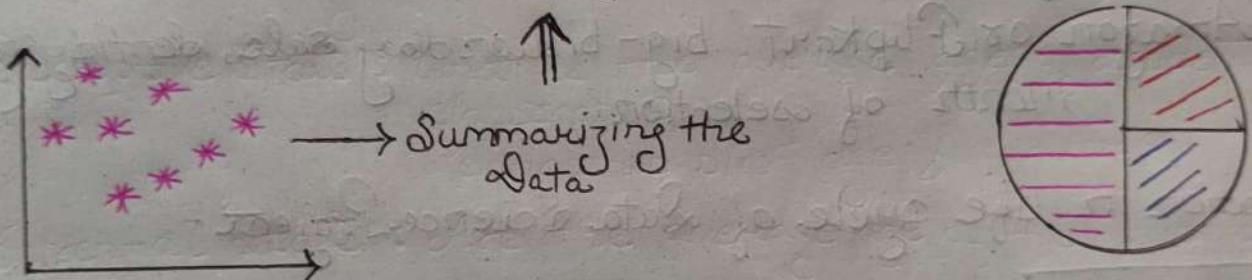
(Exploring Data Analysis)



Statistics

Analysis  
of  
Data

Descriptive Stats



$$\text{Age} = \{12, 13, 14, 18, 20, 25\} \Rightarrow \text{Average Age} \Rightarrow \text{Measure of Central Tendency}$$

↓  
Descriptive Statistics

~~Descriptive Statistics~~

\* **Statistics:** It is a science of collecting, organising & analysing the data.

→ **Data:** "Facts or pieces of information".

ex:- Age of students in classroom.

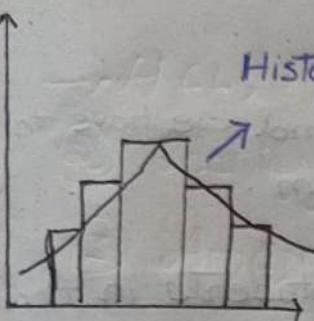
$$\{24, 25, 32, 29, 28\} \rightarrow \text{Mean, Median, Mode} \\ (\text{Standard deviation})$$

# STATISTICS

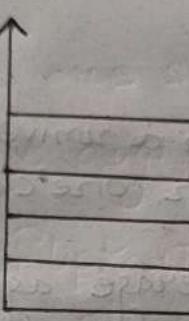
## Descriptive [EDA + FE] Stats

- ① St consists of organising & summarising the data

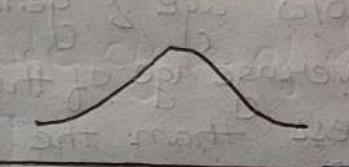
Histogram



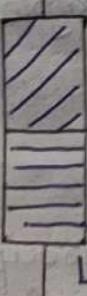
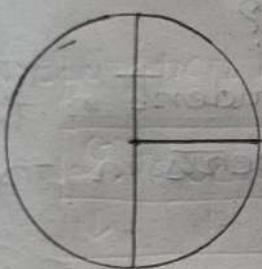
Bar Chart



Distribution



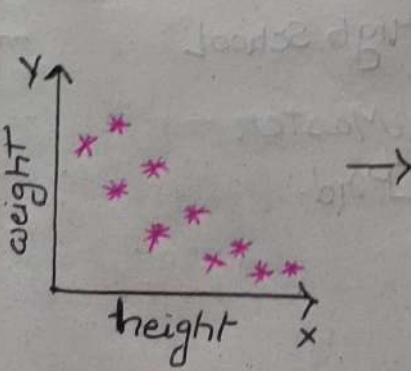
Pie



Candle Stick

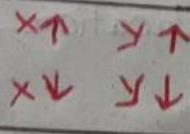


Boxplot



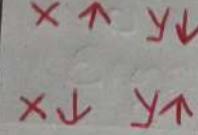
weight

height



Condition ①

[Indirectly Proportional]



Condition ②

## Inferrential Stats

St consists of collecting sample data and making conclusion abt population data using experiments

Hypothesis Testing

University  $\rightarrow$  500 people

Class A  $\rightarrow$  60 people

Sample data  $\rightarrow$  Age  $\rightarrow$  Avg. age of entire university

Hypothesis Testing  
(CI = Confidence Interval)

## Types of Testing :

- ① Z-test
- ② t-test
- ③ Chi-squared test
- ④ f-test

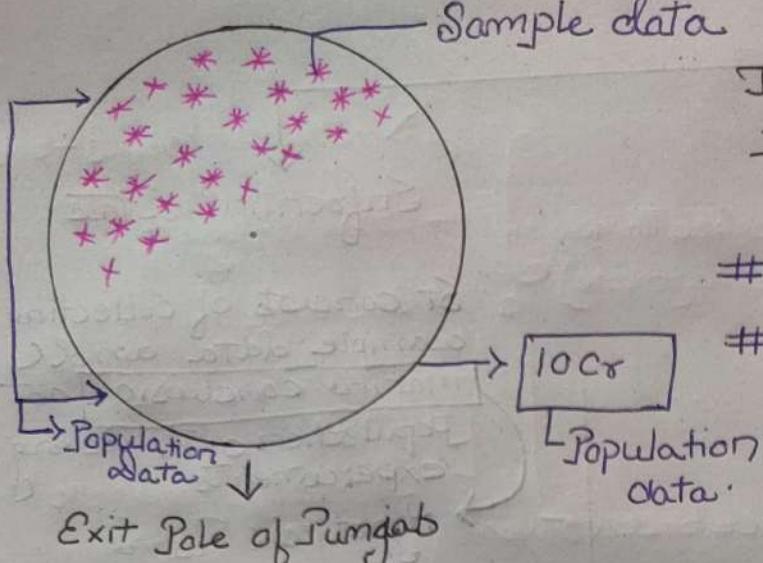
Sample Data

Conclusion

Population Data

Hypothesis Testing

## \* Sample data vs Population data :



Through sample data we can predict, which party will win or loose in this case.

# Population data ( $N$ )  
# Sample data ( $n$ )

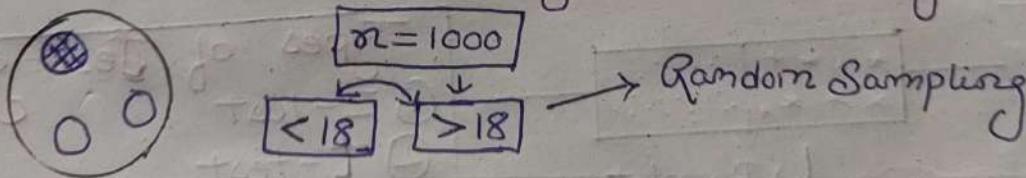
e.g.: Let say there are 20 classrooms in a university and we have to collect the age of students in one classroom.

→ **Descriptive Stats** : What is the average age of students in the classroom?  
Relationship b/w age & gender?

→ **Infrential Stats** : Are the average age of the students in the classroom less than the students in the university?

## \* Sampling Techniques :

(i) **Simple Random Sampling** : Every no. of the population ( $N$ ) has an equal chance of being selected for your sample ( $n$ )



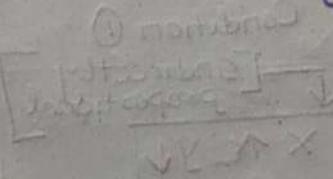
Strata → Layers → Clusters → Groups

## (ii) **Stratified Sampling** :

Gender  
Male  
Female

Blood group  
 $A^+$   
 $B^+$   
 $O^+$

Education  
High School  
Master  
P.h.d.



(iii) Systematic Sampling : → {Airport}  
 {Credit Card}  $n^{\text{th}}$  person

ex: Select every  $n^{\text{th}}$  individual out  $5^{\text{th}}$  person of population ( $N$ )  
 or Select every  $n^{\text{th}}$  individual out of population ( $N$ )

(iv) Convenience Sampling : → Only those who are interested in  
 the survey will only participate.

ex: • (Data Science Survey → General AI Survey)  
 • filling a survey form at an institution.  
 • RBI Survey - Women → Stratified + Random Sampling.

\* Variable : It is a property that can take any values.

ex - age = 14 ; Variables  $\Rightarrow$  Ages = [25, 26, 32, 34]. Collection of ages.

→ Types of Variable :

① Quantitative Variable → measured numerically {Mathematical operation}  
 ex: Age, weight, height, rainfall, temp, distance etc.

② Qualitative Variable → Categorical Variables {based on some characteristics they are grouped together}  
 ex - Gender, type, class, species etc.

### Quantitative Variable

#### Discrete Variable

eg: Whole no. - fixed  
 • no. of bank account  
 (it can't be indecimal)  
 • no. of children  
 • Pincode - fixed

#### Continuous Variable

eg. decimal values  
 height, weight, age, rainfall, speed

### # Assignment :-

- |     |  |
|-----|--|
| I   | What kind of variable is Marital Status? — Categorical |
| II  | " " " " " River length? — Continuous                   |
| III | " " " " " Movie duration? — Continuous                 |
| IV  | " " " " " Pincode? — Discrete                          |
| V   | " " " " " IQ? — Discrete                               |

## Outline

- ① Histograms
- ② Measure of Central Tendency
- ③ Measure of Dispersion
- ④ Percentiles & Quartiles
- ⑤ 5 Number Summary (Box Plot)

→ Histogram :

$$\text{Ages} = \{10, 12, 14, 18, 24, 26, 30, 35, 37, 40, 41, 42, 51, 60, 65, 68, 78, 90, 95, 100\}$$

Step-1 : Sort the numbers

Step-2 : Bins → No. of groups

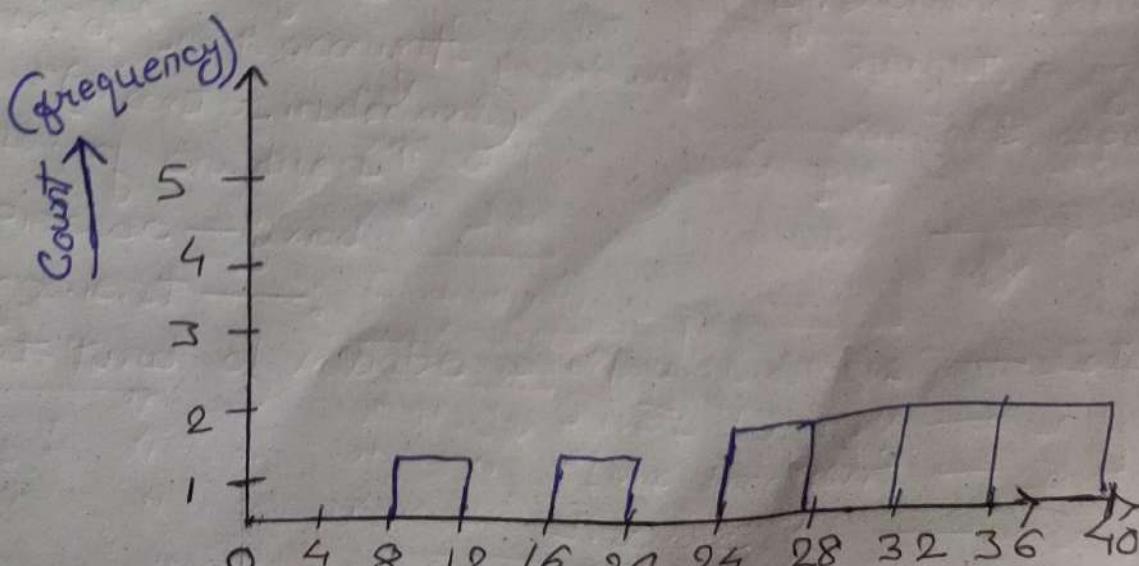
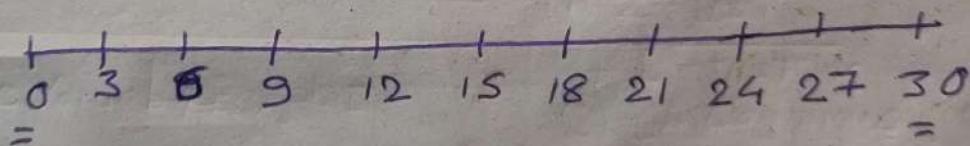
Step-3 : Bin Size - size of bins

$$[10, 20, 25, 30, 35, 40]$$

$$\text{bins} = 10$$

$$\begin{aligned} \min &= 10 \\ \max &= 40 \end{aligned}$$

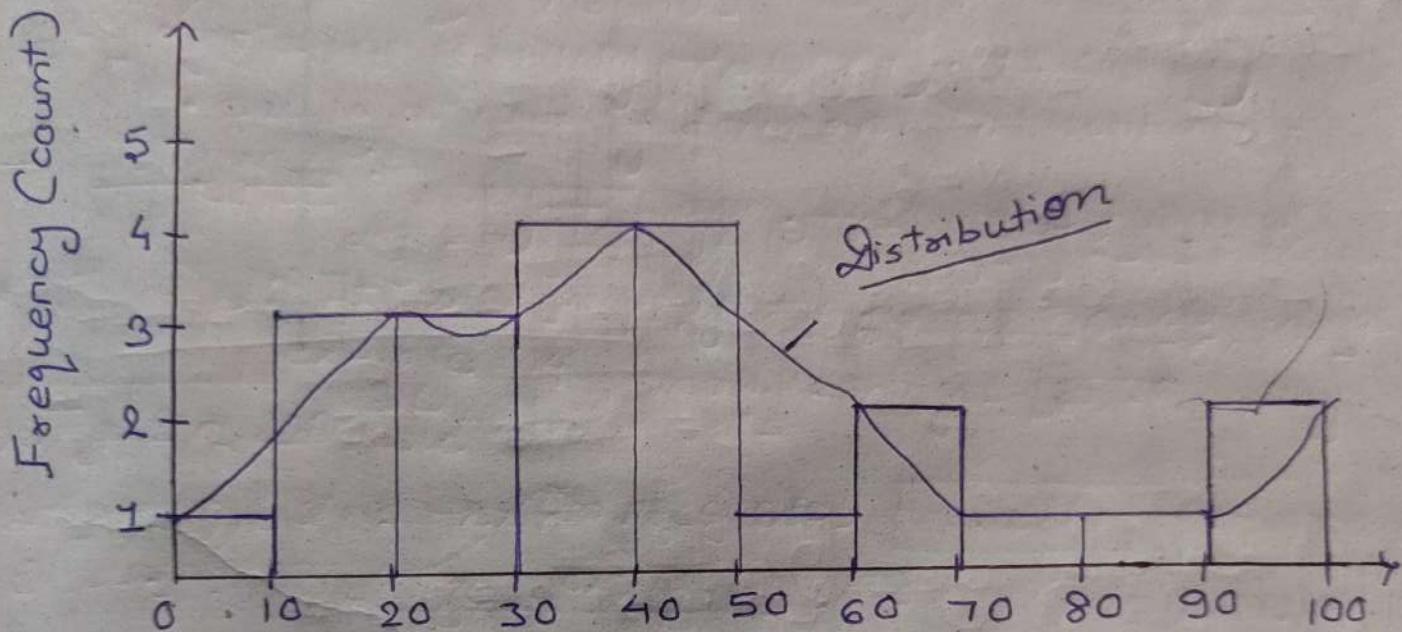
$$\frac{40 - 10}{10} = \frac{30}{10} = 3$$



Example - 1

$$\text{Ages} = \{10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100\}$$

$$\left\{ \begin{array}{l} \text{binsize} = 10 \\ \text{binsize} = \frac{100}{10} = 10 \end{array} \right. \quad \left. \begin{array}{l} \text{(Max)} \\ \text{(Min)} \end{array} \right\}$$



Example - 2

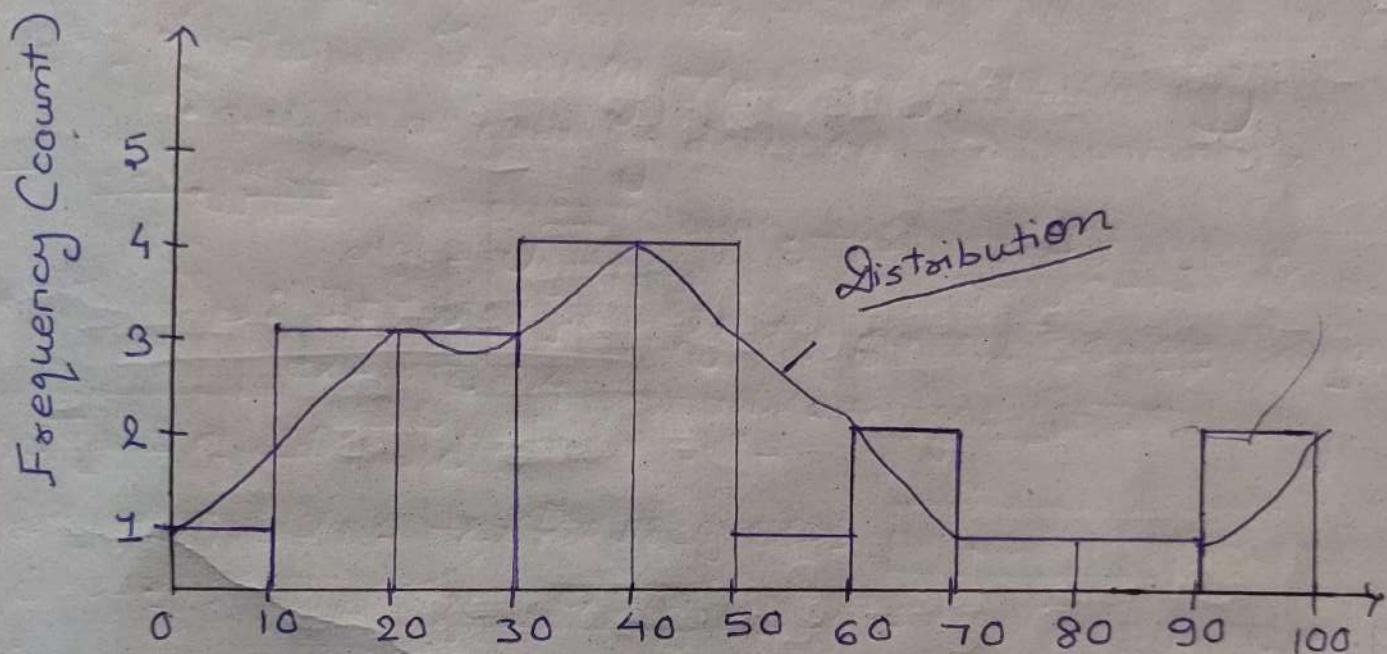
$$\text{Weight: } \{30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 77, 80, 90, 95\}$$

$$\text{bins} = 10, \text{ bin size} = \frac{95-30}{10} = 6.5$$

Example - 1

Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

$$\left\{ \begin{array}{l} \text{binsize} = 10 \\ \text{binsize} = \frac{100}{10} = 10 \end{array} \right. \quad \begin{array}{l} \xrightarrow{\text{Max}} \\ \xleftarrow{\text{Min}} \end{array}$$



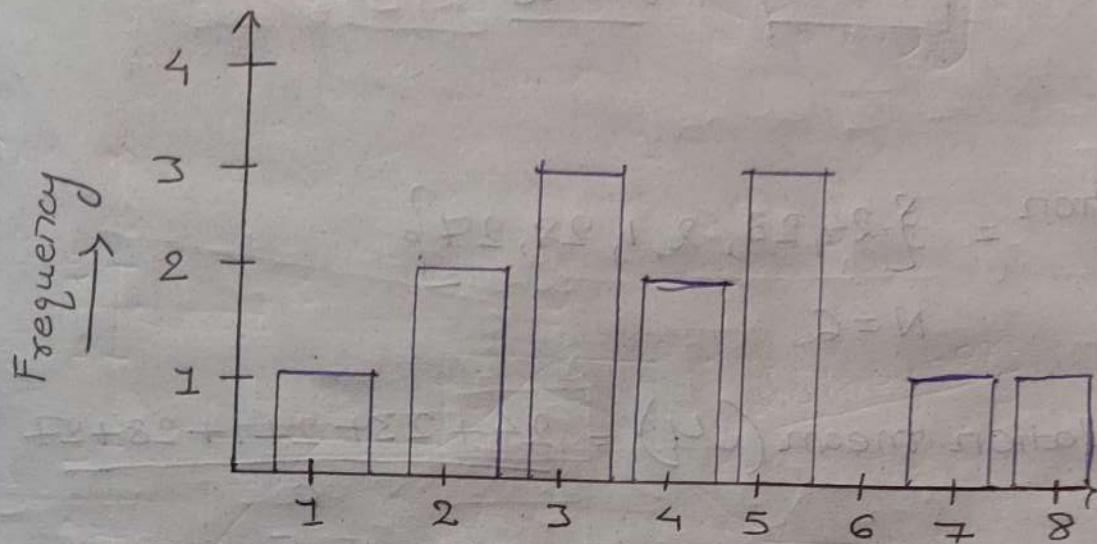
Example - 2

Weight: {30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 77, 80, 90, 95}

$$\text{bins} = 10, \text{ bin size} = \frac{95-30}{10} = 6.5 =$$

## Q) Discrete Continuous

No. of bank accounts = {2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5}



Pdf: Probability density function  $\rightarrow$  Continuous

pmf: Probability mass function  $\rightarrow$  Discrete

\* Measure of Central Tendency:

- (i) Mean
- (ii) Median
- (iii) Mode

A measure of central tendency is a single value that attempts to describe a set of data identifying the central position.

(i) Mean:  $x = \{1, 2, 3, 4, 5\}$

$$\text{Avg / Mean} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

$\rightarrow$  Population(N)

Population mean ( $\mu$ ) =

$$\boxed{\sum_{i=1}^N \frac{x_i}{N}}$$

# \* Practical Applications (Feature Engineering)

<u>Age</u>	<u>Salary</u>	<u>Family Size</u>
-	-	-
-	82	-
NAN	-	-
-	-	-
-	NAN	-
NAN	-	NAN

← Loss of Info

## ② Median

→ Steps to find the median:

- ① Sort the numbers.
- ② Find the central number

{ if the no. of elements are even we  
find the average of central elements }

{ if the no. of elements are odd we find  
the central elements. }

$$\{ 1, 2, 3, 4, 5, 6, 7, 8, 100, 120 \}$$

③ Mode : Most frequent occurring elements.

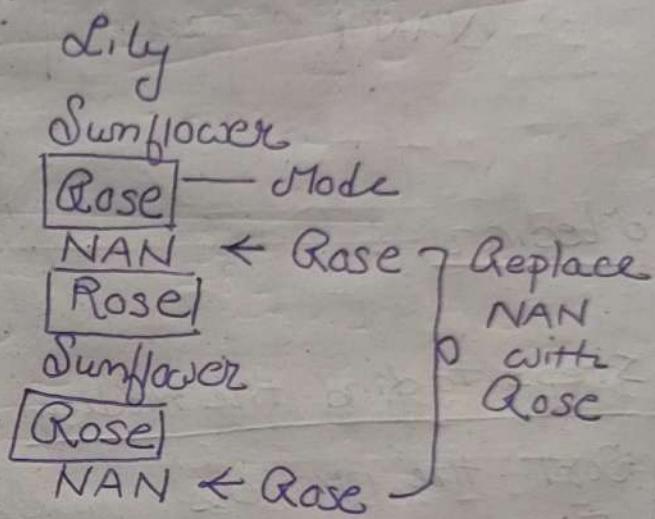
ex:

$$\{1, 2, 2, \boxed{3, 3, 3}, 4, 5\} \text{ or } \{1, 2, 2, 2, 3, 3, 3, 4, 5\}$$

$\boxed{2, 3}$  Mode

DataSet:

Types of Flower:



\* Measure of Dispersion:

- (i) Variance ( $\sigma^2$ ) → Spread of Data
- (ii) Standard deviation ( $\sigma$ )

Variance  $\rightarrow$  Population Variance ( $\sigma^2$ )  
 $\rightarrow$  Sample Variance ( $s^2$ )

(i) Population Variance :-

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

(ii) Sample Variance :

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

Ex-1 {1, 2, 3, 4, 5}

$$\text{Mean } (\mu) = \frac{1+2+3+4+5}{5} = 3$$

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}$$

$$\sigma^2 = \frac{4+1+0+1+4}{5} \Rightarrow \frac{10}{5} = 2$$

Ex-2

{1, 2, 3, 4, 5, 6, 80}

$$\text{Mean } (\mu) = \frac{1+2+3+4+5+6+80}{7} = \frac{101}{7} \Rightarrow 14.4$$

$$\sigma^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + (3-14.4)^2 + (4-14.4)^2 + (5-14.4)^2 + (6-14.4)^2 + (80-14.4)^2}{7}$$

$$\sigma^2 = 719.10$$

Note  
→ As variance increases spread of data would be higher.

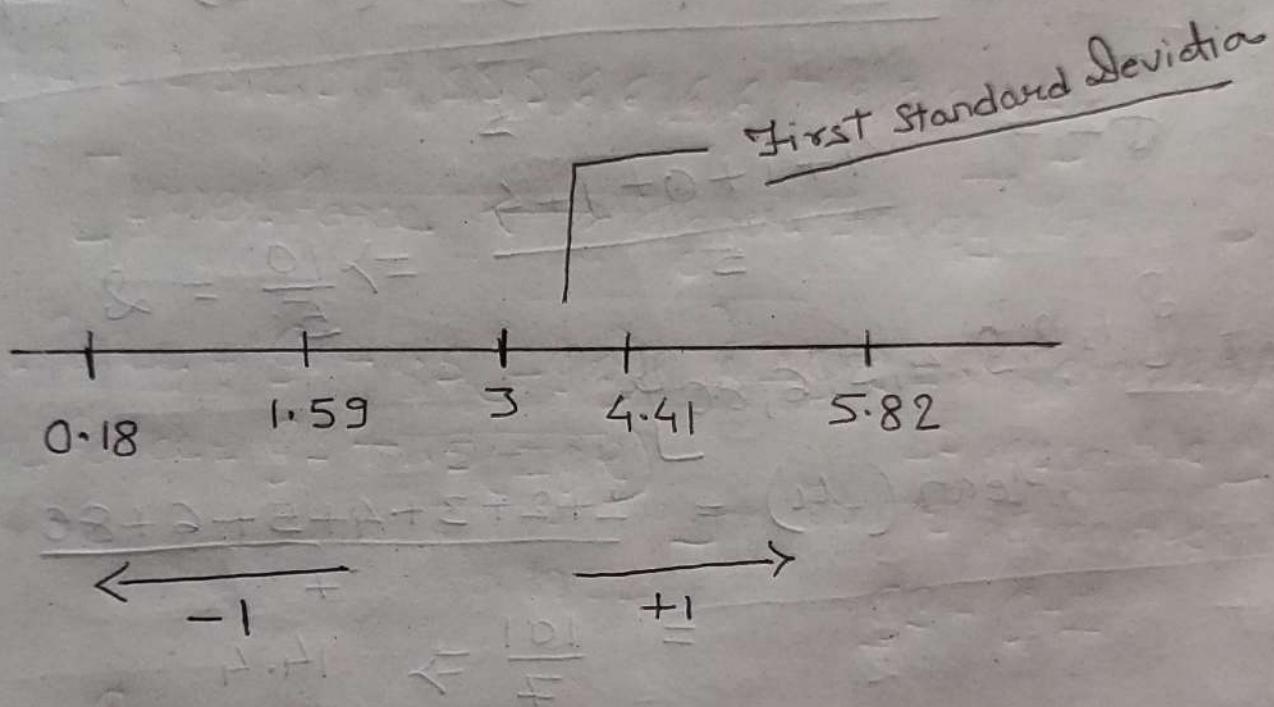
## (ii) Standard Deviation ( $\sqrt{\sigma^2}$ )

ex:-  $\{1, 2, 3, 4, 5\}$

$$\mu = 3$$

$$\sigma^2 = 2$$

$$\sigma = \sqrt{2} \Rightarrow 1.41$$



\* Percentile & Quartiles :-

Percentage =  $\{1, 2, 3, 4, 5, 6, 7, 8\}$

$$\text{Percentage of even no.} = \frac{\text{No. of even no.}}{\text{Total no. of numbers}}$$
$$= \frac{4}{8} = 0.5 = 50\%$$

Percentiles : GATE, CAT, SAT, GRE, JEE, NEET  
are getting in percentiles

Percentile is a value below which a certain percentage of observation lies.

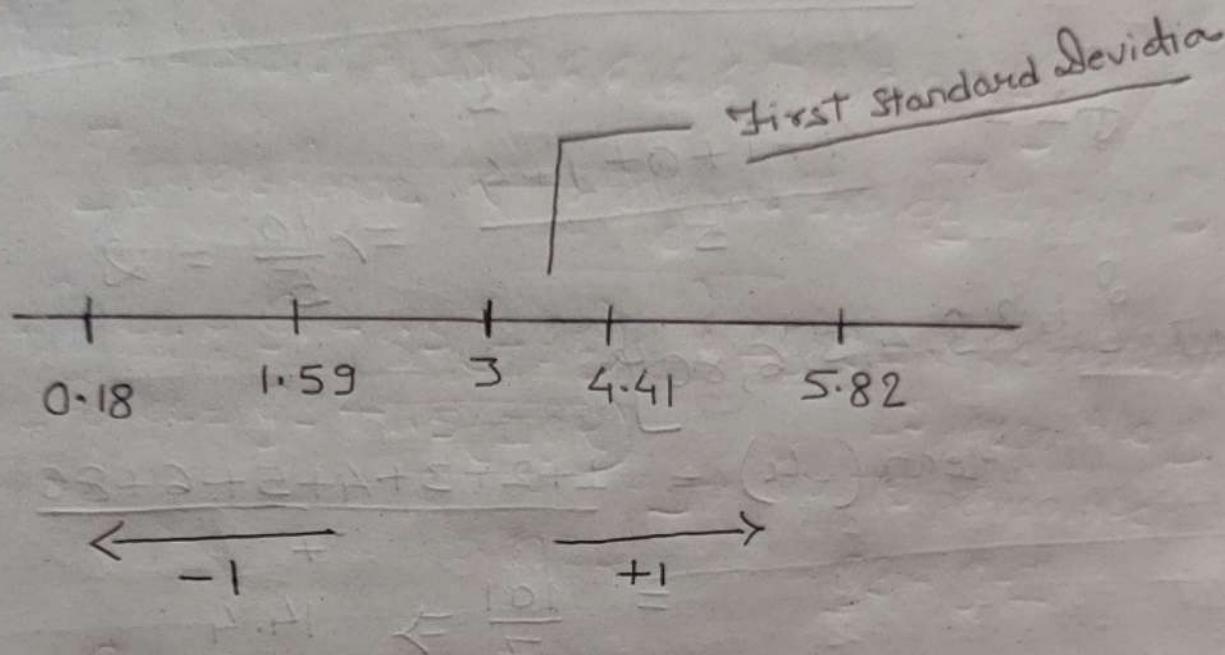
(ii) Standard Deviation - ( $\sigma$ )

ex:-  $\{1, 2, 3, 4, 5\}$

$$\mu = 3$$

$$\sigma^2 = 2$$

$$\sigma = \sqrt{2} \Rightarrow 1.41$$



\* Percentile & Quartiles :-

$$\text{Percentage} = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

$$\begin{aligned}\text{Percentage of even no.} &= \frac{\text{No. of even no.}}{\text{Total no. of numbers}} \\ &= \frac{4}{8} = 0.5 = 50\%\end{aligned}$$

Percentiles: GATE, CAT, SAT, GRE, JEE, NEET  
are getting in percentiles

Percentile is a value below which a certain percentage of observation lies.

99 percentile = It means the person has got better marks, than 99% of the entire students.

Dataset : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

Q. Find the percentile rank of 10.

Ans:

$$\text{Percentile Rank of } x = \frac{\text{No. of value below } x}{n}$$
$$= \frac{16}{20} = 80 \text{ percentile}$$

Q - What is the value that consists at 25 percentile.

$$\text{Value} = \frac{\text{Percentile}}{100} * n \quad \left\{ \begin{array}{l} \text{for even } (n+1) \\ \text{for odd } (n) \end{array} \right.$$

$$= \frac{25}{100} \times 20 = 5^{\text{th}} \text{ index}$$

Hence, 5<sup>th</sup> index is '5' in the dataset.

## \* 5 number Summary :

- (i) Minimum
- (ii) First Quartile
- (iii) Median
- (iv) Third Quartile
- (v) Maximum

Dataset

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27\}$$

[Lower Fence  $\longleftrightarrow$  Higher Fence]

$$\begin{aligned} \text{Lower Fence} &= Q_1 - 1.5(IQR) \\ \text{Higher Fence} &= Q_3 + 1.5(IQR) \end{aligned} \quad \left. \begin{array}{l} IQR = \text{Difference b/w} \\ = 25 \text{ percentile \&} \\ 75 \text{ percentile.} \end{array} \right.$$

$$Q_3 - Q_1$$

$$Q_1 = \frac{25}{100} \times 21 = 5.25 \quad \text{Index} = 3$$

$$Q_3 = \frac{75}{100} \times 21 = 15.75 \quad \text{index} = \frac{8+7}{2} = 7.5$$

$$\text{Lower Fence} = 3 - (1.5)(4.5) = -3.65$$

$$\text{Higher Fence} = 7.5 + (1.5)(4.5) = 14.25$$

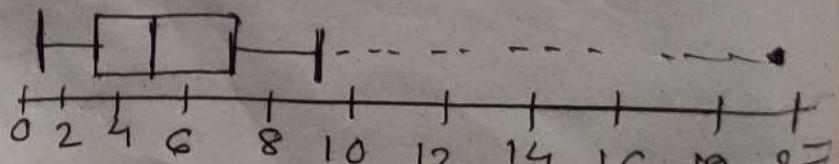
① Minimum = 1

(i)  $Q_1 = 3$

(iii) Median = 5

(iv)  $Q_3 = 7.5$

(v) Maximum = 9



## ② Standard Normal Distribution:

$X \approx$  Gaussian Distribution ( $\mu, \sigma$ )  
 [Converting to  $y$ ]  $\downarrow$   
 $y \approx$  SND ( $\mu = 0, \sigma = 1$ )

Ex:  $x = \{1, 2, 3, 4, 5\}$

$\mu = 3, \sigma = 1.41$  (Applying Z-score)

### ③ Z-Score:

$$\frac{x_i - \mu}{\sigma} \quad \left. \right\} \rightarrow \text{Standard Score}$$

$$= \frac{x_i - \mu}{\sigma} \quad ; \text{if } n=1$$

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

$$\therefore \mu = 3, \sigma = 1.414$$

$$= \frac{1-3}{1.414} = -1.414 \quad \left[ y = \{-1.414, -0.707, 0, 0.707, 1.414\} \right]$$

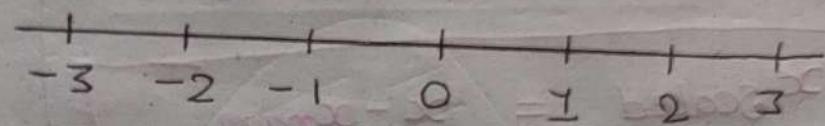
$$= \frac{2-3}{1.414} = -0.707$$

$$= \frac{3-3}{1.414} = 0$$

$$= \frac{4-3}{1.414} = 0.707$$

$$= \frac{5-3}{1.414} = 1.414$$

$$y = \{-1.414, -0.707, 0, 0.707, 1.414\}$$

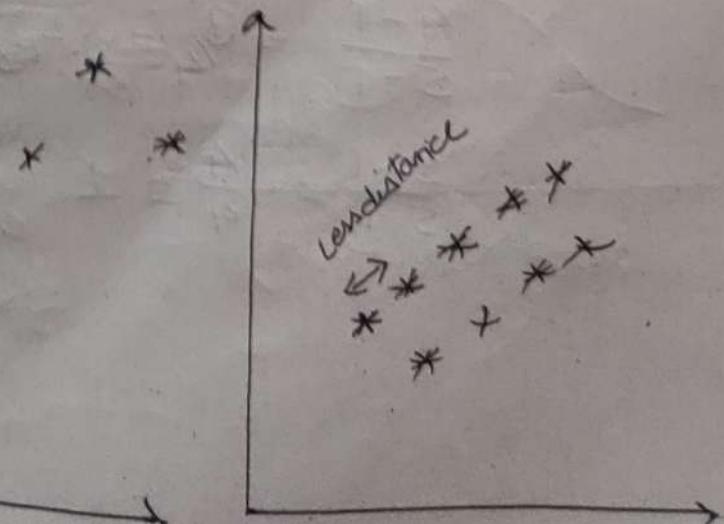
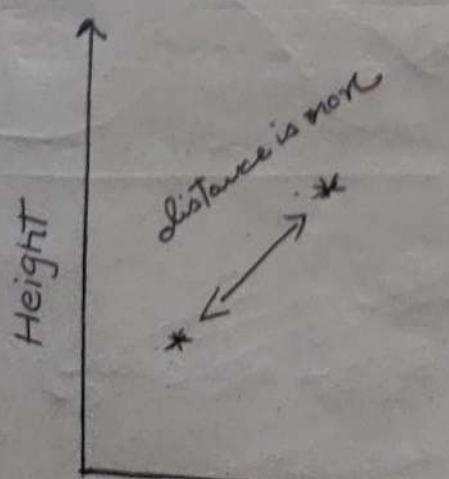


→ Why to Convert?

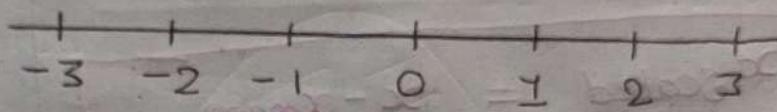
Age (years)	Weight (kg)	Height (cm)
24	72	150
26	78	160
32	84	165
33	92	170
34	87	150
34	83	180
28	80	175
29		

Machine Learning [Maths equation applied]

Algorithm → Mathematical Model



$$y = \{-1.414, -0.707, 0, 0.707, 1.414\}$$

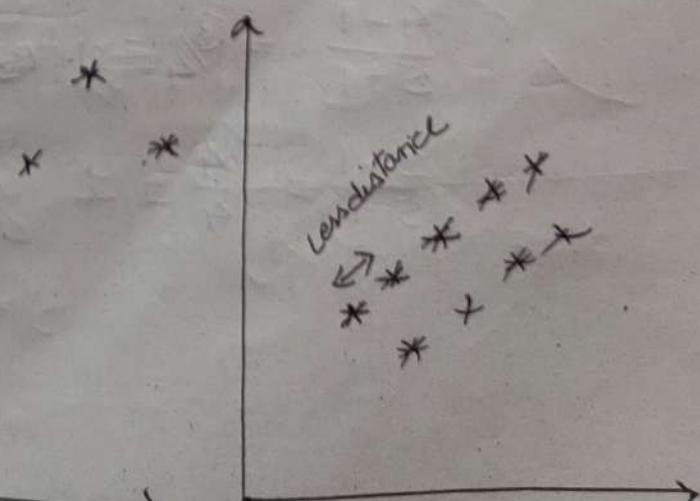
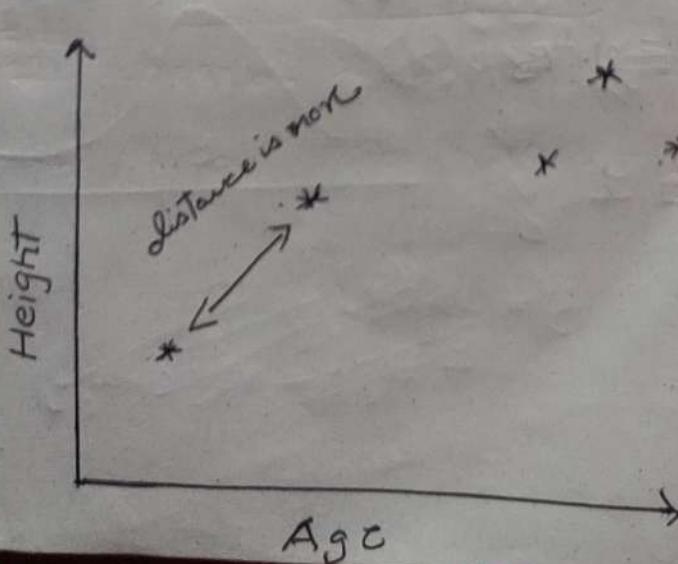


→ Why to convert?

Age (years)	Weight (kg)	Height (cm)
24	72	150
26	78	160
32	84	165
33	92	170
34	87	150
34	83	180
28	80	175
29		

Machine Learning [Maths equation applied]

Algorithm → Mathematical Model



# Z-score rule helps to put data in a same scale, this is called Standardization; where  $\mu(\text{mean}) = 0$  &  $\sigma(\text{standard deviation}) = 1$ .

### \* Feature Scaling :-

\*\* Normalization  $\Rightarrow$  We try to normalize data from lower scale to higher scale.

(i) Min Max Scalar :  $[0-1]$  [Value varies b/w 0-1]

$$x_{\text{scaled}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$

$$\begin{aligned} x &= \{1, 2, 3, 4, 5\} \\ y &= \{0, 0.25, 0.5, 0.75, 1\} \end{aligned}$$

Applying formula :

$$\frac{1-1}{5-1} = 0$$

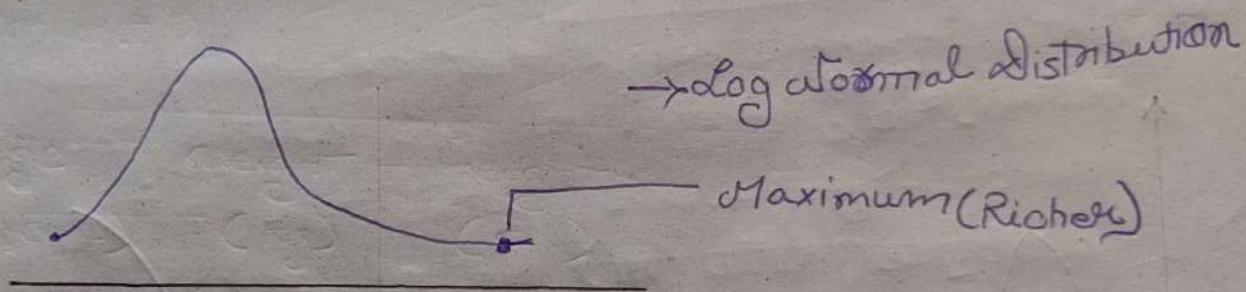
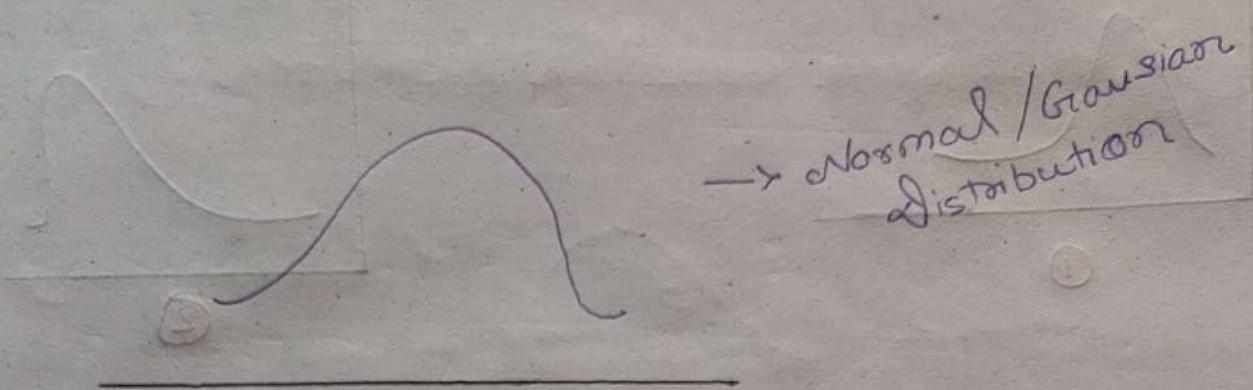
$$\frac{2-1}{5-1} = \frac{1}{4} = 0.25$$

$$\frac{3-1}{5-1} = \frac{1}{2} = 0.50$$

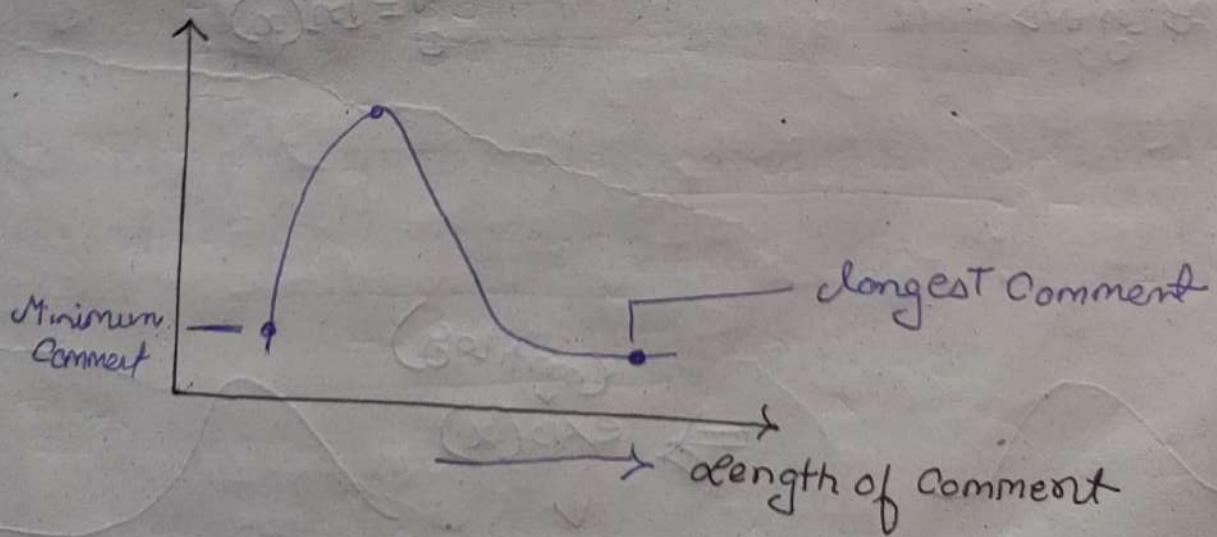
$$\frac{4-1}{5-1} = \frac{3}{4} = 0.75$$

$$\frac{5-1}{5-1} = \frac{4}{4} = 1$$

## \* Log Normal Distribution :

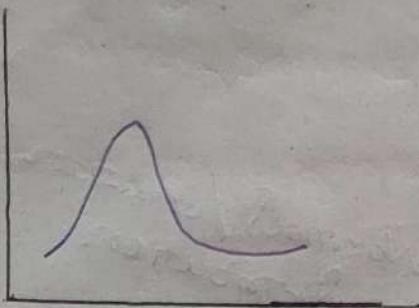


example → Wealth Distribution

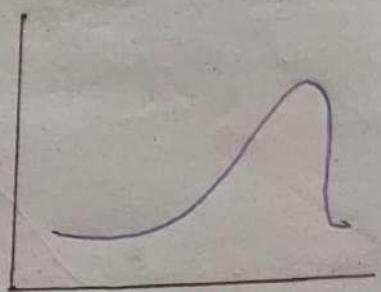


# Assignment

Q: From ascending order give the relation of mean, median, mode from given graph

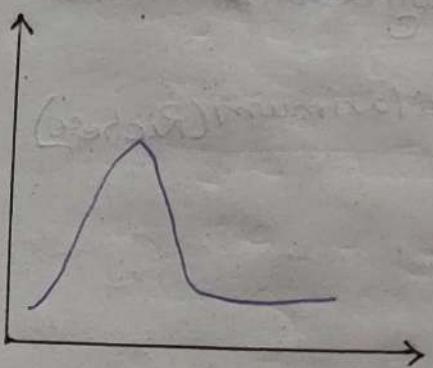


①

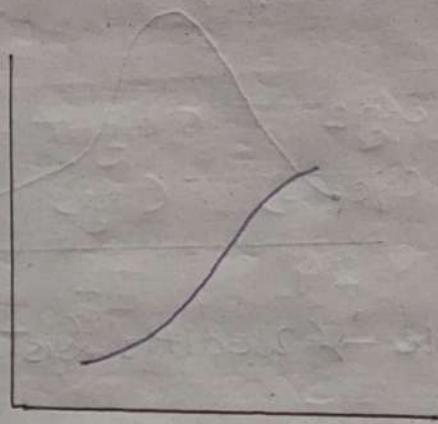


2

\*



$$x \approx \log n$$



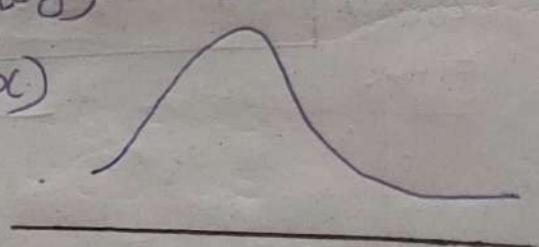
$$y = \ln(x)$$

\*



$$x \sim ND(\mu, \sigma)$$

$$\Rightarrow \text{exp}(x)$$

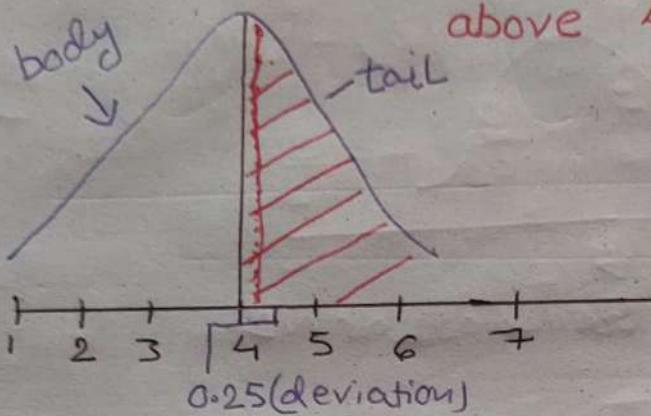


~~$y = \log(x)$~~  ~~Loge (natural log)~~

\* Inverting LND to ND: by applying:  $\log_e(x)$

\* " ND to LND: by " :  ~~$\exp(x)$~~   $\exp(x)$

Q  $x = \{1, 2, 3, 4, 5, 6, 7\}$ ; if  $\mu = 4$  &  $\sigma = 1$ , then what is the percentage of scores that falls above 4.25?



$$Z\text{score} = \frac{4.25 - 4}{1} = 0.25$$

Ans 0.25 deviation from mean.

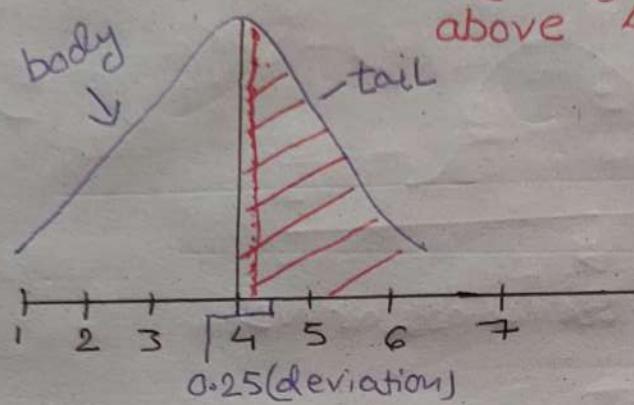
\* Z-table (Area under the table): There are two types of Z-table: (i) Negative Z-table (ii) Positive Z-table or Left-Z-table & Right Z-table respectively.

→ According to Z-table area under Z-score  $< 0.2$  would be 59% then area of the tail will be  $1 - 59\% = \boxed{0.41\%}$

Ans

- \* Inverting LND to ND: by applying:  $\text{Log}(x)$
- \* " ND to LND: by " :  ~~$\exp(x)$~~   $\exp(\ln(x))$

Q.  $x = \{1, 2, 3, 4, 5, 6, 7\}$ ; if  $\mu = 4$  &  $\sigma = 1$ , then what is the percentage of scores that falls above 4.25?



$$Z\text{score} = \frac{4.25 - 4}{1} = 0.25$$

Ans 0.25 deviation from mean.

\* Z-table (Area under the table): There are two types of Z-table: (i) negative Z-table (ii) Positive Z-table or Left-Z-table & Right Z-table respectively.

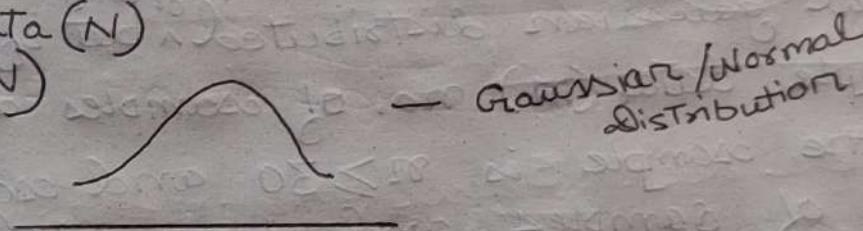
→ According to Z-table area under Z-score  $\geq 0.2$  would be 59% then area of the tail will be  $1 - 59\% = \boxed{0.41\%}$  Ans

## \* Agenda :

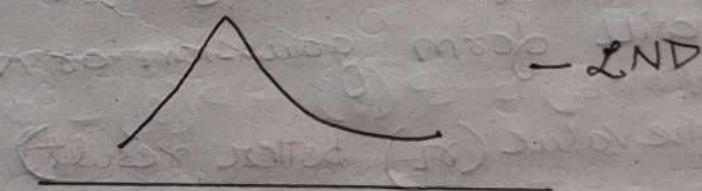
- ① Central Limit Theorem
- ② Probability
- ③ Permutation & Combination
- ④ Covariance, Pearson Correlation, Spearman Rank Correlation
- ⑤ Bernoulli Distribution
- ⑥ Binomial Distribution
- ⑦ Power law {Perito distribution?}

## ① Central Limit Theorem:

Population Data ( $N$ )  
(all graph are  $N$ )



- Gaussian/Normal  
Distribution



- LND

$N \geq 30$



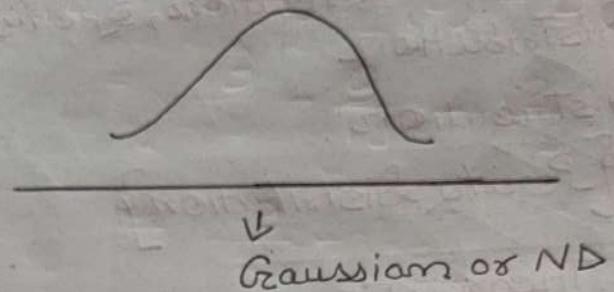
- Left Skewed

Sample Data ( $n$ )  $\leftarrow$  extracting mean from sample data

$$\left\{ \begin{array}{l} \bar{x}_1 = \{x_1, x_2, x_3, \dots, x_n\} \rightarrow \bar{x}_1 \\ \bar{x}_2 = \{x_1, x_2, x_3, \dots, x_n\} \rightarrow \bar{x}_2 \\ \bar{x}_3 = \{x_1, \dots, x_n\} \rightarrow \bar{x}_3 \\ \bar{x}_m = \{x_1, \dots, x_n\} \rightarrow \bar{x}_m \end{array} \right\} \text{sample mean } (\bar{x}_{cm})$$

$n$   
 ↓  
 Size of sample

$m$   
 ↓  
 Number of samples  
 thus, will form Gaussian distribution



② Central theorem says that whether your population data ( $n$ ) is gaussian normal distributed or is not normal gaussian distributed but if we take sample data ( $m$ ) no. of samples and if size of the sample is  $n \geq 30$  and some or no. of samples and if we take mean of all the samples and put it in histogram this will form gaussian or ND.

→ (Larger the value ( $n$ ) better result.)

③ Probability : Probability is a measure of the likelihood of an event.

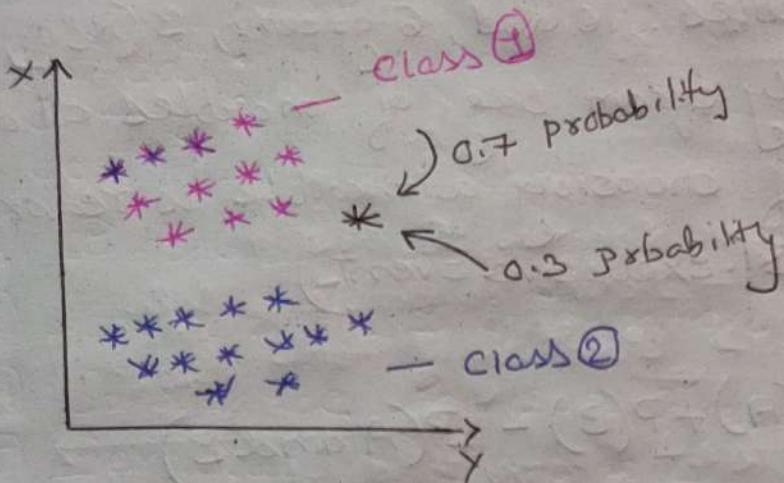
example : (i) Tossing a fair coin  $P(H) = 0.5$   
 $P(T) = 0.5$

(ii) Rolling a dice =  $P(1) = \frac{1}{6}$

$$P(2) = \frac{1}{6}$$

$$P(3) = \frac{1}{6}$$

$$P(4) = \frac{1}{6}$$



(i) → Mutual Exclusive Event : Two events are mutual exclusive, if they cannot occur at the same time.

Ex : ① Tossing a coin      ② Rolling a dice

(ii) → Non-mutual-exclusive event : Two events can occur at the same time.

Ex : (i) Picking a card from a deck of cards, two events "heart" and "King" can be selected  
 (ii) Bag of marble

Q - What is the probability of a coin landing on heads or tails?

Ans :

$$P(A \text{ or } B) = P(A) + P(B)$$

$$= \frac{1}{2} + \frac{1}{2} = 1$$

Q - What is the probability of getting 1, or 6 or 3 while rolling a dice?

Ans :

$$P(1 \text{ or } 6 \text{ or } 3) = P(1) + P(6) + P(3)$$

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{1}{2}$$

Q- What is the probability of rolling a "5" & then a "3" with a normal 6-side dice?

Ans:

$$P(1) = \frac{1}{6}, P(2) = \frac{1}{6}$$

$$P(A) \& (B) = P(A) * P(B)$$

$$= \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$

Multiplication  
Rule of  
independent event

#  $P(A \text{ or } B) \Rightarrow$

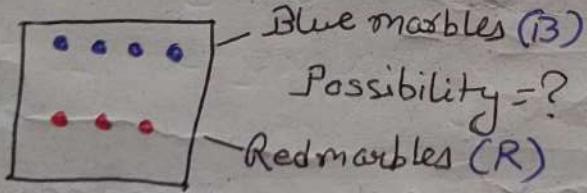
- Mutual Exclusive (ME)
- Non-Mutual Exclusive (NME)

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B) \rightarrow \text{NME}$$

$$P(A \text{ or } B) = P(A) + P(B) \rightarrow \text{ME}$$

→ Dependent Event : used in naive Bayes in ML.

ex : Bag of marbles  $P(B) = \frac{4}{7} \rightarrow$  after picking one blue marble



$$P(R/B) = \frac{3}{6} = \frac{1}{2}$$

$$P(B \text{ and } R) = P(B) * P(R/B)$$

$$= \frac{4}{7} * \frac{3}{6} = \frac{2}{7}$$

→ Independent Event :

$$P(A \text{ and } B) = P(A) * P(B)$$

Ex: Tossing a coin, rolling a dice

$\frac{1 \times 2 \times 3 \times 4 \times 5}{5!} = \frac{120}{120} = 1$

## \* Permutation:

Q:- School of children

{Dairy milk, kitkat, Milkybar, Snickers, 5-star}

$$5 * 4 * 3$$

Picking random chocolate

Possibility of total way will be

$$= 60 \text{ ways} \rightarrow (\text{Permutation})$$

→ All the possible arrangement.

$$\begin{matrix} n \\ p \\ r \end{matrix}$$

where;

$n$  = Total no. of objects

$r$  = no. of selections

$$\frac{n!}{(n-r)!}$$

then, going to upper school based ques:-

$$\frac{5!}{(5-3)!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{2 \times 1}$$

$$= 60$$

## \* Combination:

Repetition will not occur, only unique combination occurs.

$$\begin{matrix} n \\ C_R \end{matrix}$$

$$\frac{n!}{r!(n-r)!}$$

$$\text{going to upper ques} = \frac{5 \times 4 \times 3 \times 2 \times 1}{3! \times 2}$$

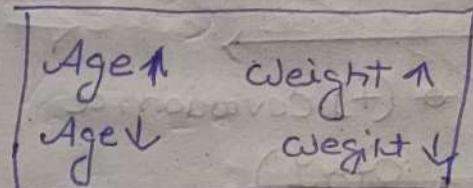
$$= 10$$

## \* Covariance :

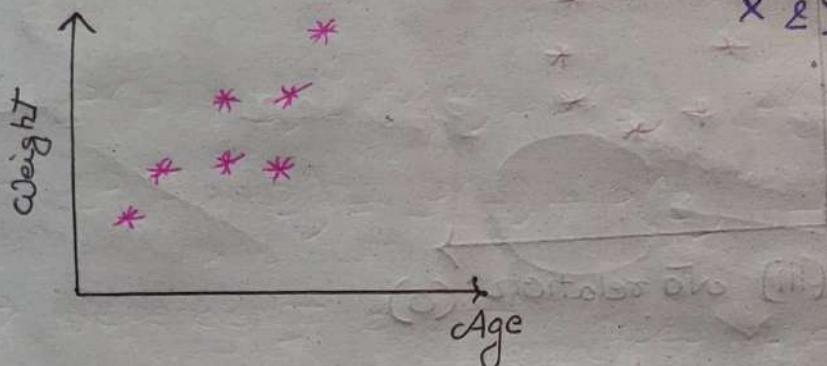
ex:-

<u>Age</u>	<u>Weight</u>
12	40
13	45
15	48
17	60
18	62

{ Feature Selection }



↑  
Quantity the relationship  
 $x \& y$  using

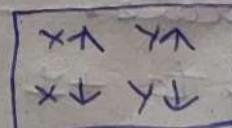


$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

it is similar to variance's formula.

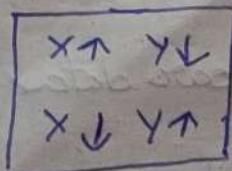
$$\left\{ \begin{array}{l} \text{Cov}(x, x) = \text{Var}(x) \end{array} \right\} \rightarrow \text{important for interview}$$

if cov is (+) then,



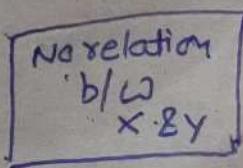
= equally proportional

if cov is (-) then,



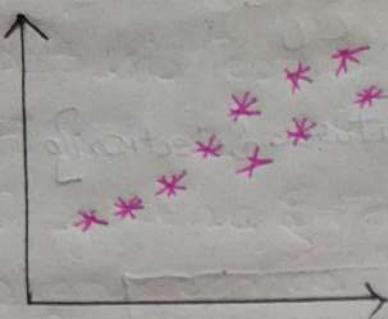
= inversely proportional

if cov is (0) then.



= no relation

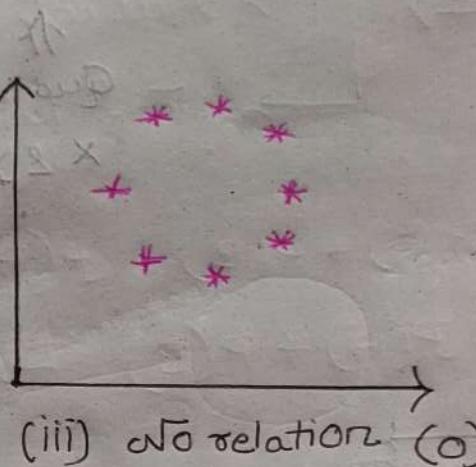
Graph → Covariance through graph:



(i)  $\rightarrow (+)$  covariance  
( $0 \text{ to } 1$ )



(ii)  $\leftarrow (-)$  covariance  
( $-1 \text{ to } 0$ )



(iii) no relation (0)

\* Pearson Correlation Co-efficient :  $(\rho)$  ~~now~~

$$\rho(x,y) = \frac{\text{Cov}(x,y)}{\sqrt{x} * \sqrt{y}}$$

# Restricting the range of covariance b/w  $(-1 \text{ to } 1)$  -  
defining the scale.

- move the value to  $+1 \rightarrow$  more  $(+)$  correlated.
- move " " " "  $-1 \rightarrow$  "  $(-)$  "

→ It only hold for linear data.

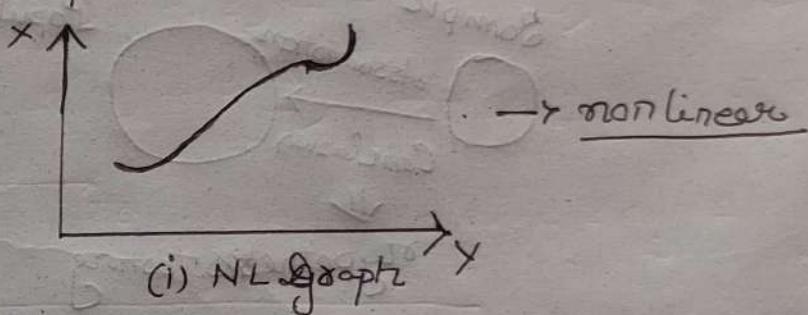
\* Spearman Rank Correlation: St only hold for ~~linear~~ data  
St assigns ranks.

non-linear

$$\gamma_S = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))}$$

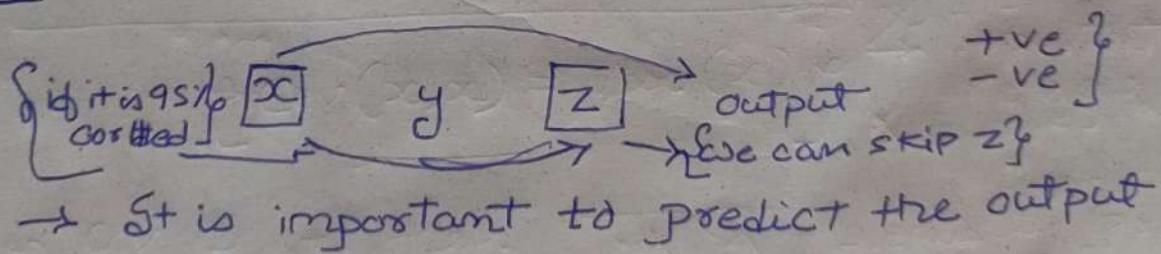
ex-(i) Rank follow in ascending order:

x	y	$R(x)$	$R(y)$
10	4	4	1
8	6	3	2
7	8	2	3
6	10	1	4

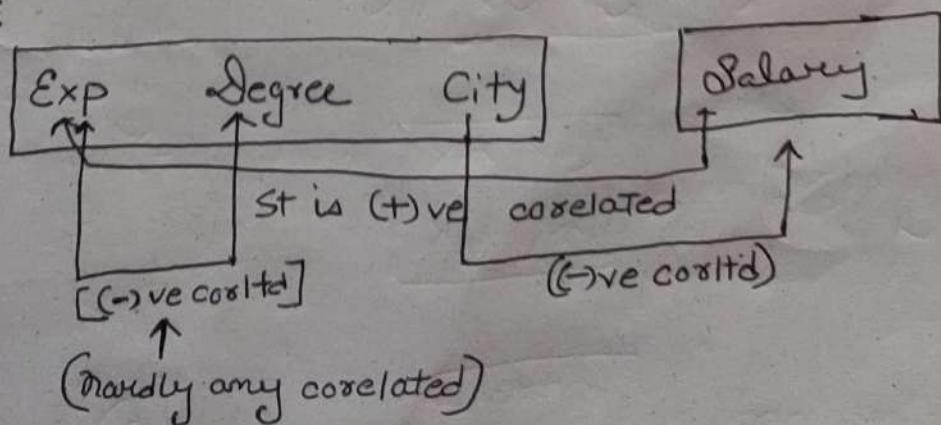


Q Why this correlation is used?

Ans



Exmp:



# We are not losing data by this.