

## 1 Intro

## 2 Introduction To linear algebra

4 - 28

→ Introduction, Scalars & vectors, Addition of vectors

Multiplication of vectors, Vector Databases - Examples

& cosines Similarity, Vectors Multiplication-Elements

wbc multiplication, Vectors Multiplication - Scalar multiplication

Introduction to matrices & Application, Matrix operation

## 3 Introduction to functions and Transformation

29 - 45

→ Introduction To functions and linear Transformation, Vector Transformation

Linear transformation, Why linear transformation fits visualization,

vector length & vector unit & introduction of projection

## 4 Inverse Functions or Transformation

46 - 55

→ Inversion function, applications of functions & inverse function

& How to find inverse of a matrix

## 5 Eigen vectors and Eigen Values

55 - 58

All you need to know about it

## 6 Equation of a Line, plane, hyperplane

59 - 61

(i) (ii)

## 7 Introduction to statistics

→ Introduction to statistics, types of statistics, population

+ Sample data, types of data, types of Sampling

and Scales measurement of data

## 8 Descriptive Statistics

76-95

→ Measure of central tendency, measure of dispersion, why

Sample variance is divided by  $N-1$ , Random variables

percentile &amp; Quartiles, 5 number of summary, histogram

+ Skewness and correlation & ~~skewness~~ Covariance

## 9 Introduction to probability

96-98

+ Addition &amp; Multiplication rule

## 10 Probability function &amp; types of distribution

99-125

→ PDF, PMF, CDF, Types of probability distribution

(Bernoulli, Poisson, Normal Gaussian Distributions)

Standard Normal distribution z-score, uniform distribution

Log normal Distribution, power law distribution,

Pareto Distribution, central limit theorem,

Estimates

(ii)

## ⑪ Inferential stats & Hypothesis Testing

→ Hypothesis Testing and its mechanism, p-value, z-test  
 hypothesis testing, Student T distribution, T-stats with  
 t test hypothesis Testing, Z test vs T-test, Type I & Type II error  
 Baye's theorem, Confidence interval And Margin of Error

126 - 137

## ⑫ Chi Square Test with solved Examples 137 - 141

What is chi-square & Goodness of fit of it

## ⑬ Anova Test with solved Examples 142 - 148

→ Anova - Assumption of Anova, types of Anova  
 partitioning of variance in Anova

## ⑭ Differential calculus

148 - 154

→ Slopes & how to calculate slope, intro to derivatives

Mathematics Notation of Derivatives with limits,

- Finding a Derivative at a point

## ⑮ Power rules & Derivative rules

155 - 159

→ Derivative rules - constant, sum, difference & scalar multiplication

Equation of Tangent of polynomials, Derivatives of Trigonometry

Logarithmic & exponential functions

(iii)

(16) Product Rules in Derivative

159

(17) Chain Rule of Derivatives

160 - 167

Chain rule of derivatives, Composition of 3 or many functions

→ Application of chain rule of derivative

(N) Application of linear algebra, stats & differential calculus

168 - 174

in Data Science

→ Simple linear regression, understanding LR equations,

Cost function of regression, Convergence Algorithm,

Multiple linear regression, performance metrics,

Overfitting & underfitting

(19) Application of linear algebra in Dimensionality Reduction

175 - 183

→ Curse of Dimensionality, Feature Selection & feature extraction

PCA Geometric intuition & Maths intuition

Cross Decomposition on Covariance Matrix

(20) Application of Derivatives in Deep learning Neural network

184 - 192

→ Perceptron & its working with advantages & disadvantages,

ANN working, Back propagation weight update using

derivatives & chain rule of derivative during  
back propagation

(iv)

# What we will learn

3 main important thing

① Linear Algebra

② Statistics (Basics for advanced)

③ Different Calculus

\* Applications of these 3 Topics in Data Science

1 Linear Algebra :- Scalar, Vectors, Vectors operations, matrices, matrix operations  
functions, linear Transformations, inverse functions, Eigen Values  
and Eigen vectors

Neural network :- Forward propagation  $\rightarrow$  Matrix Operations

Applications in Data Science

2 Statistics :-  $\rightarrow$  ML, Deep learning  $\rightarrow$  models  $\Rightarrow$  Huge Dataset  
 $\rightarrow$  Tools to learn from the data

# Statistics

Descriptive

Inferential

① Measure of central tendency

① Hypothesis Testing & P value

② Measure of Dispersion

② Z-test, t-test

③ Histograms, Box plot

③ Square Test

④ Types of distributions of data

④ ANNOVA Test

⑤ PDF, PMF, Normal Distribution, Lognormal

③ Differential Calculus :

① Derivatives, Slope

$\Rightarrow$  Visual Diagrams  $\Rightarrow$  Deriving equations.

② Tangent lines

③ Polynomial Expression

[Derivative of this Expression]

④ Trigonometric Expression

⑤ Chain rule of derivative

⑥ Composite function

} optimizations

$\Rightarrow$  chain rule

# Applications of linear Algebra, Statistics & Differential calculus in Data Sciences

- ① Simple Linear Regression , Multiple Linear Regression  $\leftarrow$  applications
- ② Dimensionality Reduction [ Principal Component Analysis ]  $\rightarrow$  Eigen values { Eigen vector }
- ③ Neural n/w Trained  $\rightarrow$  ANN  $\rightarrow$  Multi - Layered  $\downarrow$   
Artificial Neural Network

## Linear Algebra :-

Linear algebra is a branch of mathematics that focuses on the study of vectors, vector spaces [linear spaces], linear transformations and systems of linear equations. It provides a framework for understanding the properties and operations of these mathematical objects, which can be represented using matrices and vectors.

- \* It is a foundational concepts for ML, DL, NLP & Images
- \* It has more applications in Physics & mathematics.
- \* The aim is to learn Linear algebra as Computer Science Student [WS]

## Application of Linear Algebra:-

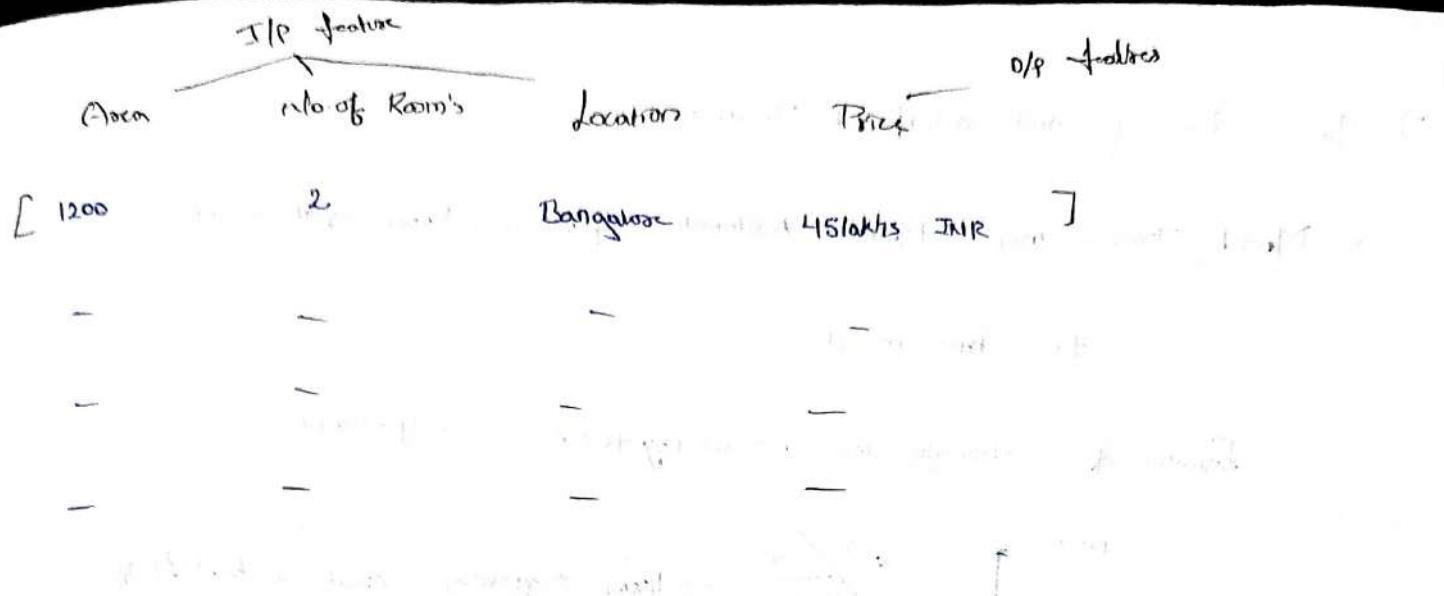
### ① Data representation & manipulation.

Dataset: House price dataset

Create a model which will be able to predict house price

Vector :- Is a list of numbers

Data is represented to model as vectors



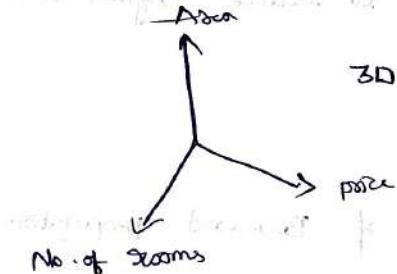
\* The model can quantify the relationships like  $x, y$

### Covariance

गोपदाता

$x \uparrow$   $y \uparrow$   
 $x \downarrow$   $y \uparrow$   
 $x \uparrow$   $y \downarrow$   
 $x \downarrow$   $y \downarrow$

} covered in  
= stats



\* Linear algebra works well with higher dimension data

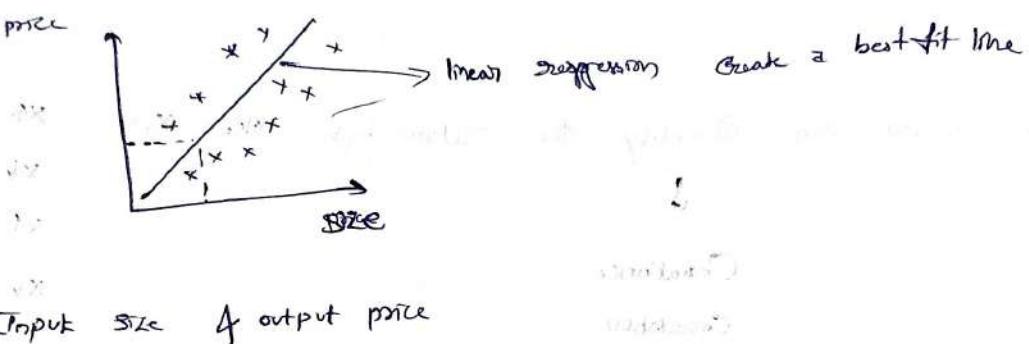
500 dimension

!! using LA concepts like Dimensionality reduction.  
2 dimensions

## ② Machine Learning and Artificial Intelligence :

- \* Model Train :- uses "Matrix Arithmetic operations, linear equations etc." to train model

Equation of a straight line  $\Rightarrow ax + by + c = 0 \Rightarrow y = mx + c$

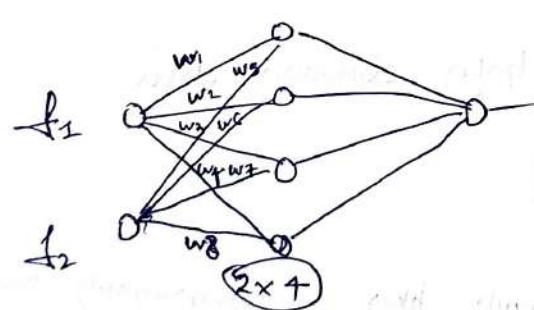


- \* Dimensionality reduction :- PCA uses linear algebra concepts like

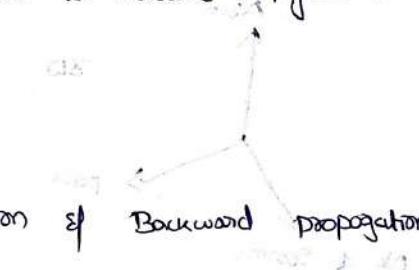
- Eigen value & Eigen vector

With this help we will be able to reduce higher dimension to lower dimension.

- \* Neural N/w :- Forward propagation & Backward propagation



$$\begin{bmatrix} f_1 \\ f_2 \end{bmatrix} \begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_5 & w_6 & w_7 & w_8 \end{bmatrix} \Rightarrow \text{Matrix multiplication}$$



Train neural network with help of GPU  $\rightarrow$  cost,

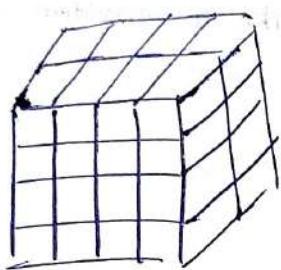
GPU train fast because parallel MO & MM

Tensorflow  $\rightarrow$  Tensor  
works

## Computer Graphics

Image

0-255 pixels



RGB image

- \* Used to perform operations such as

Scaling, rotating image, white & black image.

- \* Linear algebra is used to perform transformation

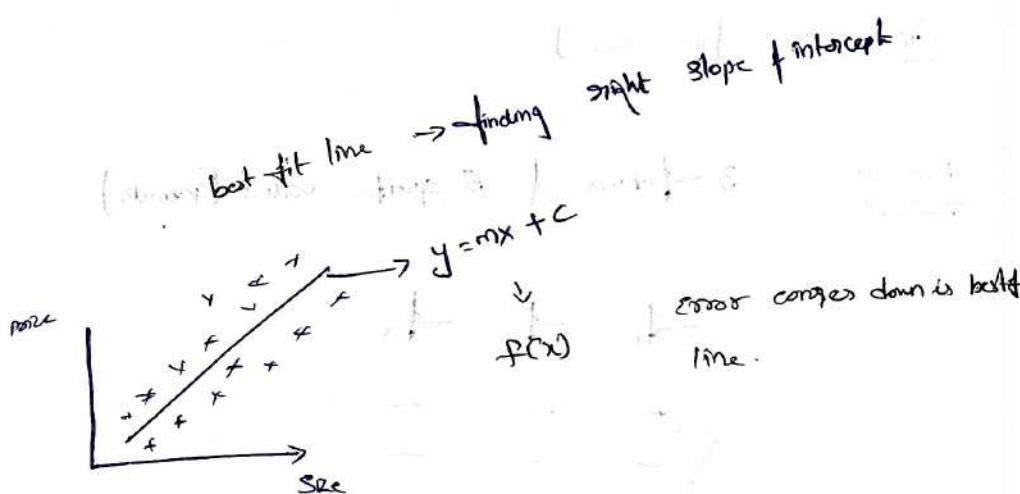
Such as Scaling, rotating & translation of

objects in Computer graphic

④

optimization

① Solving Equation:-



Gradient Descent  $\Rightarrow$  optimizor

### ③ Scalars and Vectors

Scalar :- It is a single numerical value. It represents a magnitude or quantity and has no directions.

- Ex:-
- ① Car speed =  $45 \text{ km/hr}$  Magnitude
  - ② Temperature in Celsius  $T = 25^\circ \text{C}$

Applications :- [in DS]

Dataset :- 3 features + 5 specific values [Records]

	$A_1$	$A_2$	$A_3$
1	-	-	-
2	-	-	-
3	-	-	-
4	-	-	-
5	-	-	-

Count of the Total no. of records = 5

average of feature  $A_1$  = -

} Scalars  
Quantity

Simple linear regression  $\Rightarrow y = mx + c$  intercept

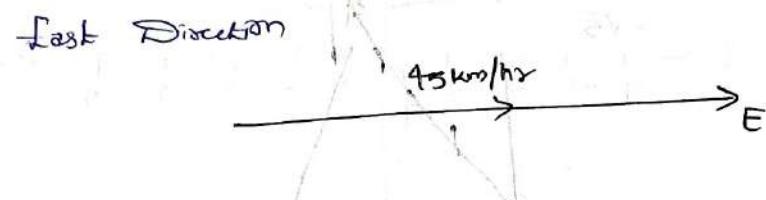
slope L scalar value

Vector: It is a numerical value which has both magnitude & direction

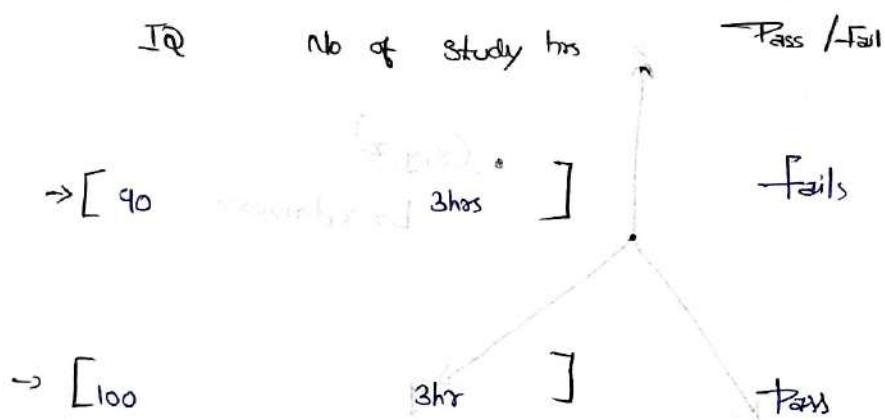
In terms of Data Science,  
↓

A vector is an ordered list of numbers. It can represent a point in space or quantity with both magnitude and direction.

Example: Speed of the car is 45 km/hr and is moving towards East direction



Example: Student marks:



→ A vector representing person I.Q + no of study hrs = [90, 3hrs]

↑  
3 hrs ↑ units  
magnitude

A vector representing persons weight over time  $[70, 72, 75, 73] \leq 4d$

No need to think they have specific directions

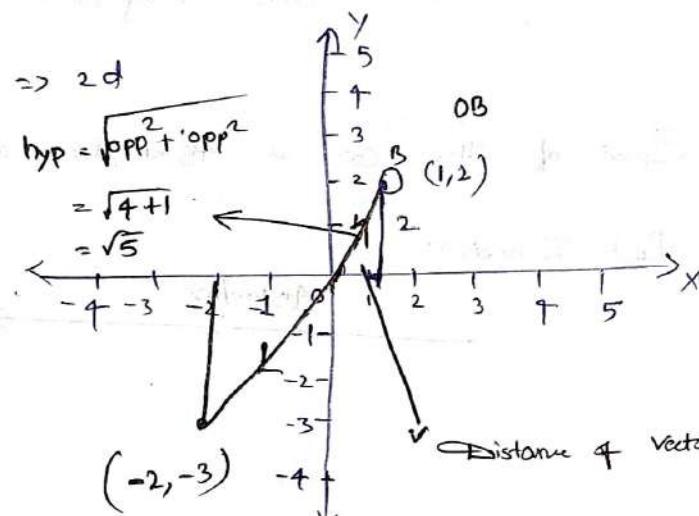
Vectors representing with respect to physics

$$A = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \Rightarrow 2d$$

$$\text{hyp} = \sqrt{\text{opp}^2 + \text{opp}^2}$$

$$= \sqrt{4+1}$$

$$= \sqrt{5}$$

(-2, -3) 

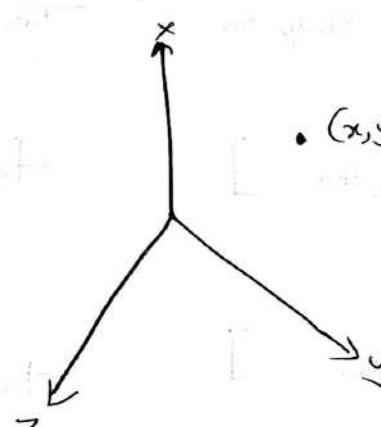
OB

$$B = \begin{bmatrix} -2 \\ -3 \end{bmatrix}$$

$$C = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

(x, y, z)

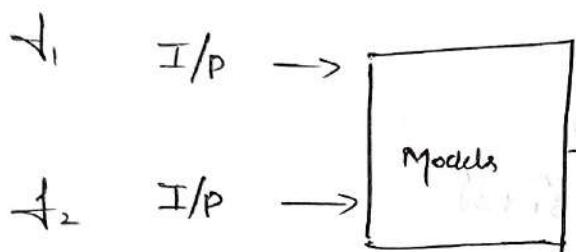
$\rightarrow$  3 dimensions



How vectors are related to Data Science?

In DL, ML we Create Some kinds of models.

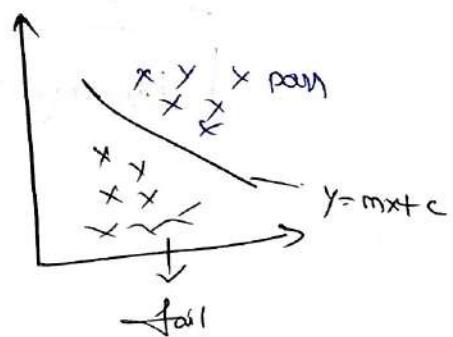
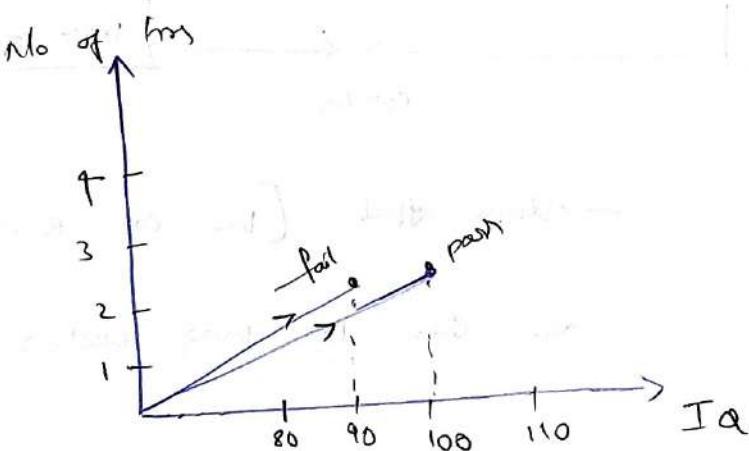
The aim of the model is to take input data & predict the output features



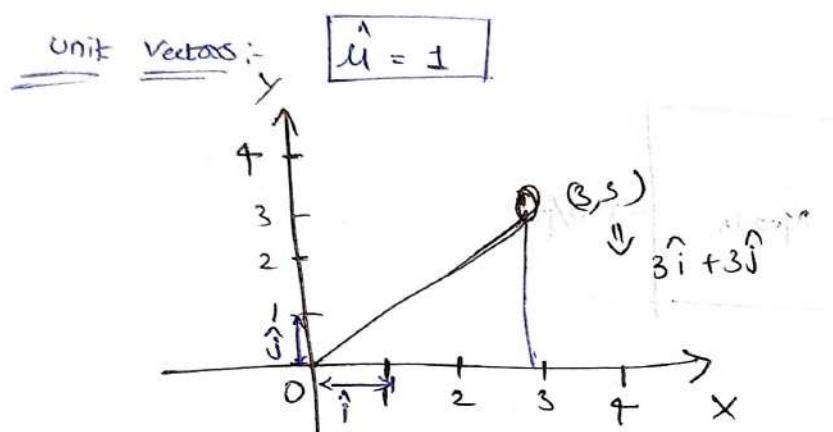
$f_1$        $f_2$       O/P ←  
I<sub>Q</sub>      No. of hours      pass/fail

$$\begin{bmatrix} 90 \\ 2 \end{bmatrix} \quad \text{Fail} \Rightarrow 0$$

$$\begin{bmatrix} 100 \\ 2 \end{bmatrix} \quad \text{Pass} \Rightarrow 1$$

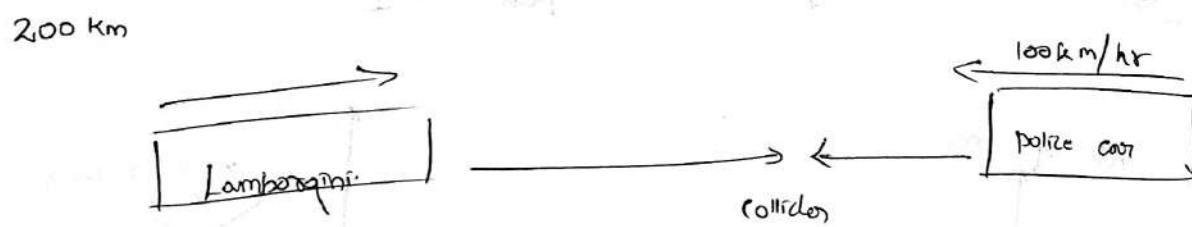


- With human eye we can see max to max 3 features but linear algebra is so amazing that [since it supports lots of matrix calculations] any number of features can be represented using it.



Examples of vectors we used daily specially in gaming industry:

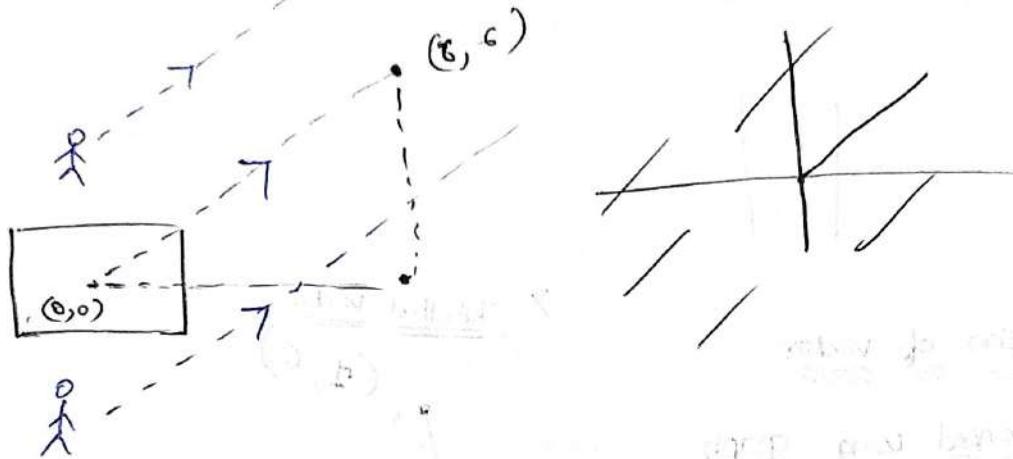
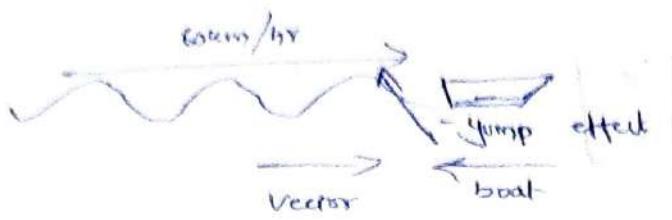
### GTA 6 :-



→ Advers effect [like car blown up, damage, etc.]

are done by using vectors

Boat



If  $(0,0)$  is origin or not that vector can move in same direction

with different origins & it can happen in any part of

Co-ordinate System

## ④ Addition of Vectors:

$$① P_1 = \begin{bmatrix} -4 \\ 3 \end{bmatrix} \quad P_2 = \begin{bmatrix} 5 \\ 3 \end{bmatrix}$$

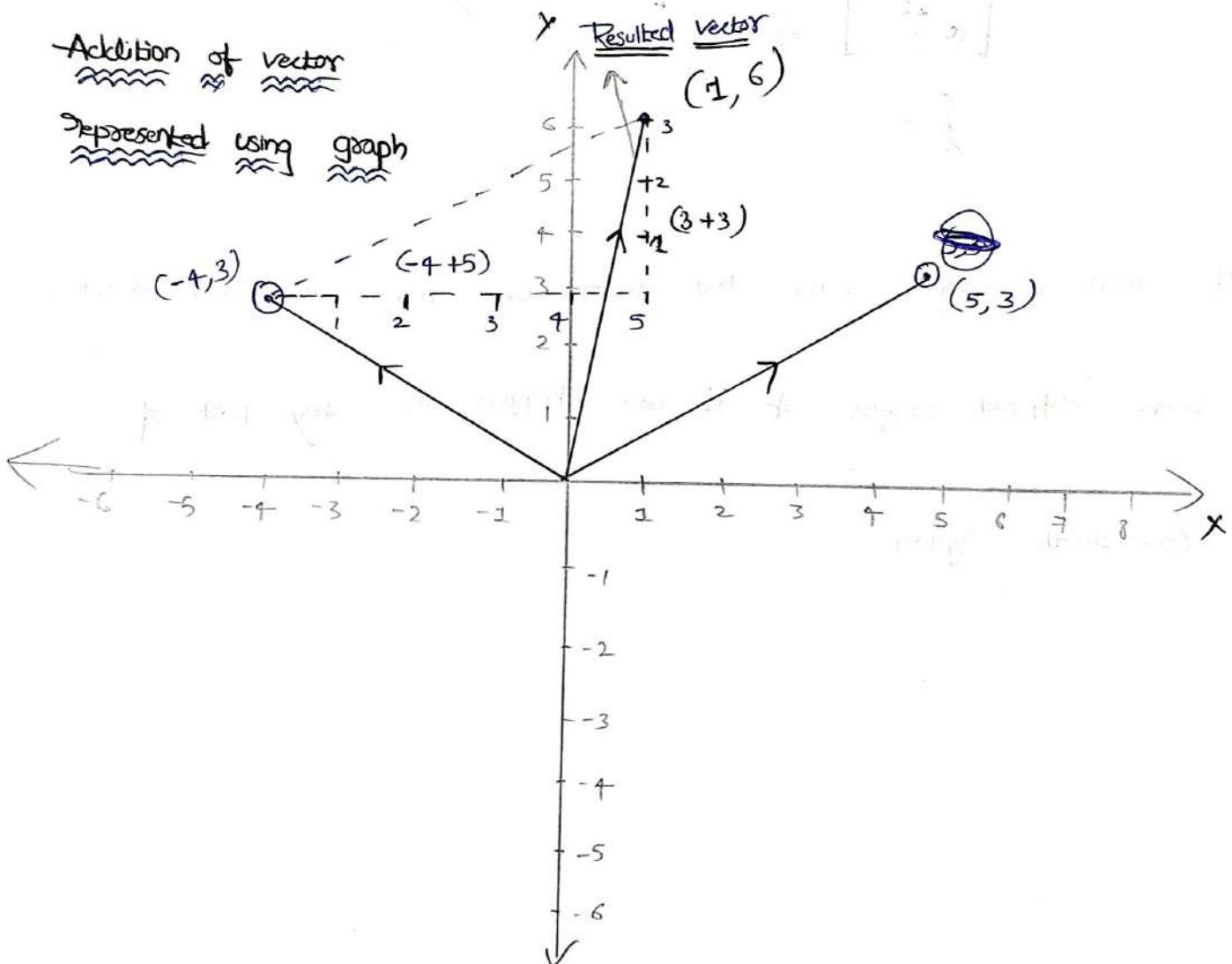
Sol

$$P_1 + P_2 = \begin{bmatrix} -4 \\ 3 \end{bmatrix} + \begin{bmatrix} 5 \\ 3 \end{bmatrix} = \begin{bmatrix} -4+5 \\ 3+3 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 6 \end{bmatrix}$$

Addition of vector

represented using graph



$$A = \begin{bmatrix} x_1 \\ y_1 \\ z_1 \end{bmatrix} \quad B = \begin{bmatrix} x_2 \\ y_2 \\ z_2 \end{bmatrix} \quad A + B = \begin{bmatrix} x_1 + x_2 \\ y_1 + y_2 \\ z_1 + z_2 \end{bmatrix} = \begin{bmatrix} x_3 \\ y_3 \\ z_3 \end{bmatrix}$$

Examples:-

Solving

a use case

in TDA for feature engineering

Sensor 1

Sensor 2

final sensor reading

$$[3, 5, 7]$$

$$[2, 4, 6]$$

$$[3, 5, 7] + [2, 4, 6] = [5, 9, 13]$$

NLP ex:- E-commerce website

Applications

① Data aggregation

The product is good

1

② Feature engineering

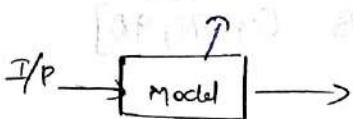
The product is bad

0

③ NLP

not understand text

④ Word Embeddings



⑤ Image processing

The text is converted to vector

Text → Vector [numerical values]

OHE

TFIDF

BOW

Word2Vec

## Word Embedding:

$$\textcircled{1} \text{ Date: } [0.2, 0.1, 0.4]$$

$$\textcircled{2} \text{ Science: } [0.3, 0.7, 0.2]$$

$$\text{Data Science} = \text{Date} + \text{Science}$$

$$= [0.2, 0.1, 0.4] + [0.3, 0.7, 0.2]$$

$$= [0.5, 0.8, 0.6] \Rightarrow \text{Data Science}$$

## Image Processing

$$\text{Color Image } [R, G, B] =$$

- Red Channel  $R = [255, 128, 0]$

$$G = [128, 156, 0]$$

$$B = [64, 78, 90]$$

$$RGB \rightarrow \text{Gray scale} = R + G + B$$

$$= \frac{[255+128+64, 128+156+78, 0+0+90]}{3}$$

white & black

$$= [ ]$$

## 4 Multiplication of Vectors

3 Types

- ① Dot product (Inner product)
- ② Element wise multiplication
- ③ Scalar multiplication

### ① Dot product:

Definition: The dot product of 2 vectors results in a scalar • it is calculated as the sum of the products of their corresponding components.

$$A = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad B = \begin{bmatrix} 4 \\ 5 \end{bmatrix}$$

$$A \cdot B = 2 \times 4 + 3 \times 5$$

$$= 8 + 15$$

$$= 23 = \text{Scalar}$$

$$A \cdot B^T = \begin{bmatrix} 2 \\ 3 \end{bmatrix} [4 \ 5]$$

$$= 2 \times 4 + 3 \times 5$$

$$= 8 + 15$$

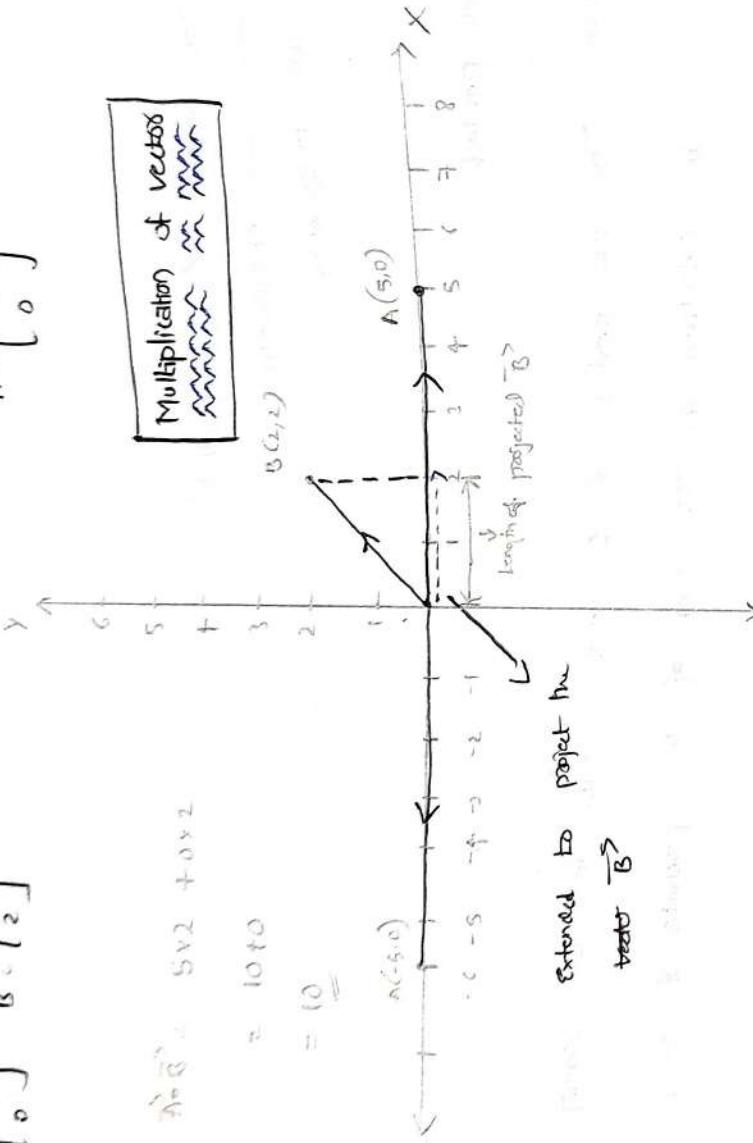
$$= 23 = \text{Scalar}$$

3

$$\textcircled{1} \quad A = \begin{bmatrix} 5 \\ 0 \end{bmatrix} \quad B = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$\textcircled{2} \quad u = \begin{bmatrix} -5 \\ 0 \end{bmatrix}$$

Multiplication of vectors



If the length of projected vector is zero the dot product is zero

$$\vec{A} \cdot \vec{B} = (\text{Length of projected } \vec{B}) \cdot (\text{Length of the vector } \vec{A})$$

$$= (2)(5) = 10 = +ve$$

(Note)

$\vec{A} \cdot \vec{B} = 0$ , project the vector to the origin

# Application of DOT PRODUCT in Data Science

→ Gen AI

→ App Systems like RAG

It is a measure used to determine how similar

## ① Cosine Similarity:

2 vectors are, it calculates cosine of the angle between 2 vectors and provide a similarity score

that ranges from -1 (dissimilar) to 1 (complete similar)

$$\cos \theta = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

watching

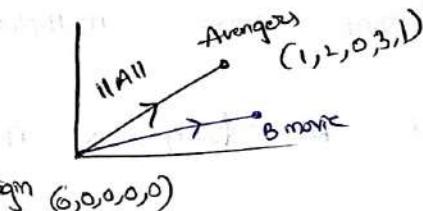
Recommendation System: Netflix account  $\Rightarrow$  Action movie

Average  $\Rightarrow$  [Drama, Action, Comedy, Romance, Story type]

$\Rightarrow$  recommendation of other action movies

graph graph 5D imagination

$$\mathbf{B} \text{ movie } = [2, 0, 1, 1, 1]$$



Step 1: Dot product of  $\mathbf{A} \cdot \mathbf{B}$ ,  $\mathbf{A} \cdot \mathbf{B} = 1 \cdot 2 + 2 \cdot 0 + 0 \cdot 1 + 3 \cdot 1 + 1 \cdot 1 = 6$

$$\text{Step 2: } \|\mathbf{A}\| = \sqrt{1^2 + 2^2 + 0^2 + 3^2 + 1^2} = \sqrt{15} = 3.872$$

$$\|\mathbf{B}\| = \sqrt{2^2 + 0^2 + 1^2 + 1^2 + 1^2} = \sqrt{7} = 2.6457$$

$$\cos \theta = \frac{6}{3.872 \times 2.6457} \approx 0.586$$

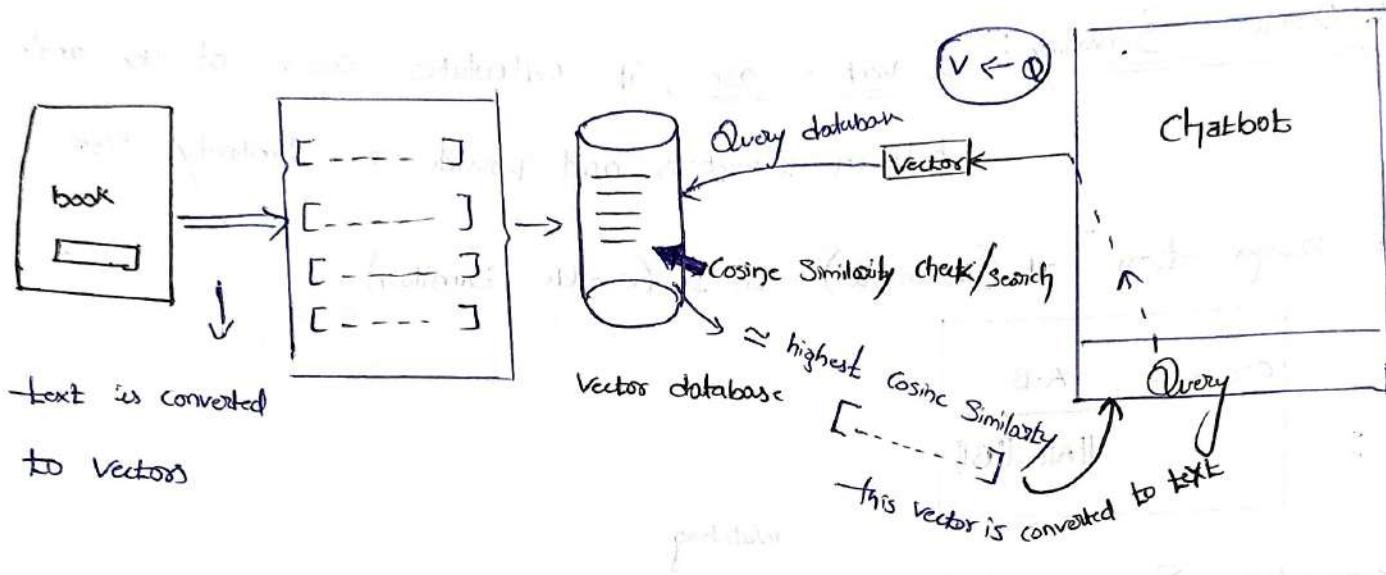
58.6% +ve Similar

$$\|\mathbf{A}\| = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

## Vector Database

→ Frequently used in GenAI & LLM models

⇒ Retrieval augmented generation (RAG) system uses Vector database



## ② Element wise multiplication

In element wise multiplication, corresponding elements of 2 vectors are multiplied to form a new vector of the same dimension.

$$A = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \quad B = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}$$

$$A \otimes B = \begin{bmatrix} 3 \\ 8 \\ 15 \end{bmatrix}$$

## Application in Data Science

### ① Feature Engineering

Product	Cost	Discount	Discount price	Final price
A	1000	0.1	100	900
B	500	0.2	100	400
C	200	0.5	30	170

Deep learning:

RNN, LSTM RNN, GRU RNN

$\otimes \oplus \Rightarrow$  Forget gate, input gate

$$\begin{bmatrix} 0.5 \\ 0.6 \\ 0.3 \end{bmatrix} \otimes \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} 0 \\ 0 \\ 0.3 \end{bmatrix}$$

gate  $\Rightarrow$  pass info or not

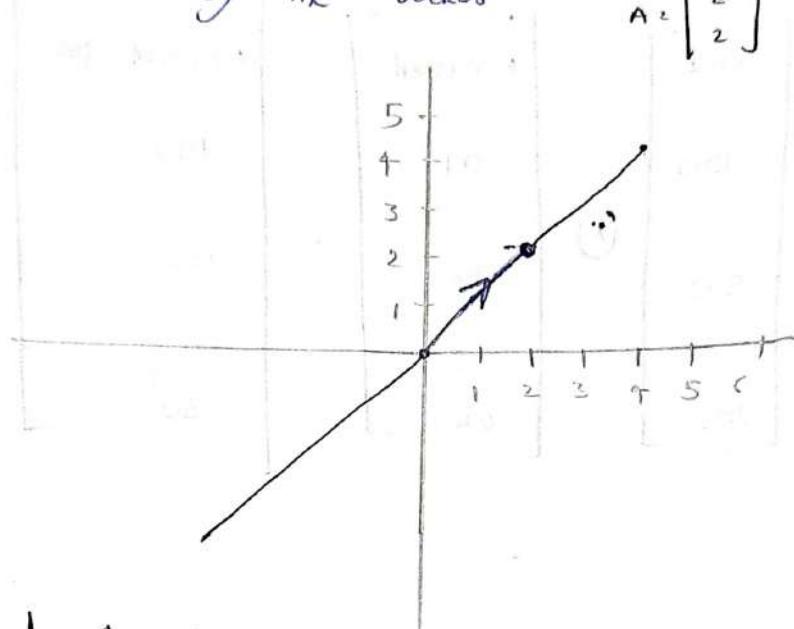
### ③ Scalar multiplication:

It involves multiplying vector by a scalar, resulting in a vector where each component is scaled by the vector.

$$A = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \quad c = 2$$

$$A = \begin{bmatrix} 3 \\ 5 \\ 7 \end{bmatrix} \quad c = 4$$

$$C = \begin{bmatrix} 12 \\ 20 \\ 28 \end{bmatrix}$$



$\Sigma x_i$ : Normalization      4: Standardization  
 ↓

Scaling data  $\Rightarrow$  units

0	-255	0	-255
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.

$\Rightarrow$  Image processing  $\Rightarrow$  Normalize pixel  $\Rightarrow$

↓

$$[0-1]$$

255	128	0	-128
.	.	.	.
.	.	.	.
.	.	.	.
.	.	.	.

If require more calculation so normalization is done

done

$$1\text{cm} = 0.01\text{m}$$

$$Ch = 0.01 [160, 170, 180]$$

$$= [1.6, 1.7, 1.8]$$

e.g.: feature

$$\text{Height} = [160, 170, 180]$$

$$\text{Scale (to meter)} = c = 0.01$$

Matrices: A matrix is a rectangular array of numbers, symbols or expressions arranged in rows and columns.

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}_{m \times n}$$

each element is represented by  $a_{ij}$   
where i denotes rows  
j denotes columns

Example of matrices in Data Science

Dataset

Math Score	Physics Score	Biology Score	$f_1$	$f_2$	$f_3$
→ 55	65	75	→ 55	65	75
→ 65	60	55	→ 65	60	55
→ 70	45	80	→ 70	45	80

$3 \times 3$

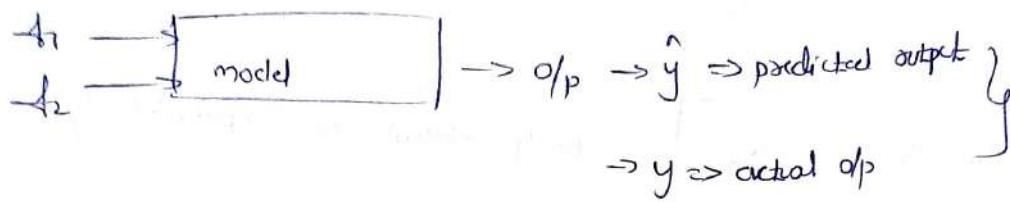
② Image in Computer Vision

$3 \times 3$  Grayscale Image

$$\underline{\underline{I}}_{\text{image}} = \begin{bmatrix} 0 & 128 & 255 \\ 255 & 128 & 0 \\ 128 & 255 & 128 \end{bmatrix}$$

$$\Rightarrow \begin{bmatrix} 0 & 128 & 255 \\ 255 & 128 & 0 \\ 128 & 255 & 128 \end{bmatrix}_{3 \times 3}$$

### ③ Confusion matrix: Accuracy of the model



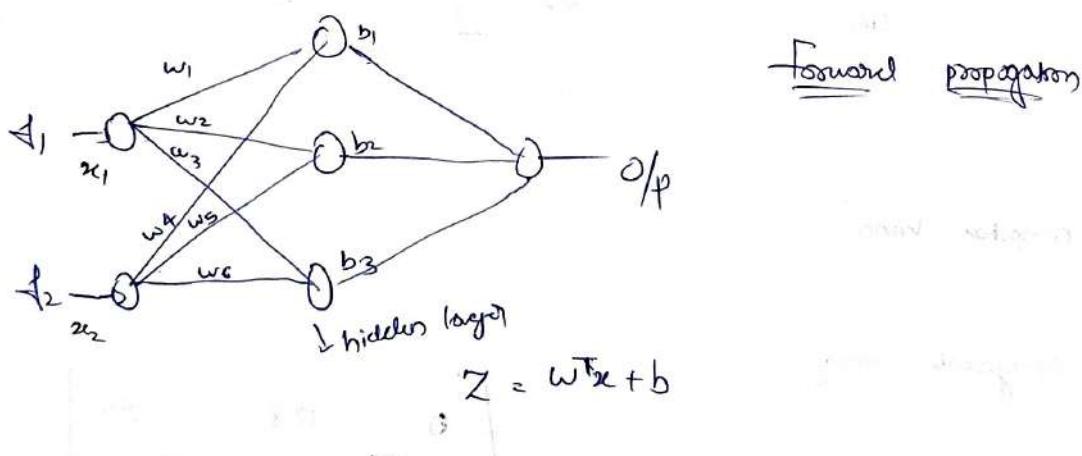
Confusion matrix = 
$$\begin{bmatrix} 50 & 10 \\ 5 & 35 \end{bmatrix}$$

50 → True positive  
 10 → False Negative  
 5 → False Positive  
 35 → True Negative

$$\frac{TP + TN}{TP + FN + FP + TN} = \text{Accuracy}$$

### ④ Neural N/w :

Matrix operation:  $\rightarrow$  used in N/w & linear regressions



$$w = \begin{bmatrix} w_1 & w_2 & w_3 \\ w_4 & w_5 & w_6 \end{bmatrix}$$

$$x = [x_1, x_2]$$

$$b = [b_1, b_2, b_3]$$

## Linear Regression:

No. of Study hours      IQ

$x_1$        $x_2$

Score dependent feature

Score independent feature

hours

Regression

4

100

90

$$y = mx + c$$

5

90

85

$$= m_1 x_1 + m_2 x_2 + c$$

Lslope Lslope or co-efficient

$$\Rightarrow m^T x + c$$

$$\left\{ \begin{array}{l} m = [m_1, m_2] \\ x = [x_1, x_2] \end{array} \right\}$$

## NLP:- Dataset

Review

Sentiment +ve/-ve

The food is good

0

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \end{bmatrix}$$

$$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

The food is bad

1

we can apply different ML algorithms  
to train this data

## Matrix operation:

To manipulate & analyze multidimension data efficiently:

① Matrix Addition & Subtraction

② Scalar Matrix Multiplication

③ Matrix Multiplication.

## ① Matrix operation:

→ To manipulate & analyze data

### ① Matrix - Addition and Subtraction:

Add or Subtract corresponding elements of 2 matrices of the same dimensions.

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix} \Rightarrow \text{STORE A}$$

$3 \times 3$

$$B = \begin{bmatrix} 4 & 5 & 6 \\ 7 & 8 & 9 \\ 1 & 2 & 3 \end{bmatrix} \Rightarrow \text{STORE B}$$

$3 \times 3$

$$A + B = \begin{bmatrix} 1+4 & 2+5 & 3+6 \\ 4+7 & 5+8 & 6+9 \\ 7+1 & 8+2 & 9+3 \end{bmatrix}$$

	Product A	Product B	Product C
Day 1	1	2	3
Day 2	4	5	6
Day 3	7	8	9

$$= \begin{bmatrix} 5 & 7 & 9 \\ 11 & 13 & 15 \\ 8 & 10 & 12 \end{bmatrix}$$

## ② Scalar multiplication:

Scalar multiplication involves multiplying every element of a matrix by a scalar value.

$$B = cA$$

Eg: Suppose we have a matrix representing product prices in dollars and we want to adjust these prices for inflation by a factor of 1.05.

$$\text{Original } P = \begin{bmatrix} 10 & 20 & 30 \\ 15 & 25 & 35 \\ 20 & 30 & 40 \end{bmatrix}$$

Scalar multiplication:

$$P_{\text{adjusted}} = 1.05 \cdot P = 1.05 \begin{bmatrix} 10 & 20 & 30 \\ 15 & 25 & 35 \\ 20 & 30 & 40 \end{bmatrix} = \begin{bmatrix} 10.5 & 21 & 31.5 \\ 15.75 & 26.25 & 36.75 \\ 21 & 31.5 & 42 \end{bmatrix}$$

Used for adjusting salary after inflation:

S/W Eng	HR	Accountant	6%
45k	30k	40k	1.06%
50k	35k	45k	
-	-	-	
-	-	-	

### ③ Matrix multiplication

**Operation:** It involves the dot product of rows of the first matrix with columns of the second matrix.

For 2 matrix  $A(m \times n)$  &  $B(n \times p)$ , the result matrix is  $C(m \times p)$

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix} \Rightarrow 1 \times 2 + 2 \times 3 + 3 \times 4 = 2 + 6 + 12 \\ = [20]$$

$$A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{2 \times 3}$$

$$B = \begin{bmatrix} 7 & 9 & 11 \\ 8 & 10 & 12 \end{bmatrix}_{2 \times 3}$$

$m \times n$

$$B^T = \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix}_{3 \times 2}$$

$n \times p$

$$C = A \cdot B = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{2 \times 3} \begin{bmatrix} 7 & 8 \\ 9 & 10 \\ 11 & 12 \end{bmatrix}_{3 \times 2} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}_{2 \times 2} = \begin{bmatrix} 58 & 64 \\ 139 & 154 \end{bmatrix}$$

$$C_{11} = (1 \times 7) + (2 \times 9) + (3 \times 11) = 58$$

$$C_{12} = (1 \times 8) + (2 \times 10) + (3 \times 12) = 64$$

$$C_{21} = (4 \times 7) + (5 \times 9) + (6 \times 11) = 139$$

$$C_{22} = (4 \times 8) + (5 \times 10) + (6 \times 12) = 154$$

# Functions And Linear Transformations

Functions: A function is a mathematical relationship that uniquely associates elements of one set (called the domain) with elements of another set (called the codomain). In simpler terms, a function maps inputs to outputs in a specific way.

Notation: A function  $f$  mapping elements from set  $X$  (domain) to set  $Y$  (codomain) is denoted by  $f: X \rightarrow Y$ .

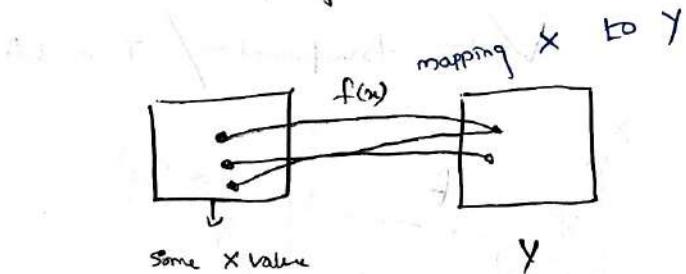
If  $x$  is an element of  $X$ , then  $f(x)$  is the corresponding element in  $Y$ .

Example:  $f(x) = 2x + 3$   $\Rightarrow$  maps each real number  $x$  to a real number

$$x=5 \quad x \xrightarrow{f} y$$

$$f(5) = 2 \times 5 + 3 = 13$$

$f(x) \rightarrow$  mapping  $5 \in \mathbb{R}$  to  $13 \in \mathbb{R}$



$$f: x \xrightarrow{\curvearrowright} y$$

Vector

$$F: \begin{bmatrix} x \\ y \\ z \end{bmatrix} \in \mathbb{R}^3 \xrightarrow{\quad} \in \mathbb{R}^3 \quad f(x) = \begin{bmatrix} xy \\ yz \end{bmatrix}$$



$$f(x)$$

It mapped 3D vector to 2D vector

Eg.: Dimensionality Reduction  $\Rightarrow$  we use functions for transformation

to reduce the dimensions

$$\begin{array}{c} \text{Domain} \\ \uparrow \\ f: \mathbb{R}^3 \rightarrow \mathbb{R}^2 \\ \uparrow \text{Codomain} \end{array}$$

$f$  is function which is mapping 3D real number to 2D real number

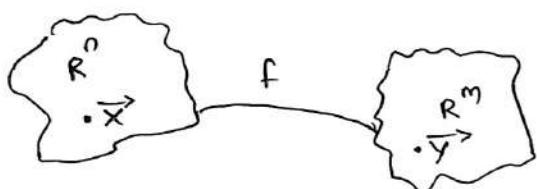
## Vector

### Transformation:

$$f: x \rightarrow y$$

Basis may have different set of vectors

$$\vec{x} \rightarrow \vec{y}$$



can be denoted as

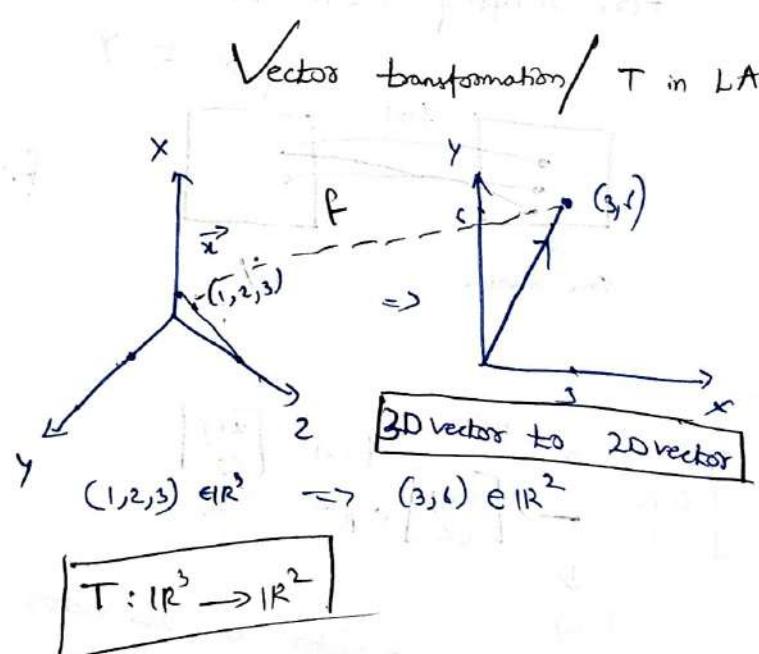
$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \quad x_1, x_2, x_3, \dots, x_n \in \mathbb{R}$$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f(x, y, z) = (x+y, 2z)$$

$$f \left( \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \right) = \begin{bmatrix} 3 \\ c \end{bmatrix}$$

rotation



Definition: Vector transformations refers to operations that maps vectors from one space to another, often changing their magnitude, direction, or both. These transformations are typically described using matrices & are fundamental in various fields, including physics, engineering, computer graphics, and Data Science.

In gaming industry this kind of these Transformations are heavily used.

### ① Ex: Scaling:

Scaling is a transformation that changes the magnitude of vector while keeping their direction same.

Scaling is used for Data normalization & Computer graphics to resize the image.

$\Rightarrow$  open point  $\rightarrow$  image  $\rightarrow$  Resize

$$\text{Scaling matrix} = \begin{bmatrix} k & 0 \\ 0 & k \end{bmatrix}$$

### ② Rotation: Transformation that turns vectors around the origin

$$V = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \in \mathbb{R}^2$$



$$V' = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \rightarrow \text{Showing a 90 degree rotation}$$

Used in  
Deep learning  
image processing  
Robotics

~~Rotation~~: Rotation will be used in image processing  $\Rightarrow$  Rotating image.

Robotics  $\Rightarrow$  adjusting robot orientation

3D graphics  $\Rightarrow$  Rotating objects

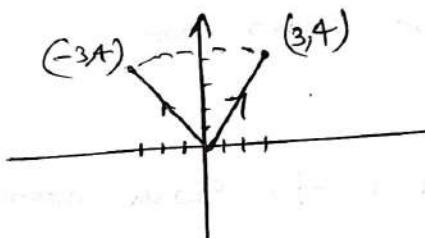
### ③ Reflection

[vector  $T$  is used in it]

Transformation that flips vectors over a specific axis or plane.

$$v = \begin{bmatrix} 3 \\ 4 \end{bmatrix} \Rightarrow \text{Across the } y \text{ axis}$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -3 \\ 4 \end{bmatrix}$$

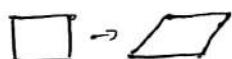


uses

Mirror images

Analyzing wave reflections

### ④ ~~Shearing~~: Shearing



Shearing is a linear transformation that skews the shape of an object by shifting points in a one direction, proportional to their position on another axis. It changes shape without necessarily changing area or volume.

$$V = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad VT = V$$

Horizontal shear:  $\begin{bmatrix} 1 & k \\ 0 & 1 \end{bmatrix} \xrightarrow{k=2} \begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix} \xrightarrow{\cdot \begin{bmatrix} 1+4 \\ 0+2 \end{bmatrix}} \begin{bmatrix} 5 \\ 2 \end{bmatrix}$

Vertical shear:  $\begin{bmatrix} 1 & 0 \\ k & 1 \end{bmatrix} \xrightarrow{\begin{bmatrix} 1 & 0 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}} \begin{bmatrix} 1 \\ 4 \end{bmatrix}$

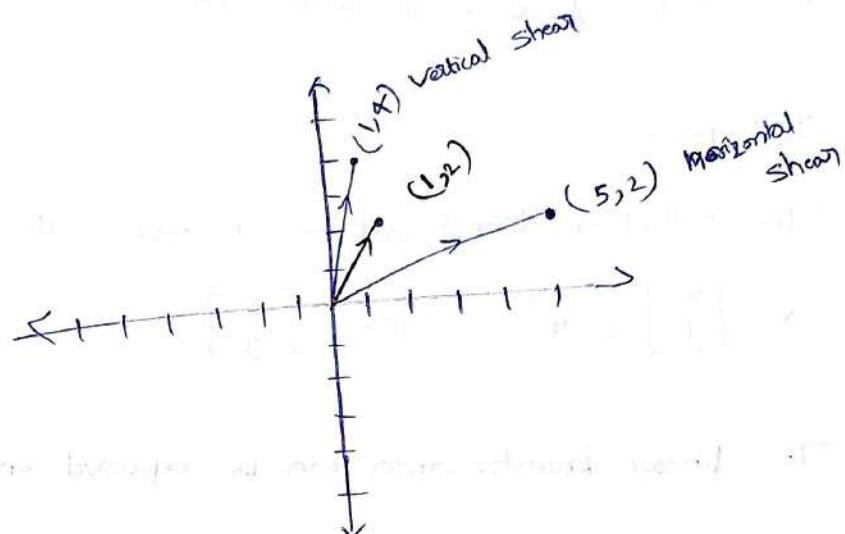


uses:

① Computer graphics & Animation

② ML & Data Augmentation

③ Structural Engineering



### \* Linear Transformation:

A linear transformation is a function between two vector spaces that

preserves the operations of vector addition and scalar multiplication. This

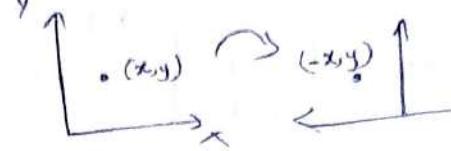
means that if  $T$  is a linear transformation from a vector space  $V$

to a vector space  $W$ , then for any vectors

## 2 Important properties of LT

① Additivity  $T(u+v) = T(u) + T(v)$

② Homogeneity  $T(cu) = cT(u)$



$T: V \rightarrow W \rightarrow$  Linear Transformation

$u, v \in V$  &  $c$  is a Scalar Value

### Ex: Reflection:

The reflection transformation  $T$  across the  $y$ -axis maps a vector

$$x = \begin{bmatrix} x \\ y \end{bmatrix} \in \mathbb{R}^2 \quad T(x) = \begin{bmatrix} -x \\ y \end{bmatrix}$$

The linear transformation can be expressed as  $\rightarrow [T(x) = Ax]$

$$A = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} x & y \end{bmatrix}_{1 \times 2} \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}_{2 \times 2} = \begin{bmatrix} -x \\ y \end{bmatrix}_{1 \times 2}$$

### ① Checking additivity:

Let  $u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$  and  $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$  be two vectors in  $\mathbb{R}^2$

$$T(u+v) = T(u) + T(v)$$

$$u+v = \begin{bmatrix} u_1 + v_1 \\ u_2 + v_2 \end{bmatrix}$$

$$T(u+v) = A(u+v)$$

$$T(u+v) = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 + v_1 \\ u_2 + v_2 \end{bmatrix} = \begin{bmatrix} -(u_1 + v_1) \\ u_2 + v_2 \end{bmatrix}$$

— (1)

R.H.S

$$T(u) + T(v)$$

$$T(u) = Au = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} -u_1 \\ u_2 \end{bmatrix}$$

$$T(v) = Av = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} -v_1 \\ v_2 \end{bmatrix}$$

$$R.H.S = T(u) + T(v) = \begin{bmatrix} -u_1 - v_1 \\ u_2 + v_2 \end{bmatrix}$$

— (2)

$$L.H.S = R.H.S \quad \checkmark$$

② Checking homogeneity:

$$\text{Let } u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \in \mathbb{R}^2 \text{ & } c \text{ be a scalar}$$

→ Homogeneity requirement:

$$T(cu) = cT(u)$$

$$L.H.S \quad T(cu) = A(cu)$$

$$\therefore = A(c \begin{bmatrix} u_1 \\ u_2 \end{bmatrix})$$

$$= A \begin{bmatrix} cu_1 \\ cu_2 \end{bmatrix}$$

$$= \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} cu_1 \\ cu_2 \end{bmatrix} = \begin{bmatrix} -cu_1 \\ cu_2 \end{bmatrix} = 0$$

R.H.S

$$\begin{aligned}CT(u) &= c(A \cdot u) \\&= c \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = c \begin{bmatrix} -u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} -cu_1 \\ cu_2 \end{bmatrix} \quad \textcircled{2}\end{aligned}$$

$$L.H.S = R.H.S \quad \textcircled{1}$$

It is linear transformation.

Example that does not follow linear transformation

$$T: \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

$$u \quad v$$

$$b = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$T(x) = x + b$$

↓

fixed vector

$$T(x) = x + \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

① check Additivity

$$T(u+v) = T(u) + T(v)$$

$$u = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad v = \begin{bmatrix} 4 \\ -1 \end{bmatrix}$$

$$T(u+v) = T\left(\begin{bmatrix} 2 \\ 3 \end{bmatrix} \begin{bmatrix} 4 \\ -1 \end{bmatrix}\right) = T\left(\begin{bmatrix} 6 \\ 2 \end{bmatrix}\right)$$

$$= \begin{bmatrix} 6 \\ 2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 7 \\ 3 \end{bmatrix} \Rightarrow \text{R.H.S}$$

L.H.S

$$T(u) + T(v)$$

$$T(u) = T\left(\begin{bmatrix} 2 \\ 3 \end{bmatrix}\right) = \begin{bmatrix} 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 4 \end{bmatrix}$$

$$T(v) = \cancel{\left(\begin{bmatrix} 1 \\ -1 \end{bmatrix}\right)} \quad \cancel{\left[\begin{bmatrix} 1 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix}\right]} = \cancel{\begin{bmatrix} 4 \\ 0 \end{bmatrix}} = T\left(\begin{bmatrix} 4 \\ -1 \end{bmatrix}\right) = \begin{bmatrix} 4 \\ -1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 0 \end{bmatrix}$$

$$T(u) + T(v) = \begin{bmatrix} 8 \\ 4 \end{bmatrix} \Rightarrow \text{L.H.S}$$

$$\text{R.H.S} \neq \text{L.H.S}$$

Checking Homogeneity:

$$T(cu) = cT(u)$$

$$u = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \quad c = 2$$

$$T(cu) = T\left(\begin{bmatrix} 4 \\ 6 \end{bmatrix}\right) = \begin{bmatrix} 4 \\ 6 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 7 \end{bmatrix} \Rightarrow \text{L.H.S}$$

$$cT(u) = 2\left(\begin{bmatrix} 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) = 2 \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 6 \\ 8 \end{bmatrix} \Rightarrow \text{R.H.S}$$

$$T(x) = x + 1 \text{ is not a linear transformation}$$

→ fails both additivity & homogeneity properties.

## In/by Linear Transformation:-

Linear transformations are fundamental in data science for several reasons. They provide a mathematical framework for manipulating & analyzing data, which is crucial for various data processing tasks, model building & interpretation.

→ Here are the some key reasons why linear transformations are important in data science.

### ① Dimensionality Reduction:-

#### Principal Component Analysis (PCA):

PCA is a widely used technique for reducing the dimensionality of datasets while retaining as much variance as possible. It involves finding a set of orthogonal axes (principal components) & projecting the data onto these axes. The transformation of data points in the original space to the new space defined by the principal components is a linear transformation. This helps in

While we are reducing dimension, we need to make sure that we ~~capture~~ a lot of variance [information that should not be lost & it talks about how data is spreaded]

Reducing computational cost

Mitigating the curse of dimensions

Visualizing high-dimensional data

## ② Feature engineering

linear transformations can be used to create new features from existing ones.

For example, interactions between features can be captured through linear combinations,

which can be used in ML models to improve predictive performance.

Techniques like linear algebra, stepwise regression, ridge regression & linear discriminant analysis (LDA)

all rely on linear transformations to find meaningful feature representations.

The fundamental step in feature engineering is Normalization & standardization.

## ③ Data preprocessing:

Normalization & Standardization:

Linear transformations are used to scale data, making it suitable for ML models. Standardization transformation data to have a mean of zero & a standard deviation of one, while normalization scale data to specific range (e.g.,  $[0,1]$ ). These transformations are essential for ensuring that all features contribute equally to the model, especially in algorithms like gradient descent.

## ④ Neural Networks

In Neural Network, especially deep learning model, the layers consist of linear transformation followed by non-linear activation functions. The weights in a neural network can be seen as a series of linear transformations that map input data to intermediate layers & eventually, to the output layer.

Crucial for learn complex patterns in data

## ⑤ Image & Signal Processing:

In image & signal processing, linear transformations are used extensively.

Ex: Convolutional filters: in image processing can be seen as linear transformation applied to local regions of a image. Fourier transforms, which decompose signals into sinusoidal components are linear transformations that convert time domain signals in to frequency-domain representation.

## ⑥ Understanding & Interpretation:-

Linear transformations simplify complex relationships between variables into linear relationships, which are easier to understand & interpret.

For example: linear regression provides a clear model of how each feature affects the target variable through linear coefficients, making it easier to explain to stakeholders.

## ⑦ Optimization & Solving Systems of Equations:

Linear transformations are used to solve systems of linear equations, which is a common problem in data analysis & optimization.

Techniques like matrix inversion & the use of pseudo-inverse are essential for finding solutions in linear regression & other linear models.

## 8) Theoretical Foundations:-

Many advanced machine learning algorithms & statistical techniques have linear algebra and linear transformations at their core, understanding these fundamentals is crucial for grasping more complex topics like Support Vector Machines.

### Linear Transformation Visualization

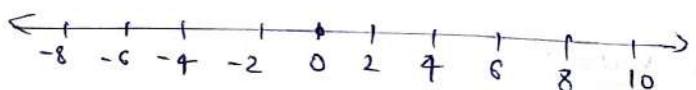
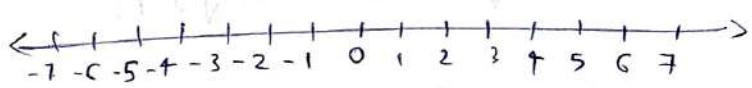
$T: \mathbb{R} \rightarrow \mathbb{R}$  → Satisfy to properties of LT

1D

$$T(x) = 2x$$

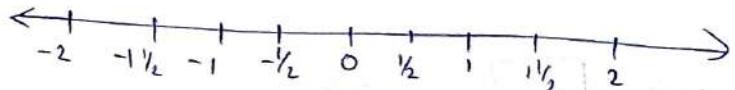
$$\downarrow \rightarrow f(x) = 2x$$

we can suppose it as functions  
for understanding



$$T(x) = \frac{1}{2}x$$

$$T(x) = -3x$$



one of the very important property of LT is origin must remain fixed

Property

- ① origin must be fixed
- ② All lines must remain lines

The better way to understand is by

2D matrix

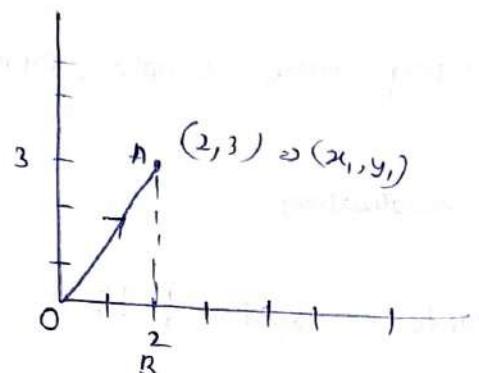
## Magnitude and unit vector:

↓  
Vector length

↳ Vector has a length of 1

$$\vec{x} \in \mathbb{R}^n \quad \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\vec{A} = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \in \mathbb{R}^2$$



$$OA = \|\vec{A}\| = \sqrt{OB^2 + OA^2} = \sqrt{2^2 + 3^2} = \sqrt{4+9} = \sqrt{13}$$

$$\|\vec{x}\| = \sqrt{x_1^2 + x_2^2 + x_3^2 + \dots + x_n^2}$$

magnitude / length of vector

## Unit vector:

$$\|\vec{u}\| = 1 \Rightarrow \hat{u} \rightarrow \text{cap}$$

$$\vec{v} = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{bmatrix} \rightarrow \vec{u} \Rightarrow \|\vec{u}\| = 1 \quad \text{Ex: } \vec{v} = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} \in \mathbb{R}^3$$

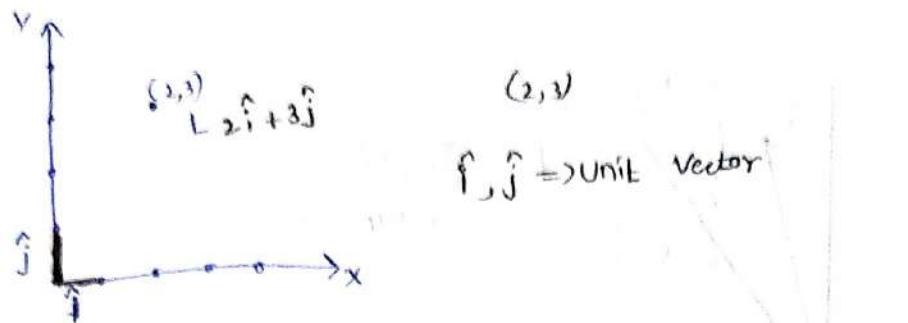
$$\boxed{\vec{u} = \frac{1}{\|\vec{v}\|} \cdot \vec{v}}$$

scalar multiplication

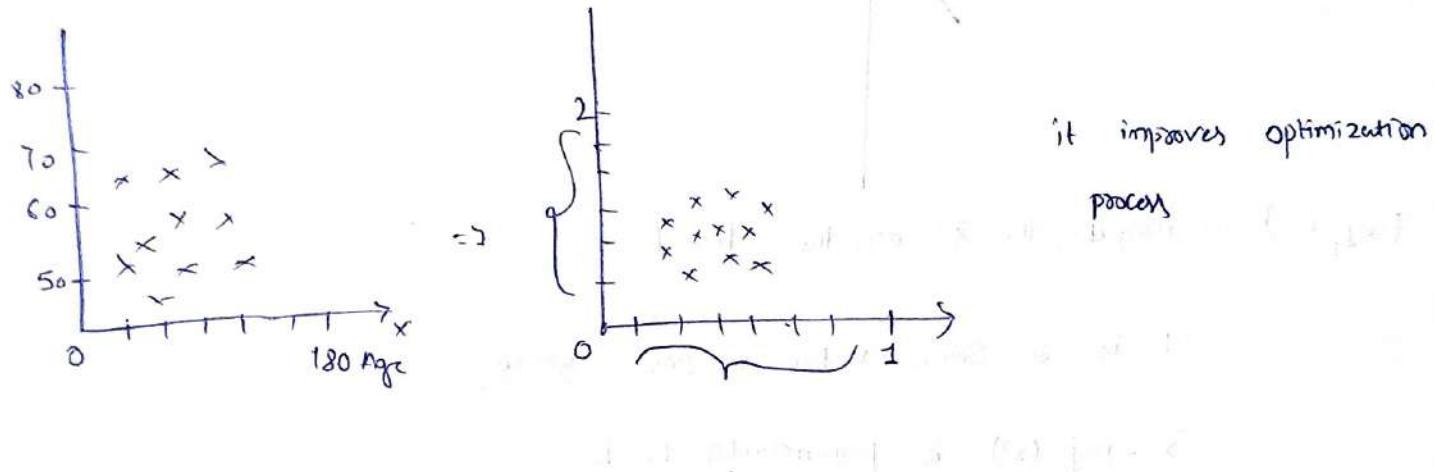
$$\|\vec{v}\| = \sqrt{1^2 + 2^2 + 0^2} = \sqrt{1+4+0} = \sqrt{5}$$

$$\vec{u} = \frac{1}{\sqrt{5}} \cdot \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{5}} \\ \frac{2}{\sqrt{5}} \\ 0 \end{bmatrix}$$

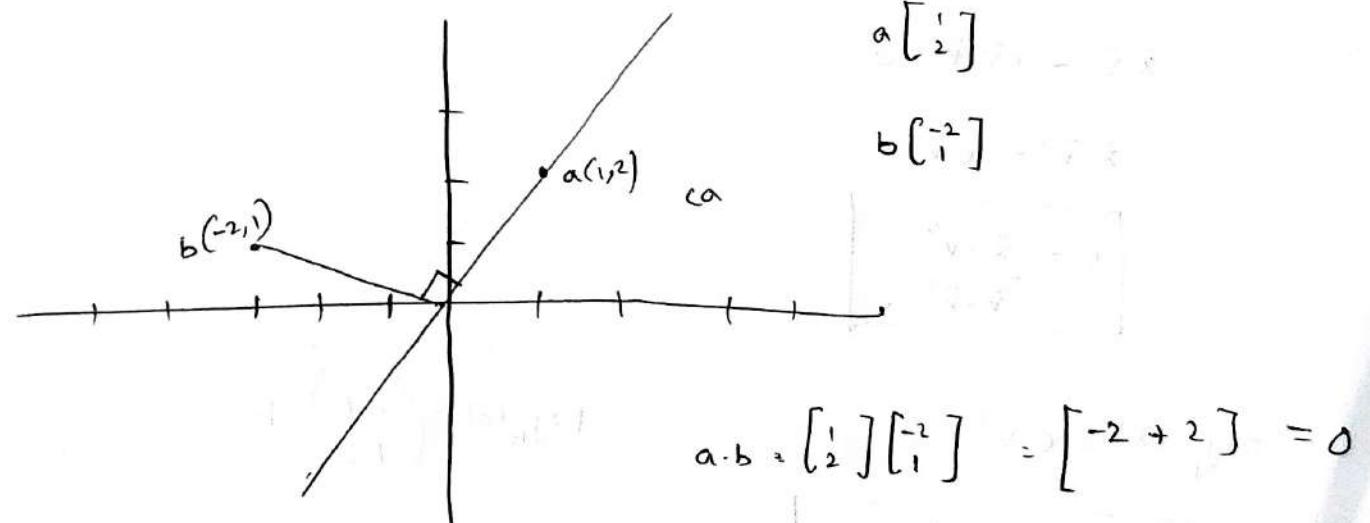
$$\|\vec{u}\| = \sqrt{\left(\frac{1}{\sqrt{5}}\right)^2 + \left(\frac{2}{\sqrt{5}}\right)^2 + 0^2} = \sqrt{\frac{1}{5} + \frac{4}{5} + 0} = \sqrt{\frac{5}{5}} = \sqrt{1} = 1 = \hat{u}$$



Normalization:



Introduction to projection:



$$a \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$b \begin{bmatrix} -2 \\ 1 \end{bmatrix}$$

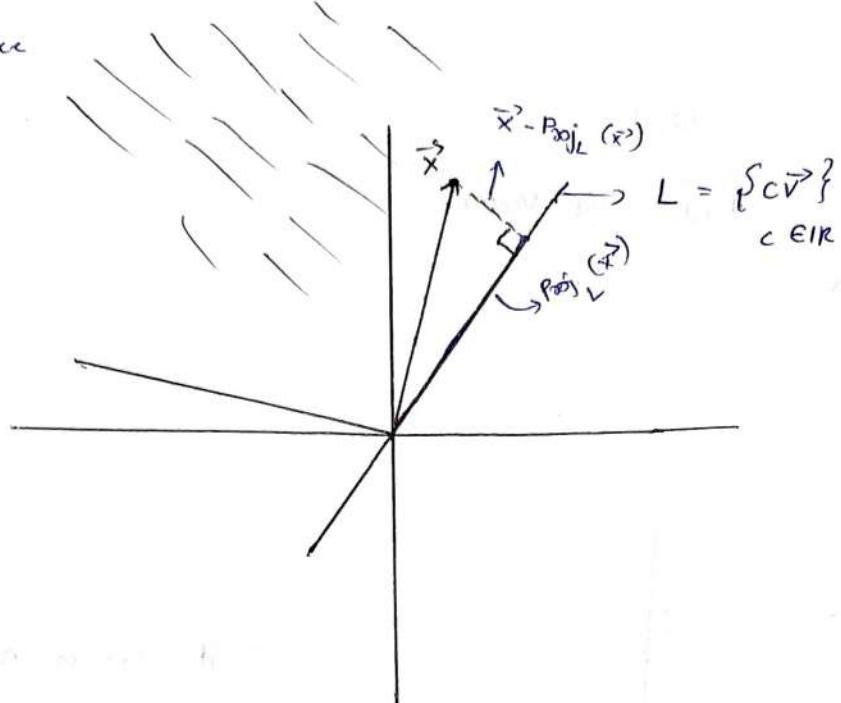
$$a \cdot b = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \begin{bmatrix} -2 \\ 1 \end{bmatrix} = \begin{bmatrix} -2 + 2 \end{bmatrix} = 0$$

$$a \cdot b = a_1 b_1 + a_2 b_2$$

$$a = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

$a \cdot b = 0 \Rightarrow$  Dot product is always 0.

Source



$\text{Proj}_L(\vec{x}) \Rightarrow$  project the  $\vec{x}$  on the line L

If is a some vector is line where.

$\vec{x} - \text{Proj}_L(\vec{x})$  is perpendicular to L

$$(\vec{x} - c\vec{v}) \cdot \vec{v} = 0$$

$$\vec{x} \cdot \vec{v} - c\vec{v} \cdot \vec{v} = 0$$

$$\vec{x} \cdot \vec{v} = c\vec{v} \cdot \vec{v}$$

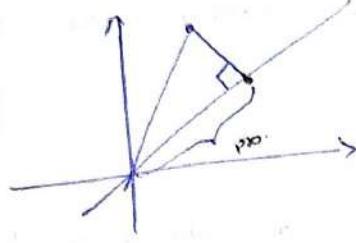
$$c = \frac{\vec{x} \cdot \vec{v}}{\vec{v} \cdot \vec{v}}$$

$$\text{Proj}_L(\vec{x}) = c\vec{v}$$

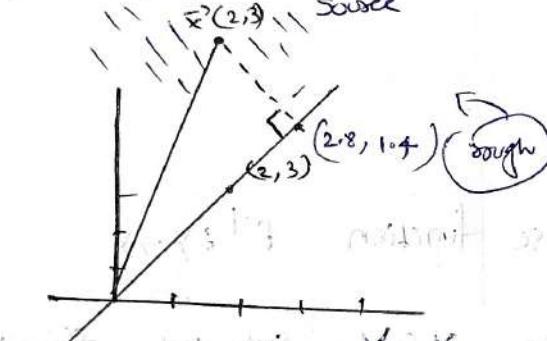
$$\text{Proj}_b(a) = \left( \frac{a \cdot b}{b \cdot b} \right) \cdot b$$

$$\boxed{\text{Proj}_L(\vec{x}) = \left( \frac{\vec{x} \cdot \vec{v}}{\vec{v} \cdot \vec{v}} \right) \cdot \vec{v}}$$

Problem:  $\vec{x} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$  &  $\vec{v} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$



$$L = \left\{ c \begin{bmatrix} 2 \\ 1 \end{bmatrix} \mid c \in \mathbb{R} \right\} \quad \vec{x} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$



$$\text{Proj}_L(\vec{x}) = \left( \frac{\begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 1 \end{bmatrix}}{\begin{bmatrix} 2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} 2 \\ 1 \end{bmatrix}} \right) \cdot \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$= \frac{7}{5} \cdot \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} \frac{14}{5} \\ \frac{7}{5} \end{bmatrix} = \begin{bmatrix} 2.8 \\ 1.4 \end{bmatrix}$$

## Inverse of a function:

A inverse of a function is a function that "reverses" the effect of the original function.

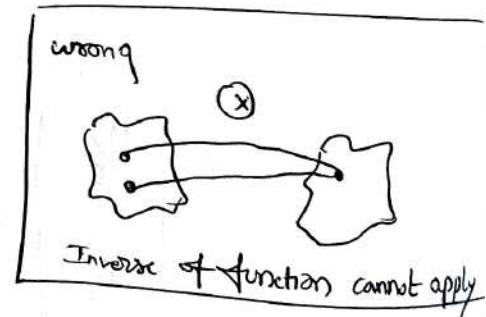
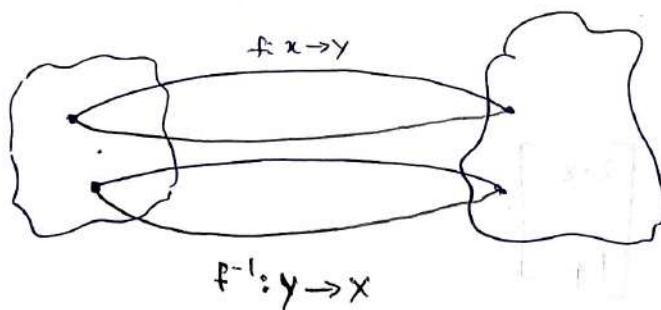
If you have a function  $f$  that maps an element  $x$  from set  $X$  to an element  $y$  in a set  $Y$ , then inverse function  $f^{-1}$  map  $y$  back to  $x$ .

### Defination:

Given function  $f: X \rightarrow Y$ , then inverse function  $f^{-1}: Y \rightarrow X$   
for every  $y \in Y$ , there is a unique  $x \in X$  such that  $f(x) = y$

The inverse function  $f^{-1}$  satisfies the following condition

- 1) For all  $x \in X$ :  $f(f^{-1}(y)) = y$
- 2) for all  $y \in Y$ :  $f^{-1}(f(x)) = x$



These conditions imply that applying the function & then its inverse will return the original value

## Identity function

$$I_x : x \rightarrow x \quad \Rightarrow \quad I_x(a) = a$$

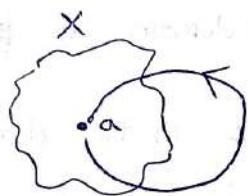
$a \in x$

$$x = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad v = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

$$Iv = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$

for a set  $x$ , the identity function  $I_x$  is defined as:

$$I_x(a) = a \quad \text{for all } a \in x$$



### Properties of identity function

① Preservation: Does not alter any element. If  $x$  is the domain, then the image of  $x$  under the identity  $f_n$  is  $x$ .

② Linearity: Identify ~~for~~  $f_n$  is a linear transformation

$$\star I(u+v) = I(u) + I(v)$$

$$\star I(cu) = c I(u) = cu$$

③ Identity matrix:  $n \times n \Rightarrow$  All the diagonal elements will be 1 and others 0's

$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

$3 \times 3$

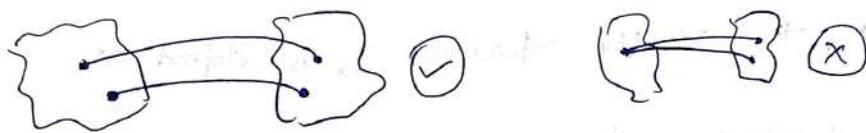
$$I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

4 Inverse: The identity function  $f_n$  is its own inverse.

Existence & uniqueness:

A function  $f$  has an inverse if & only if it is bijective.

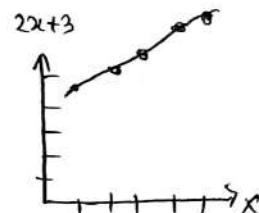
ii) Injective (one to one): Different elements in the domain map to different elements in the codomain.



iii) Surjective (on to): Every element in the codomain is the image of at least one element in the domain.

Eg: Linear function

$$f(x) = 2x + 3$$



Find the inverse:

$$y = 2x + 3 \text{ for } x$$

$$y = 2x + 3$$

$$y - 3 = 2x$$

$$x = \frac{y-3}{2}$$

The inverse function

$$f^{-1}(y) = \frac{y-3}{2}$$

Verification:

$$\textcircled{1} \quad f(f^{-1}(y)) = f\left(\frac{y-3}{2}\right) = 2\left(\frac{y-3}{2}\right) + 3 = y$$

$$\textcircled{2} \quad f^{-1}(f(x)) = f^{-1}(2x+3) = \frac{(2x+3)-3}{2} = x$$

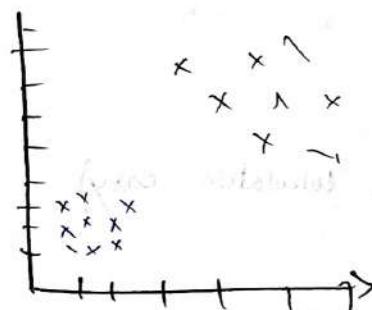
## (Application of Inverse function In Data Science:-)

### ① Normalization And Standardization :-

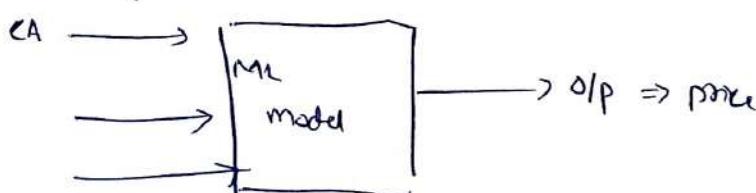
Ex: House Price Dataset

Independent features				Dependent features
Carpet Area	No. of Room	Area (sqft)	price	
[1500]	[4]	[1800]	[50 lakhs]	
1800	3	2100	60 lakhs	

Some Scaling is done to the vectors in a Same Scale



instead of having huge magnitude we try to reduce it



{0 - 1}  
 mathematical operations }  
 operations ⇒ quickly ⇒ optimization

## Standardization

We take a Data feature [vector] which will transform to the new feature of having  $\mu=0$  &  $\sigma=1$

No of Rooms

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

$\in \mathbb{R}$

Standardization

Transformation

$$z = \frac{x_i - \mu}{\sigma}$$

No. of Rooms

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

$\in \mathbb{R}$

$f(x)$

$f^{-1}(x)$

No. of rooms  $^T$

$$\begin{bmatrix} -1.5/2 \\ -0.5/2 \\ 0.5/2 \\ 1.5/2 \end{bmatrix}$$

$$z_1 = \frac{1 - \mu}{\sigma}$$

$$\mu = \frac{1+2+3+4}{4} = \frac{10}{4} = 2.5$$

$$\sigma = 2 \quad (\text{not correct answer, just to make calculation easy})$$

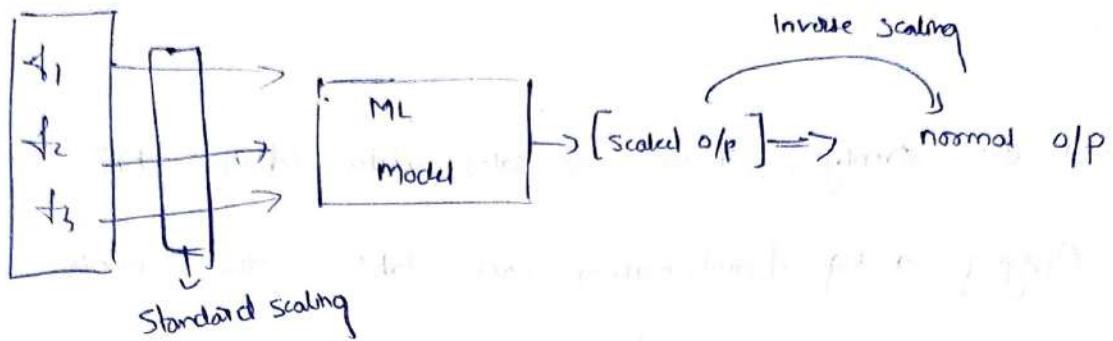
Original Transformation

$$z = \frac{x - \mu}{\sigma}$$

$$x = \frac{-1.5}{z} \cdot z + 2.5 = 1 =$$

Inverse Transformation

$$x = z\sigma + \mu$$



Uses :-

After training a ML model on standardized data, the predictions are often descaled back to the original scale to interpret the result in a meaningful way. For instance, if house prices were standardized → the inverse transformation would convert the standardized predictions back to the original price scale.

\* Normalization :-

Feature Scaling with Min Max Normalization

$$\rightarrow \text{original Transformation} : z = \frac{x - \min(x)}{\max(x) - \min(x)} \rightarrow T: x \rightarrow y$$

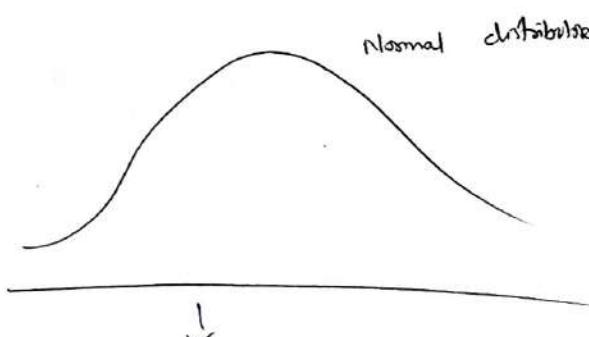
$$\rightarrow \text{Inverse Transformation} : x = \frac{z(\max(x) - \min(x))}{\max(x) - \min(x)} \rightarrow T^{-1}: y \rightarrow x$$

\* Distribution of Data

\* logarithmic distribution



$\Rightarrow$   
 $\Rightarrow$



$$\text{OT} : y = \log(x)$$

$$\text{IT} : x = e^y$$

Because ML model need ND

### Uses case:

In financial data analysis, income or sales data often exhibit skewness. Applying a log transformation can stabilize the variance & make patterns more visible. After model prediction, the inverse log transformation is applied to interpret the result on the original scale.

### ⑤ Data Encryption & Decryption:

Encryption function:  $E(p) = c$  (where  $p$  is plain text &  $c$  is a cipher text)

Decryption function:  $D(c) = p$

### Use case:

Sensitive data like personal information, financial records, & medical data are encrypted before storage or transmission.

Decryption is applied to retrieve the original information.

→ How to find inverse of a matrix

## ① Determinant

## ② How to Inverse

Eg:  $2 \times 2$  matrix

Find  $A^{-1}$  & verify using a transformation

$$A = \begin{bmatrix} 4 & 7 \\ 2 & 6 \end{bmatrix}$$

$$T: \mathbb{R} \rightarrow \mathbb{R}$$



$$Ax = y$$



$$A^{-1}$$

Find the inverse of A

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

$$A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

$$\text{using example } A = \begin{bmatrix} 4 & 7 \\ 2 & 6 \end{bmatrix} \Rightarrow A^{-1} = \begin{bmatrix} 6 & -7 \\ -2 & 4 \end{bmatrix}$$

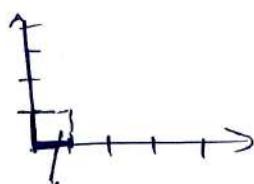
The determinant is a scalar value that can be computed from a

Square matrix. It provides important information about the matrix, such as

whether the matrix is invertible (ie has an inverse), & it also has

geometric interpretations, such as describing the scaling factor of

linear transformations represented by the matrix.



$\Delta \text{area} = \text{Determinant}$

Determinant :  $\begin{bmatrix} a & b \\ c & d \end{bmatrix} \Rightarrow \boxed{ad - bc} \Rightarrow$  scalar  $\Rightarrow$  determinant

$$A = \begin{bmatrix} 4 & 7 \\ 2 & 6 \end{bmatrix}$$

$$\det(A) = 24 - 14 = 10$$

Step 1 :  $\det(A) \neq 0$

Non zero  $\Rightarrow$  inverse of matrix exists  $\Leftrightarrow \det(A) \neq 0$

Matrix A is invertible

Step 2 find the inverse

$$A^{-1} = \frac{1}{\det(A)} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix} \Rightarrow \text{adjoint matrix}$$

$$= \frac{1}{10} \begin{bmatrix} 6 & -7 \\ -2 & 4 \end{bmatrix} = \begin{bmatrix} 0.6 & -0.7 \\ -0.2 & 0.4 \end{bmatrix}$$

Step 3 Verification using a vector

$x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow x$  using A then use  $A^{-1}$  to recover the original vector

Transformation using A

$$y = Ax = \begin{bmatrix} 4 & 7 \\ 2 & 6 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 4+7 \\ 2+6 \end{bmatrix} = \begin{bmatrix} 11 \\ 8 \end{bmatrix}$$

Recovering  $x$  using  $A^{-1}$

$$x = A^{-1}y = \begin{bmatrix} 0.6 & -0.7 \\ -0.2 & 0.9 \end{bmatrix} \begin{bmatrix} 11 \\ 8 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

Thus  $A^{-1}$  successfully recover the original vector  $x$

### Eigen vectors And Eigen values

Eigenvalues & Eigen vectors are fundamental concepts in linear algebra that have numerous applications in various fields such as - Physics, computer science, and data science. They provide insights into the properties of linear transformations represented by matrices

#### Definition

Eigen value ( $\lambda$ ) :- A scalar that indicates how much an eigen vector is stretched or compressed during linear transformation

Eigen vector : A non zero vector that only changes in scale (not direction) when a linear transformation is applied.

\* To calculate eigen vector given by  $AV = \lambda V$

For square matrix  $A$ , an eigen vector and its corresponding eigen value  $\lambda$  satisfy the above equation

$$A = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix}$$

① Find eigen values

$$\det(A - \lambda I) = 0$$

$$A - \lambda I = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} = \begin{bmatrix} 4-\lambda & 1 \\ 2 & 3-\lambda \end{bmatrix}$$

$$\det \begin{bmatrix} 4-\lambda & 1 \\ 2 & 3-\lambda \end{bmatrix} = 0 \Rightarrow (4-\lambda)(3-\lambda) - 2 = 0$$

$$\Rightarrow 12 - 4\lambda - 3\lambda + \lambda^2 - 2 = 0$$

$$\Rightarrow \lambda^2 - 7\lambda + 10 = 0$$

② Solve the equation

$$\lambda^2 - 7\lambda + 10 = 0$$

$$\boxed{\lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}}$$

$$a=1, b=-7 \quad \& \quad c=10$$

$$\lambda = \frac{7 \pm \sqrt{49 - 40}}{2 \times 1} = \frac{7 \pm \sqrt{9}}{2} = \frac{7 \pm 3}{2}$$

$$\lambda_1 = \frac{10}{2} = 5 \quad (\text{or}) \quad \lambda_2 = \frac{4}{2} = 2$$

Eigen values of A       $\lambda_1 = 5$      $\lambda_2 = 2$

③ Find the eigen vectors:

$$(A - \lambda I) v = 0$$

for  $\lambda_1 = 5$

$$A - 5I = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix} - \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} = \begin{bmatrix} 4-5 & 1-0 \\ 2-5 & 3-5 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 2 & -2 \end{bmatrix}$$

$$(A - \lambda I) v = 0$$

$$\begin{bmatrix} -1 & 1 \\ 2 & -2 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = 0$$

$2 \times 2 \quad 2 \times 1$

$$-x + y = 0$$

$$\begin{cases} -x + y = 0 \\ 2x - 2y = 0 \end{cases}$$

$\boxed{\begin{array}{l} y=x \\ x=y \end{array}} \rightarrow$  Eigen vector corresponding to Eigen value  $\lambda_1 = 5$

for

$$\lambda_2 = 2$$

$$A - 2I = \begin{bmatrix} 4-2 & 1 \\ 2 & 3-2 \end{bmatrix} = \begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix}$$

$$(A - \lambda I) v = 0$$

$$\begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0$$

$$2x + y = 0$$

$$\boxed{y = -2x} \rightarrow \text{Eigen vector}$$

corresponding to  $\lambda_2 = 5$

$$A = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix} \quad \lambda_1 = 5 \quad \lambda_2 = 2$$

for  $\lambda_1 = 5$        $v_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

for  $\lambda_2 = 2$        $v_2 = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$

$$A \vec{v} = \lambda \vec{v}$$

↓  
how much eigen value that the vector  $\vec{v}$  is stretched

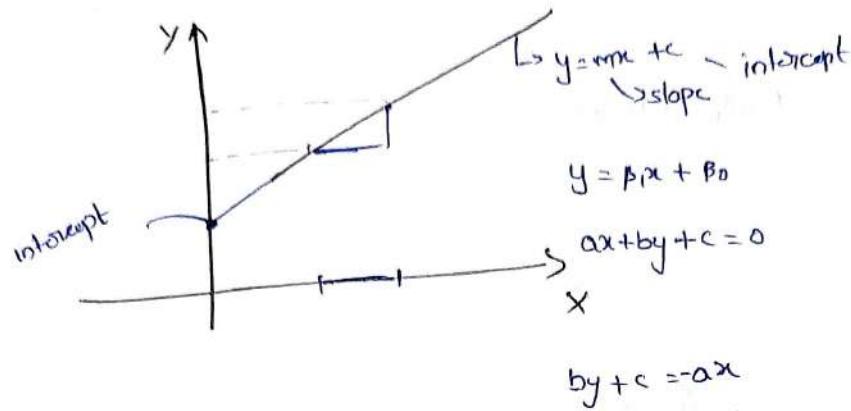
These eigenvectors & eigenvalues describe how the matrix A scales and rotates vectors in its transformation. Eigen values indicates the factor by which the eigenvectors are stretched or compressed, and eigenvectors provide the directions in which this stretching or compression occurs.

### Applications:

Principal component analysis  $\rightarrow$  Dimensionality reduction

# Equation of line, 3d plane and hyperplane (in Dimension)

(2D)



$$w_1x_1 + w_2x_2 + b = 0$$

$$\boxed{w^T x + b = 0}$$

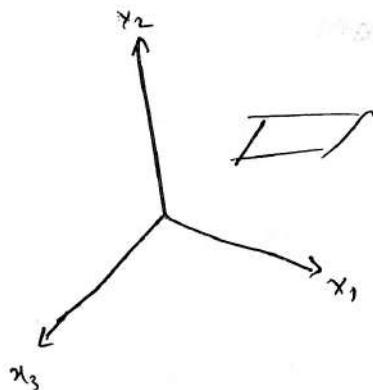
eq of straight line

$$y = -\frac{a}{b}x - \frac{c}{b}$$

$$\downarrow \quad \downarrow$$

$$y \quad c$$

(3D)



$$w_1x_1 + w_2x_2 + w_3x_3 + b = 0$$

$$\boxed{w^T x + b = 0}$$

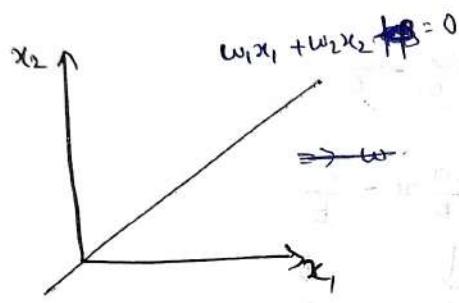
$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \cdot x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

## n-Dimension plane

$$w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b = 0$$

$$w^T x + b = 0$$

If can pass through origin  $\Rightarrow b = 0$



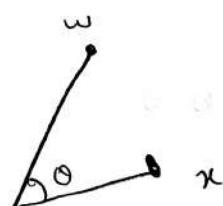
Equation of a straight line passing through origin

$$w^T x = 0$$

~~#~~ 3.142

Equation of a plane =  $T_h \Rightarrow w^T x = 0$

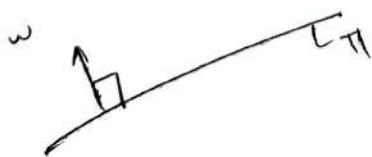
$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \\ \vdots \\ w_n \end{bmatrix} \quad \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$



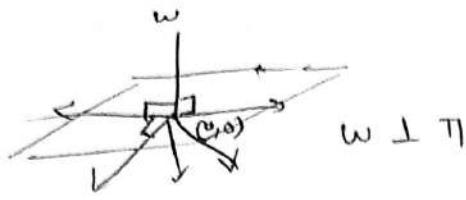
$$w \cdot x = w^T x = \|w\| \|x\| \cos \theta = 0$$

$$\theta = 90^\circ$$

$$\cos \theta = 0$$



$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$



Intercept = 0

## Statistics :-

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

Decision making

Data : facts or pieces of information  $\Rightarrow$  measured, collected, Analyzed

Ex:- weight of student in the class \* IQ of the Student of class

{ 60kg, 50kg, ... }

{ 100, 90, 99, ... }

<u>city</u>	Area	No. of rooms	popn
Hyd	1000	2	45 lakhs
New York	1250	2.5	50 lakhs

→ Data Analyst → Report → Visualization → Meaning Decision  
L project

Analysis can be done to the data with the help of statistics

## Application: Typ. statistics

- ① Data exploration & summarize
- ② Model build & validation
- ③ Statistical Analysis  $\Rightarrow$  Different analysis we do on sample data to make conclusion about population
- ④ Hypothesis Testing
- ⑤ Optimization & efficiency
- ⑥ Reporting

## Types of statistics

- ① Descriptive Statistics
- ② Inferential Statistics

### ① Descriptive Statistics:-

Descriptive statistics involves methods for summarizing & organizing data to make it understandable. This type of statistics helps to describe the basic features of the data in a study.

## II Measure of central tendency

[mean, median, mode]

## II Measure of Dispersion

[variance, standard deviation]

## III Data Distribution

i) Histograms

[to understand probability distribution function & probability density function]

ii) Box plot

iii) pie chart

iv) PDF, PMF

## IV Summary statistics

i) five number summary

$Q_1, Q_2, Q_3, \text{ Maximum}$

## 2) Inferential Statistics:

Inferential statistics involves methods for making predictions or inferences about a population based on a sample of data. It allows for hypothesis testing, estimation & drawing conclusion

① Hypothesis Testing

② P value

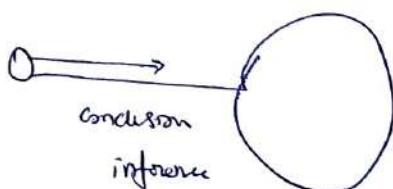
③ Confidence interval

④ statistical analysis Test      ① Z test

② T test

③ ANOVA  $\rightarrow$  F test

④ Chi square



## Type of Statistics

### Key Concepts

Example:

## Descriptive Statistics

Measures of central tendency  
(mean, median, mode), Measures  
of Dispersion (range, variance,  
standard deviation), Data  
Distribution (Histograms,  
Box plots), Summary  
statistics (five-number summary)

Mean source of the students,

Range of temperature,

Histogram of ages

## Inferential Statistics

Hypothesis Testing (Null &  
Alternative hypothesis, P-value),  
Confidence intervals, Regression  
Analysis (Simple & multiple  
Linear regression), ANOVA,  
Chi-Square Test

P-value in test ~~source~~ source

Comparison, 95% confidence

interval for average height,

Predicting house price,

Comparing test scores of

different schools, Association

between gender & product

Preference.

Eg: Let say there are 20 statistics class in your college, & you have

Collected the height of students in the class

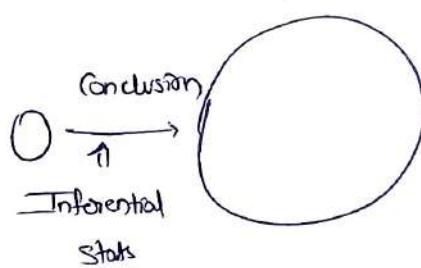
→ heights are recorded [175cm, 180 cm, 140cm, 135cm, 160cm, 120cm]

Descriptive question:

"What is the average height of the entire classroom" → Measure of Central Tendency

Inferrential Question

"Are the heights of the sample students in classroom similar to what you expect in the entire college?"



③ Population

Sampling

Data:

Population: A population is the entire set of individuals or objects of interest in a particular study. It includes all members of a defined group that we are studying or collecting information on.

## Characteristics:-

- ① Complete Set: Contains all the observations of interest  
② Parameter: A numerical value summarizing the entire population

Ex: ① Population mean ( $\mu$ )

② Population variance ( $\sigma^2$ )

### Example:

① population in a school study

\* All students enrolled in a school

use case

Determine the avg height of student,

Population mean

② population in market Research

\* All consumers in a city

use cases

To understand the purchasing behaviour of all consumers

③ population in a medical study

\* All the patients with a specific disease

use cases

To study the effect of a drug.

## Sample data:-

- \* A Sample is a Subsets of the population that is used to represent the entire group. Sampling involves Selecting a group of individuals or observations from the population to draw conclusions about the whole population.

## Characteristics:-

- ① Subset : Represents a portion of the population.
- ② Statistic : A numerical value Summarizing the Sample data  
[Sample mean , Sample variance ]
- ③ Random Sampling : Sampling should be randomly Selected to avoid bias

## Example:-

- ① Sample in a School study
  - (i) A group of 50 students from School
  - (ii) Usecase :- Estimate the average height of students in a school
- ② Sample in Market research ;
  - \* 500 consumer from the city
  - \* Behaviour → population

$\Rightarrow$  ③ Sample in a Medical Study

- \* 200 patients
- \* Tests the effectiveness of the drug

### Types of Sampling Techniques:-

① Probability Sampling

② Non probability Sampling

① Probability Sampling :- All individual have chance to ~~not~~ being selected

a) Sample Random Sampling:

The Sampling data is collected randomly

$\Rightarrow$  every no of the population has an equal chance of being selected

Ex:-

- Selecting people randomly

- Draw names random from a class of student

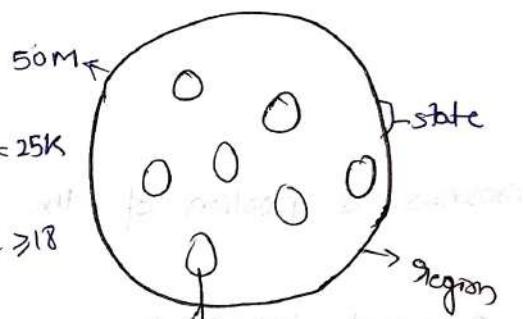
b) Systematic Sampling

Select every  $n^{th}$  member of the population after a random starting point.

Ex:- Airport  $\rightarrow$  credit card  $\rightarrow$  5<sup>th</sup> person, 10<sup>th</sup>, 15<sup>th</sup>

• Feedback Survey  $\rightarrow$  Selected every 11<sup>th</sup> number  $\rightarrow$  Feedback Survey

Ex:- Exit poll A, B, C, D



### ③ Stratified Sampling:

Divide the population into strata (group) based on specific characteristics & then randomly sampling from each strata.

Eg:- Divide employees by department & then randomly select a proportional number from each department to form a survey sample.

$$\text{Eg:- Age} \rightarrow \begin{array}{c} <12 \\ \uparrow \\ 12-18 \\ \uparrow \\ >18 \end{array} \quad \left\{ \text{politics} \right\}$$

### ④ Cluster Sampling:

Divide the population in to clusters, randomly selecting clusters, then Sampling all the members from the selected clusters.

Eg:- Randomly Selecting several schools from a district & Surveying all faculty within these schools.

### ⑤ Multi-stage Sampling:

Combining several Sampling method usually involved Selecting clusters, then randomly Sampling within those clusters.

Eg:- Randomly Selecting cities, each selecting city randomly selecting households to survey.

## ② Non Probability Sampling:

Not all the members have equal chance of being Selected

### (a) Convenience:

Selecting individual who are easy to reach

Eg:- Surveying students in a classroom because they are readily available

\* Consecutive Sampling: Similar to Convenience Sampling with slight variation

The researcher picks a single person or group of people for sample

### (b) Judgmental (purposive) Sampling:

Select individual based on the Research's judgement  $\Rightarrow$  useful or representation

Eg:- choose experts in a field to participate. ~~Select~~ focuses on specific characteristic / condition

### ③ Snowball Sampling:

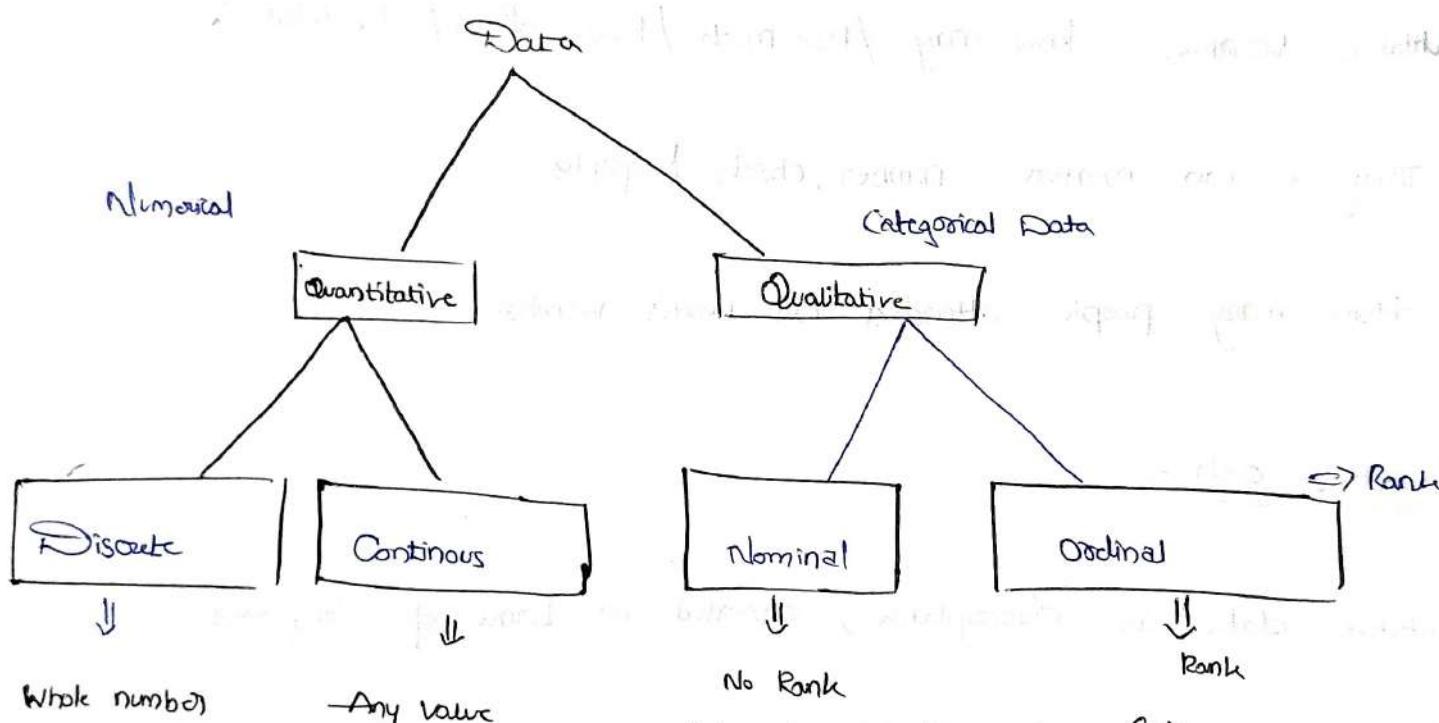
Existing ~~st~~ participants recruit new participants  
Referrals - used for hidden or hard to reach population creating a chain of

Eg:- Survey members of a rare disease.

d) Data Sampling: Sample represents certain characteristics in % like Age, group, gender, caste

e.g.: choosing 60% woman & 40% man

### Types of data:



Eg:

i) No of childrens  
in a family

ii) Marks obtained  
by students  
of a class 10

Eg:

Weight  
Height  
Temperature

Space  
=> It can be  
negative / positive  
decimal

Eg:- Gender

M,F

Blood group

Pincush

Eg:- Customer feedback

Good bad better  
1 3 2

Qualitative data:

Quantitative

It refers to any information that can be quantified - that is, numbers.

If it can be counted or measured & given a numerical value, it's quantitative in nature.

Specific to objective measure of numerical facts

Quantitative variable: how many / how much / how often / The what

Thing u can measure number, charts & graphs

Ex How many people attended last week's webinar

Qualitative data:

Qualitative data is descriptive, expressed in terms of languages

rather than numerical values

Subjective explanatory measures of qualities & characteristics

why or how

How do customers feel about their customer service experience,

## \* Scales of Measurement of Data :-

The Scales of measurement describe the nature of information within the value assigned to variables.

### 4 primary scales of measurement

① Nominal Scale

② Ordinal Scale

③ Interval

④ Ratio

#### 1 Nominal Scale :

This scale classifies data into distinct categories that do not have an intrinsic order.

Qualitative data / Categorical data

#### Characteristics :-

- i) Data is categorized based on label, names or qualities
- ii) Categories are **mutually exclusive**
- iii) No logical order among categories [No rank]

any one will be true

- Ex:-
- Gender      • Colors
  - M              Red → 5 50%
  - F              Blue → 4 40%
  - Pink → 1 10%

### ② ordinal scale :-

This scale classifies the data into categories that can be ranked or ordered characteristics.

- (i) Data is categorized & Ranked in a specific order
- (ii) The interval between rank are not necessarily equal

Example :-

Education level	Ranks	Customer feedback	Rank
High School	1	Satisfied	1
Bachelor	2	Very Satisfied	2
Masters	3	Not Satisfied	0
Doctorate	4		

So economic status :-

Low	2
Middle	1
High	0

### ③ Interval scale :-

The interval scale not only categorizes & orders but also specifies the exact difference between intervals. It lacks a true zero point.

## Characteristics:-

- ① Data is ordered with Consists Interval between values
- ② Allows for meaningful comparison of differences [ratio cannot be measured]
- ③ No true zero point

Ex:- IQ Score:

90, 10, 100,

IQ ≠ 0

Ex:- Student marks in a class

4 Ratio scale:-

\* The order matters

\* Differences are measurable

\* Contains a 0 starting point

0, 90, 60, 30, 45

$$\text{Ratio} = \frac{90}{30} = \frac{3}{1}$$

3:1

Ex:- weight

10, 20, 30, 40

Income :- 10,000 1P

20000

30000

40000

1P

1P

## Assignment:-

i) Length of different rivers in the world? { Ratio }

ii) Favorite food based on Gender? { Nominal }

iii) Marital status? { Nominal }

iv) IQ Measurement { Ratio }

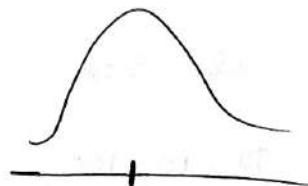
## <sup>Centre</sup> Measure of $\uparrow$ tendency:

Measure of central tendency are statistical metrics that describes the centre point or type value of a data set. They provide a single value that summarizes a set of data by identifying the central position within that dataset.

- 1 Mean or Average
- 2 Median
- 3 Mode

$$\rightarrow \text{Age: } [24, 28, 15, 20, 30]$$

↓  
Centre point



1 Mean: Mean is the sum of all values divided by the number of values.

Population mean ( $\mu$ )		Sample mean ( $\bar{x}$ )
Population size ( $N$ )		Sample size ( $n$ )

$n \leq N$

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

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$

\* Here  $X$  is a random variable

$$X = \{5, 8, 12, 15, 20\}$$

$$N=5$$

$$\mu = \frac{5+8+15+12+20}{5} = \frac{60}{5} = 12$$

### x characteristics:

- Affected by the extreme outliers
- Used for interval & ratio data

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

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

$$x = \{1, 2, 3, 4, 5, 100\}$$

$$\bar{x} = \frac{1+2+3+4+5+100}{6} = 23$$



### Median:

The median is the middle value in a dataset when the values are arranged in ascending or descending order.

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

$$x = \{3, 4, 1, 5, 2, 100\} \Rightarrow [1, 2, 3, 4, 5, 100]$$

$$\text{No. of element} = 5 \quad [\text{odd}]$$

$$\text{No. of element} = 6 \quad [\text{even}]$$

$$\Rightarrow \left(\frac{n+1}{2}\right)^{\text{th}} \text{ term}$$

$$\Rightarrow \frac{\left(\frac{n}{2}\right)^{\text{th}} \text{ term} + \left(\frac{n+1}{2}\right)^{\text{th}} \text{ term}}{2}$$

$$\text{Median} = \left(\frac{5+1}{2}\right)^{\text{th}} \text{ term}$$

$$\Rightarrow \text{median} = \frac{3+4}{2}$$

$$= 3^{\text{th}} \text{ term}$$

$$= \underline{\underline{3.5}}$$

### Characteristics:

- Not affected by extreme outliers

- Used for ordinal, interval & ratio data

### ③ Mode

The mode is the value that appears most frequently in a dataset.

Dataset: 2, 4, 4, 9, 9, 8, 7, 4, 10

$$\text{mode} = 4$$

$$\Rightarrow 6, 6, 3, 17, 4, 9, 3$$

mode = 6, 3 (bimodal) similarly multimodal

\* Characteristic:-

- ① Not affected by extreme values
- ② used for nominal, ordinal, interval & ratio data.

Choosing the Appropriate Measures:

1 Mean:- Best used when data is Symmetrically distributed without outliers provides a mathematical average, which is useful for further statistical calculations

2 Median: Best used when data is skewed or contains outliers provides the middle value, which better represents the centre of a skewed data

3 Mode: Best used for Categorical data to identify the most common category. Also useful for identifying the most frequent value in, ordinal, interval or ratio data

This missing values are replaced with the mean / median to

Real world Application:

Get accurate output

Used in Feature engineering



Nominal data

mode ↑

Gender

Degree

Age	Weight	Salary	Gender	Degree	Handling values
24	70	70k	M	BE	
25	80	70k	F	-	Nominal
27	95	80k	F	-	+ ordinal
24	-	55k	M	PMS	
32	-	45k	F	BE	
-	60	-	M	Master	
-	75	-	M	BS	
40	90	35k	-	BE	

## Measure of Dispersion:-

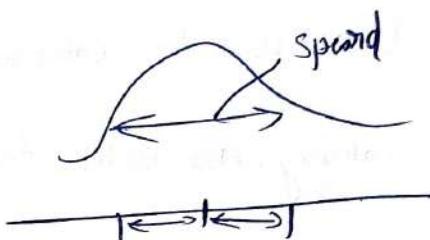
Measures of dispersion describe the spread or variability of a dataset.

They indicate how much the values in a dataset differ from the

Central tendency

Common measure of Dispersion:-

- ① Range
- ② Variance
- ③ Standard Deviation
- ④ Interquartile range (IQR)



1 Range: Range is the difference between the maximum & minimum value in a dataset

$$\text{Range} = \text{Maximum value} - \text{Minimum Value}$$

Eg: Ages  $\{14, 13, 10, 20, 25, 75, 15\}$

$$\text{Range} = 75 - 10 = 65$$

Characteristics:-

- ① Simple to calculate
- ② Sensitive to outliers
- ③ Rough measure of dispersion

$$\text{Weight} = \{34, 40, 45, 30\}, 100\}$$

$$\text{Range} = 45 - 30 = 15$$

$$\text{Range} = 100 - 30 = 70$$

## 2 Variance:-

Variance measures the average squared deviation of each value from the mean. It provides a sense of how much the values in a dataset vary.

Vary

Population variance

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

Sample Variance

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

$x_i$  → Data points

$x_i$  → Data points

$\mu$  → population mean

$\bar{x}$  → Sample mean

$N$  → population size

$n$  → sample size

Sample: Size of a flower petals

{5, 8, 12, 15, 20} → Variance of this distribution

$$N = 5$$

$$\mu = \frac{5+8+12+15+20}{5} = \frac{60}{5} = 12$$

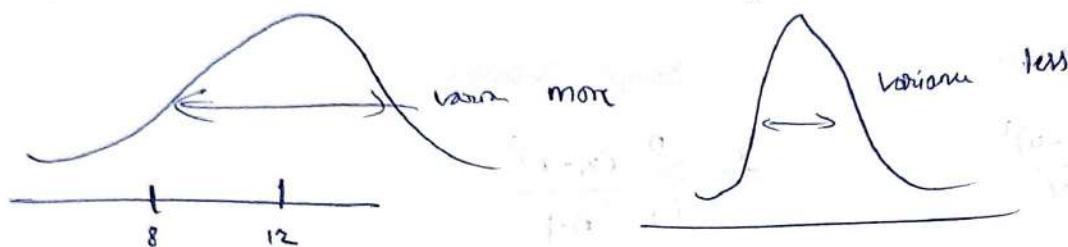
$$\text{Variance} = \frac{(12-5)^2 + (12-8)^2 + (12-12)^2 + (12-15)^2 + (12-20)^2}{5}$$

$$= \frac{49 + 16 + 0 + 9 + 64}{5}$$

$$= 27.6$$

## Characteristics:-

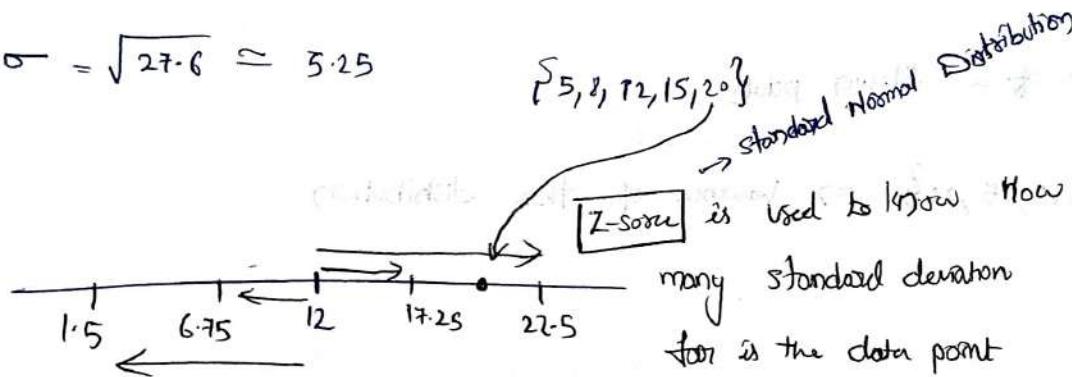
- \* provide a precise measure of variability
- \* units are squared of the original data units
- \* More sensitive to outliers than the range



## (3) Standard Deviation

The Standard Deviation is the square root of the Variance

$$\sigma = \sqrt{27.6} \approx 5.25$$



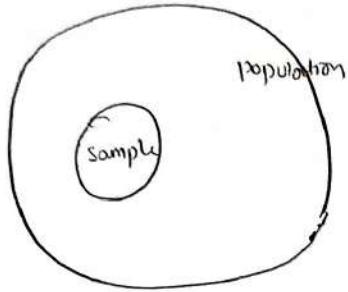
## Characteristics:-

- ② provide a clear measure of spread in the same units as the data

## (2) Sensitive to outliers.

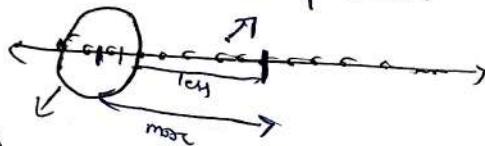
## Sample Variance:

If we use Sample variance formula as  $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$



$$\text{Data: } \{x_1, x_2, x_3, \dots, x_n\}$$

$$\bar{x} \text{ (mean)} + s^2 \text{ (variance)}$$



(Sample) if divide by n then  ~~$\frac{1}{n-1}$~~

Sample mean will be in middle & it will be far away  
from the population mean/variance

From the population mean/variance

\* So the true population variance is underestimated.

Thus we use: 
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$
 & It is called  $s^2$

\* Bessel's correction (using  $n-1$ ) compensates for the estimation

\* "Lose one degree of freedom" when we estimate  $\bar{x}$

you are using information(n) from sample itself, not the true population mean.  
This makes the Sample variance a biased estimator of the population variance if you divide by n

Using  $n-1$  corrects that bias & gives you an unbiased estimate of the population variance.

## Random Variables (X) :-

A random variable is a mathematical function that assigns numerical values to the outcomes of a random experiment

$$Y = 5x + 5$$

Tossing coin

$$x = \{0, 1\}$$

Rolling a die

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

Rolling a fair die

Types of Random variables

① Discrete random variable :- A Discrete random variable is used to quantify the outcome of a random experiment. It takes a countable number of possible outcomes. It is also called a stochastic variable

Generally counted on  $0, 1, 2, 3, \dots$

def: It is a type of a variable whose value depends upon the numerical outcomes of a certain random experiment

### Continuous Variable:

It is a variable that takes any real value within a given range or interval, including fractions & decimals. Its possible values are uncountably infinite.

Ex :- Temperature on a given day

26.4°C, 26.45°C, etc

### Percentiles and Quartiles:

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

# No. of odd numbers = 3

$$\text{Percentage of odd numbers in a group} = \frac{3}{6} \times 100 = 50\%$$

Percentile: A percentile is a value below which certain percentage of observations lie.

$$\{2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 9, 9, 10\} \quad n=9$$

$$\frac{3+4}{2} = \text{Percentile} = \frac{\text{Number of values below } x}{n} \times 100$$

$$= \frac{11}{14} \times 100 = 78\% \text{ of value 9} \quad 25\% \text{ is } 3.75$$

$$\Rightarrow \text{Value} = \frac{\text{Percentile} \times (n+1)}{100}$$

$$= \frac{25}{100} \times (15) \\ = 3.75 \\ \neq 3.5$$

$$\frac{3+4}{2} = 3.5$$

Quartiles

25%      50%      75%

25% = First Quartile

50% = 2<sup>nd</sup> Quartile

75% = 3<sup>rd</sup> Quartile

## 5 Number Summary

used in FF A to remove outliers

① Minimum

② 1<sup>st</sup> Quartile (25%)

Ex: 1, 2, 2, 2, 3, 3, 4, 5, 5, 6, 6, 6, 7, 8, 8, 9, 27

outlier

③ Median

④ 3<sup>rd</sup> Quartile (75%)

⑤ Maximum

\* The first step to remove outlier is to define

lower fence & higher fence value

lower fence  $\rightarrow$  higher fence value

- Any thing lower than the lower fence value is removed as outlier.

- Any thing higher than the <sup>higher</sup> fence value will be removed as outliers

$$\text{Lower fence} = Q_1 - 1.5 \times (\text{IQR})$$

$$\text{Higher fence} = Q_3 + 1.5 \times (\text{IQR})$$

$$Q_1 = \frac{\text{Percentile}}{100} \times (n+1) = \frac{25}{100} \times (19+1) = 5^{\text{th}} \text{ position}$$

$$Q_3 = \frac{75}{100} \times 20 = 15^{\text{th}} \text{ position}$$

$$IQR = Q_3 - Q_1 = 7 - 3 = 4$$

Lower fence =  $Q_1 - 1.5(IQR)$  higher fence =  $Q_3 + 1.5(IQR)$

$$= 3 - 1.5(4)$$

$$= 7 + 1.5(4)$$

$$= 3 - 6 = \boxed{-3}$$

$$= 7 + 6 = \boxed{13}$$

After removing outliers

Ex:- 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9

$$\text{Minimum} = 1$$

$$1^{\text{st}} \text{ Quartile} = 3$$

$$\text{Median} = 5$$

$$3^{\text{rd}} \text{ Quartile} = 7$$

$$\text{Maximum} = 9$$

This information indicates the 5 number Summary

\* Based on the particular value we can define a

box plot. (Important)

\* What kind of point you can actually use to see or visualize the outliers? box plot

\* One more plot is Whisker plot:

## Histogram And Skewness

A histogram is a graphical representation of the numerical data. It is an estimate of the probability distribution of a continuous variable and is used to visualize the shape, central tendency & variability of a dataset.

$$\text{Ages} = \{11, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50\}$$

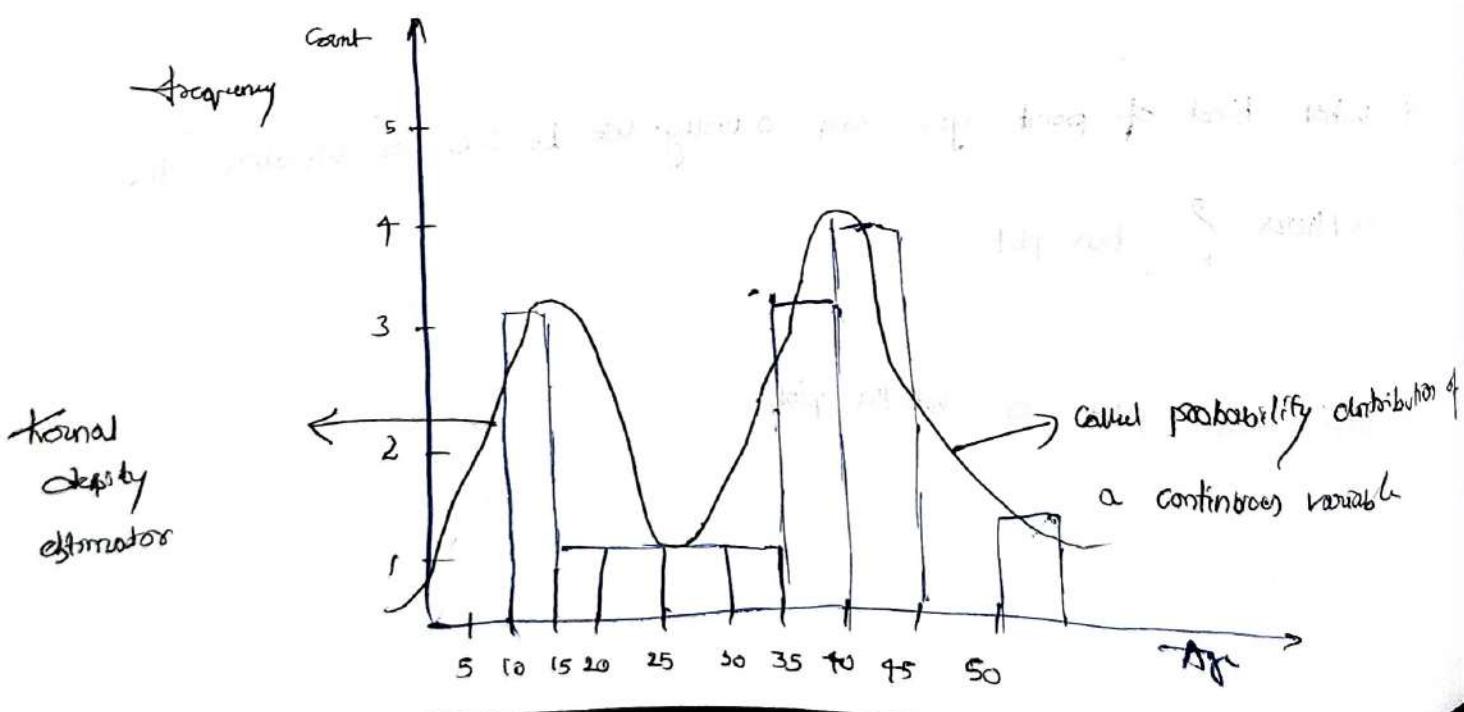
$\Rightarrow$  Histogram

- ① Bin refers to a interval or range of values in to which data points are grouped

$$0 - 50$$

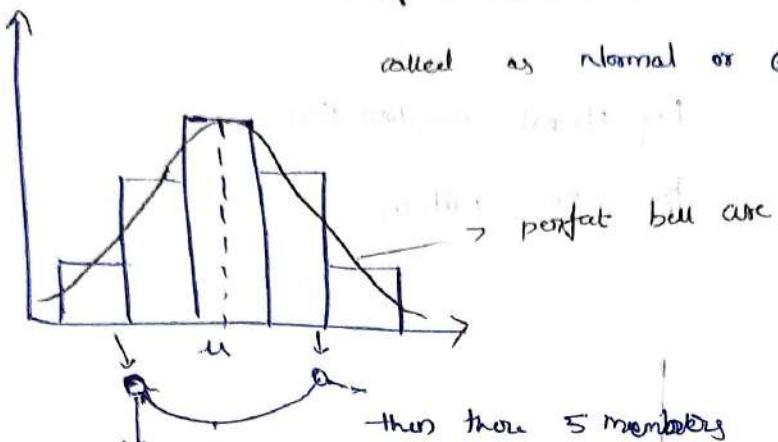
② No. of bins = 10  $\Rightarrow \frac{0-50}{10} = 5$  bin size.

$$\text{Bin} \rightarrow [0-5, 5-10, 10-15, \dots, 45-50]$$



Skewness:-

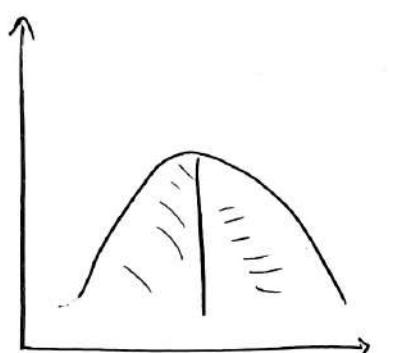
- If a distribution is having perfect bell curve then it is called as normal or Gaussian distribution.



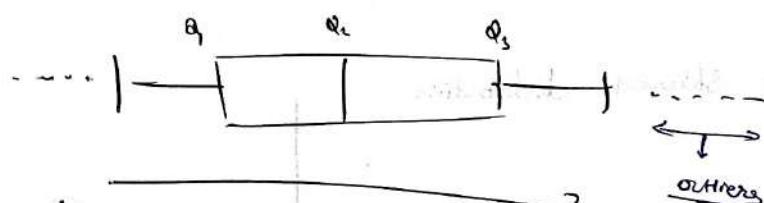
It has 5 members

It is Symmetrical distribution & It has no skewness, perfect in shape.

50% of PD one side & 50% of PD other side



Box plot



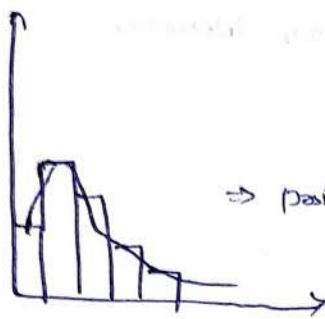
⑧

- The mean, median, mode will all perfectly at the centre

Mean = Median = Mode

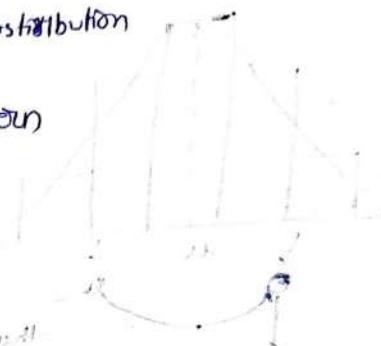


### ② Right Skewed:



Log Normal Distribution

have same pattern



Mean is mostly impacted

$$\text{mean} \geq \text{median} \geq \text{mode}$$

$$Q_3 - Q_2 \geq Q_2 - Q_1$$

Box plot

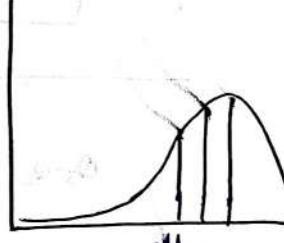


=

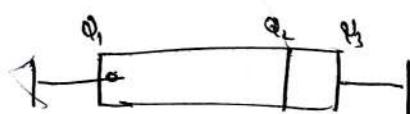
### ③ Left Skewed distribution

Relation

$$\text{mean} \leq \text{median} \leq \text{mode}$$



⇒ negative skewed



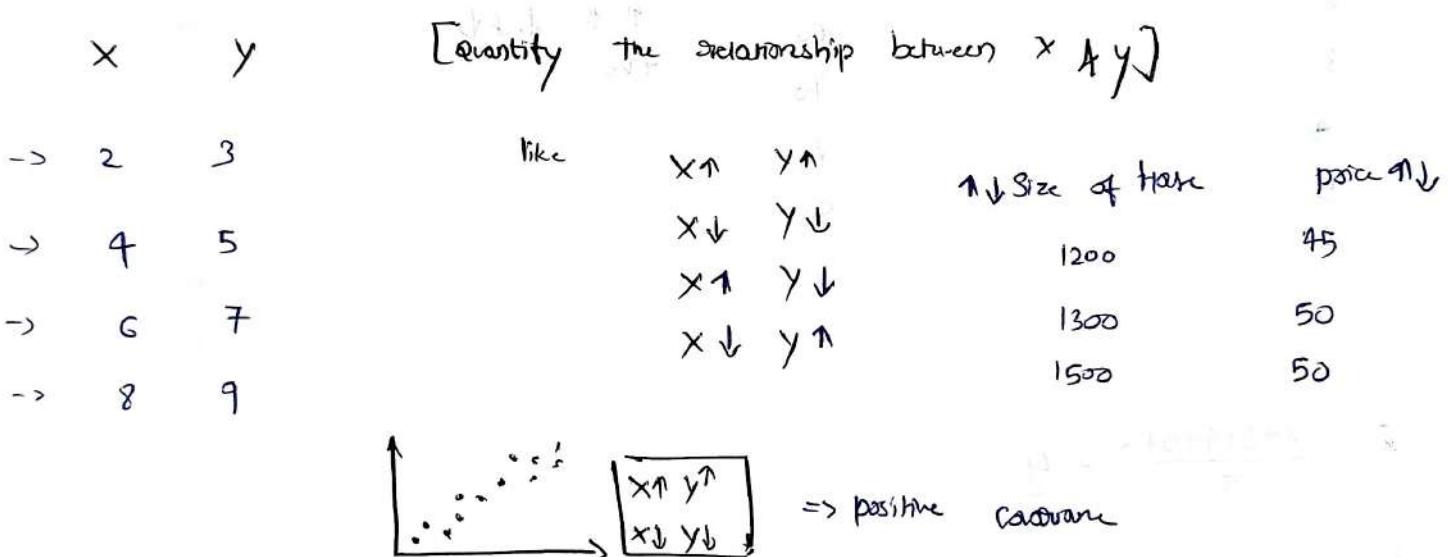
$$Q_2 - Q_1 \geq Q_3 - Q_2$$

## Covariance and Correlation:

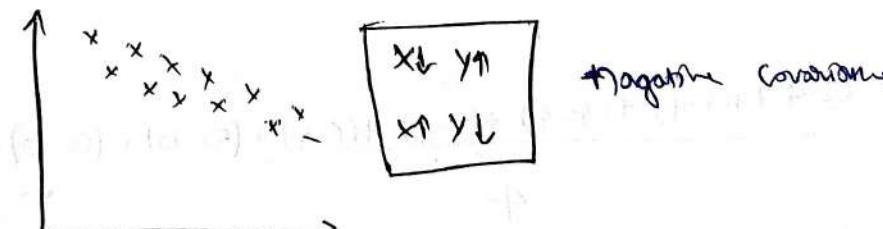
Covariance & Correlation are two statistical measures used to determine the relationship between two variables. Both are used to understand how changes in one variable are associated with change in another variable.

### Covariance:

Covariance is a measure of how much random variables change together. If the variables tend to increase & decrease together, the covariance is positive. If one tends to increase when the other decreases, the covariance is negative.



$x$	$y$
7	10
6	12
5	14
4	16



## Covariance

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

$x_i \rightarrow$  Data point of RV  $X$

$\bar{x} \rightarrow$  Sample mean of  $X$

$y_i \rightarrow$  Data point of RV  $y$

$\bar{y} \rightarrow$  Sample mean of  $y$

$$\boxed{\text{Cov}(x,x) = \text{Var}(x)}$$

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

Student:

Hours Studied ( $X$ )

Exam Score ( $Y$ )

2

50

3

60

4

70

5

80

6

90

70

60

80

70

90

60

$$\bar{x} = \frac{2+3+4+5+6}{5} = 4$$

$$\bar{y} = \frac{50+60+70+80+90}{5} = 70,$$

$$\begin{aligned} \text{Cov}(x,y) &= \frac{(2-4)(3-4)+(3-4)(4-4)+(4-4)(5-4)+(5-4)(6-4)+(6-4)(70-70)+(50-70)+(60-70)+(70-70)+(80-70)+(90-70)}{4} \\ &= 20 \end{aligned}$$

This +ve covariance indicate the no of hours studied increased the exam score.

x	y
7	10
8	12
9	14

$x \uparrow y \downarrow$

$\Rightarrow -ve$

$Cov(A, B)$  we cannot say whether  
 $Cov(B, C)$  if not have limit

Advantages:

Disadvantages:

- ① Quantify the relationship between

$x \wedge y$

- ① Covariance does not have a

specific limit value

$$Cov(x, y) \Rightarrow -\infty \text{ to } \infty$$

$\rightarrow -1 \text{ to } 1 \text{ limits } 0.96 < 0.98$

- ② Correlation:

Pearson Correlation Coefficient

Spearman Rank Correlation

- ① Pearson Correlation Coefficient  $\Rightarrow [-1 \text{ to } 1]$

$$\rho_{x,y} = \frac{Cov(x, y)}{\sigma_x \cdot \sigma_y}$$

$\Rightarrow$  it limits to  $-1 \text{ to } 1$  if we can say

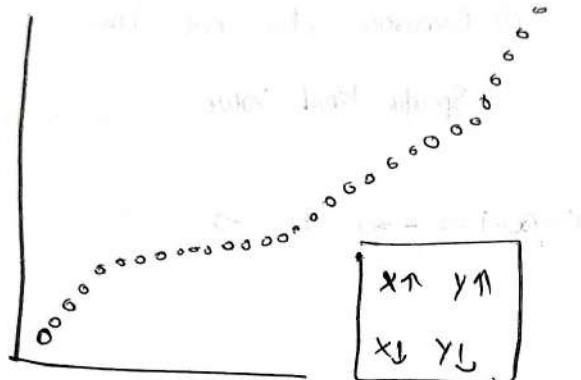
which is Strong

- The more the value towards  $+1$  the more +ve correlated  $x \wedge y$
- The more the value towards  $-1$  to more -ve correlated  $x \wedge y$

### (ii) Spearman Rank Correlation:

Spearman correlation = 1

Pearson correlation = 0.98 (Because it is non linear)



$x \ y \ R(x) \ R(y)$

	1	2	3	4	5	6	7	8
$R(x)$	1	2	3	4	5	6	7	8
$R(y)$	1	2	3	4	5	6	7	8
$x$	0	1	2	3	4	5	6	7
$y$	0	1	2	3	4	5	6	7

$$r_s = \frac{\text{Cor}(R(x), R(y))}{\sigma(R(x)) \sigma(R(y))}$$

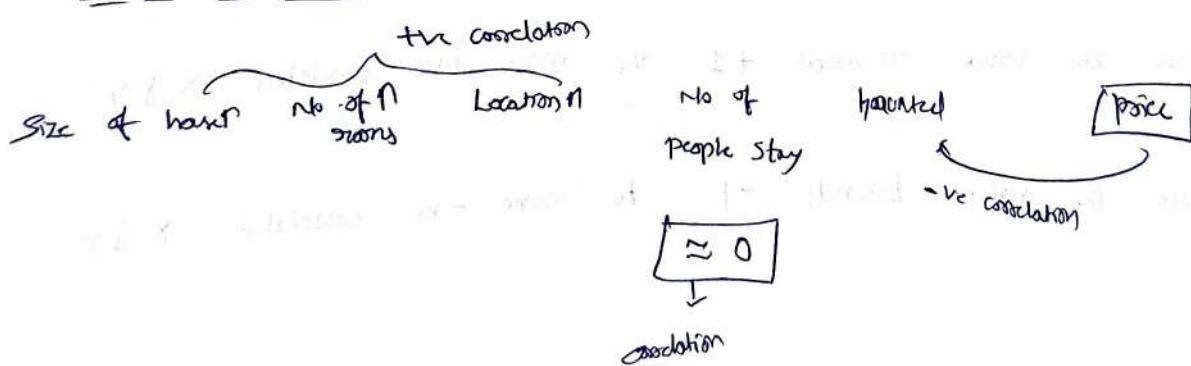
$$\sigma(R(x)) = \sqrt{\frac{1}{n} \sum_{i=1}^n (R(x)_i - \bar{R(x)})^2}$$

Spearman correlation of 1 results when 2 variables being compared are monotonically related, even if relation is not linear.

$\Rightarrow X \uparrow \rightarrow Y \uparrow$

$X \downarrow \rightarrow Y \downarrow$

Used in feature selection:



# Probability

① Introduction

② Addition Rule (For mutually exclusive event)

③ Addition Rule (For non mutually exclusive event)

④ Multiplication Rule (Independent & Dependent Event)

① Probability:- It is about determining the likelihood of an event

Eg:- • Toss a coin

$$P(H) = \frac{1}{2} \quad \{H, T\}$$

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

• Rolling a dice  $\{1, 2, 3, 4, 5, 6\}$

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

Mutually exclusive events:

Two events are mutually exclusive if they cannot occur at the same time

$$P(H \text{ or } T) = P(H) + P(T)$$

{Addition rule for mutually exclusive events}

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

Eg: Rolling a die  $\{1, 2, 3, 4, 5, 6\}$

$$P(1 \text{ or } 5) = P(1) + P(5)$$

$$= \frac{1}{6} + \frac{1}{6} = \frac{2}{6} = \frac{1}{3}$$

## \* Non Mutual Exclusive Events:

Two event are mutually exclusive if they can occur at the same time

Eg:- Taking a card from deck

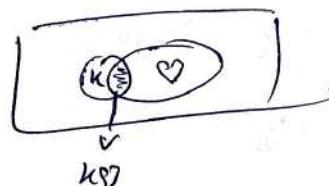
$\boxed{K}$  again  $\boxed{K}$  can get

$$P(K \text{ or } Q) = P(K) + P(Q) - P(K \cap Q)$$

↑  
Non mutual exclusive events

$$= \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

$$= \frac{16}{52}$$



Multiplication rule :- {Independent & dependent events}

2 events are independent, if they do not affect one another

Eg:- Tossing a coin {H & the Tail}

$$P(H) = \frac{1}{2} \quad P(T) = \frac{1}{2}$$

Eg:- Rolling a dice

Dependent events:

2 events are Dependent if they effect each other

e.g.: Take = taking a card & then Taking a Queen card

$$P(K) = \frac{4}{52}$$

$$P(Q) = \frac{4}{51}$$

Multiplication rule:

① Independent event {Taking a card}

$$P(H \text{ and } T) = P(H) * P(T)$$

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

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

② Dependent Event

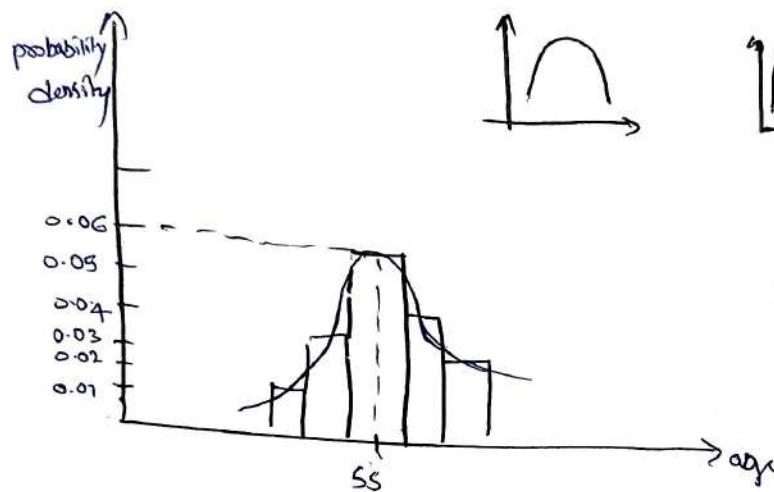
$$P(K \text{ and } Q) = P(K) * P(Q|K)$$

$$= \frac{4}{52} * \frac{4}{51}$$

## Probability Distribution Function:

Probability distribution function → helps us to understand how describe how the probabilities are distributed over the values of a random variable

$\text{Age} = \{ \dots \}$  → continuous random variable



Two main types of probability distribution function

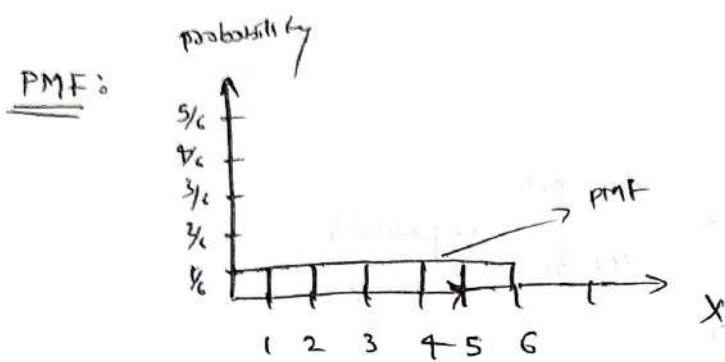
① Probability Mass function (PMF) :- used for discrete

② Probability Density function (PDF) :- used for continuous

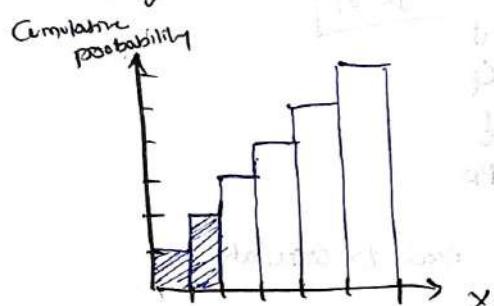
① Probability Mass function : [Discrete Random Variable]

Eg:- Rolling a dice  $\{1, 2, 3, 4, 5, 6\} \Rightarrow$  fair dice

$$P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = \frac{1}{6}$$



Cumulative Density function:-

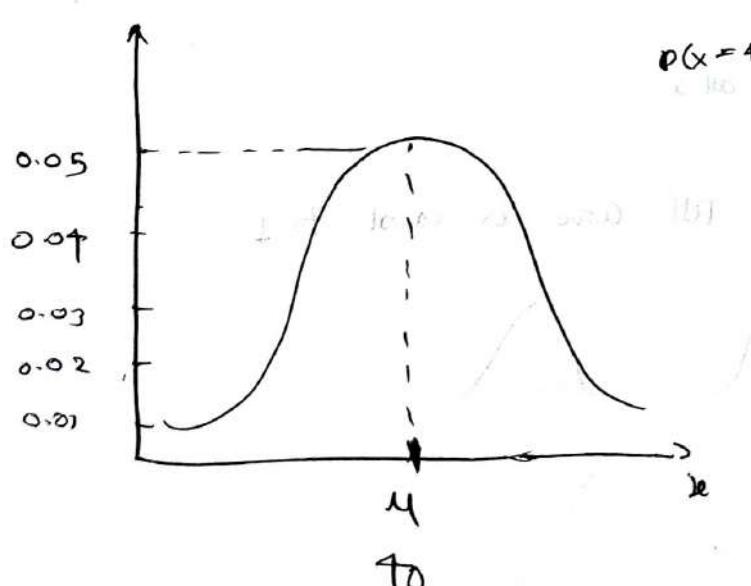


$$P(X \leq 2) = P(X=1) + P(X=2) \\ = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}$$

Cumulative density function we need to combine all the probabilities as

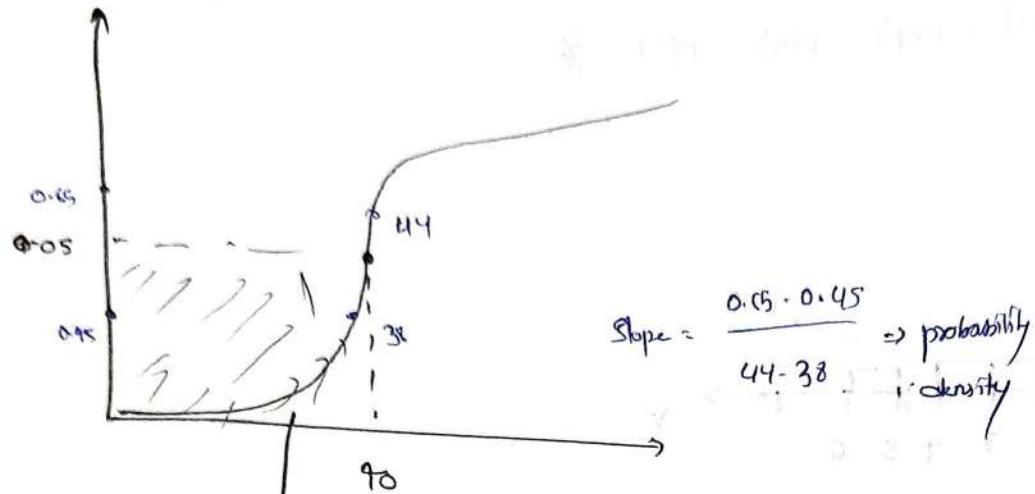
we go from 1 to 6 and add up to get the cumulative density function.

② Probability Density Function: [Distribution of continuous random variables]



$$P(X=10) = 0.05 \quad (\text{X})$$

Cumulative probability



$P(x) = 5\%$  to fall in  
this area This area is given by PDF

$$\text{Slope} = \frac{x_2 - x_1}{y_2 - y_1}$$

↓  
G  
↓  
PD

In order to calculate probability density function we need to calculate

Gradient of CPF / slope of that point

Probability density function = Gradient of Cumulative Density function

Properties

① Non negative  $f(x) \geq 0$  for all  $x$

② The total area under the Ptf curve is equal to 1

$$\int_{-\infty}^{\infty} f(x) dx = 1 \Rightarrow$$

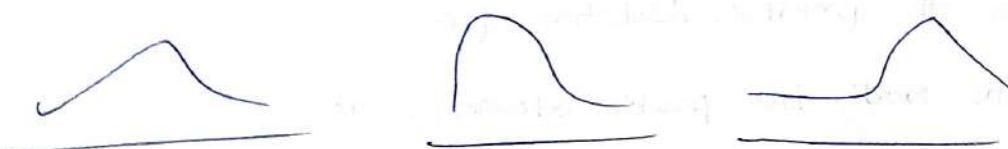
With respect to different distribution

$f(x)$  function has going to change

## Different types of distribution:-

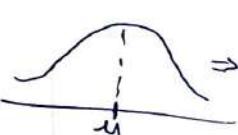
Types of Probability Distribution (pdf, pmf, cdf)

Age, weight, salary



① Bernoulli Distribution: outcomes are binary (pmf)  $\Rightarrow$  Discrete random variable

② Binomial Distribution  $\rightarrow$  pmf

③ Normal / Gaussian Distribution  $\rightarrow$  pdf  $\Rightarrow$   Assumptions

④ Poisson Distribution pmf

⑤ Log Normal Distribution pdf

⑥ Uniform Distribution pmf

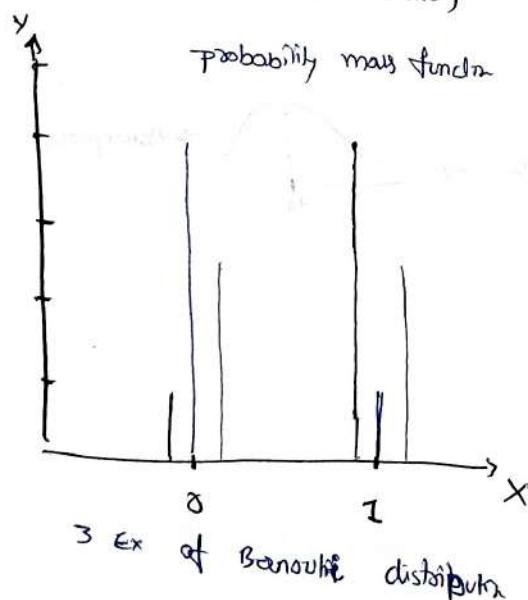
Dataset  $\rightarrow$  HPPD

Size of house	No. of rooms	Location	Floor	Sea side	Price
CRV	Discrete		1	0/1	Continuous
			Discrete	pmf	pdf
			pmf		

Bernoulli Distribution:-

The Bernoulli's distribution is the simplest discrete probability distribution. It represents the probability distribution of a random variable that has exactly two possible outcomes: success (with probability  $p$ ) and failure (with probability  $1-p$ ). It is used to model binary outcomes, such as a coin flip or yes/no question.

Bernoulli Distribution



$$P(x=0) = 0.2 \quad + P(x=1) = 0.8$$

$$P(x=0) = 0.8 \quad 4 P(x=1) = 0.2$$

$$P(x=0) = 0.5 \quad \& P(x=1) = 0.5$$

## ① Discrete Random Variable

### ② Outcomes are Binary

Ex: Student will pass/fail

$$P(X = \text{pass}) = 0.4$$

$$P(X = \text{fail}) = 1 - 0.4 = 0.6$$

Parameters:  $0 \leq p \leq 1$

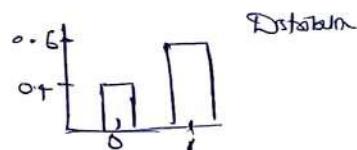
$$q = 1 - p$$

Suppose  $X \in \{0, 1\} \Rightarrow 2 \text{ outcomes}$

### PMF

$$k=1 \text{ pass} = 40\% \Rightarrow P$$

$$\text{fail} = 60\% \Rightarrow q$$



$$\text{PMF} \left\{ \begin{array}{ll} q = 1-p & \text{if } k=0 \\ p = P & \text{if } k=1 \end{array} \right.$$

$$\text{PMF} = p^k * (1-p)^{1-k}$$

if  $k=1$

$$P(k=1) = p^1 (1-p)^{1-1} = p$$

If  $k=0$

$$P(k=0) = p^0 (1-p)^{1-0} = 1-p = q$$

## Mean:

The expected value of a Bernoulli random variable  $X$  is

$$E[X] = p \quad K = \{0, 1\}$$

$$E[X] = \sum_{k=0}^1 k \cdot P(X=k) \quad p = 0.4 \quad n = 0.4$$

$$= 0 \cdot 0.40 + 1 \cdot 0.60$$

$$= 0 + 0.6$$

$$= 0.6 = q$$

$$= p$$

## Mode

~~Mode~~  $\rightarrow$  ~~P~~  $\rightarrow$  ~~Bin~~

Median of Bernoulli

$$\text{median} \begin{cases} 0 & \text{if } p < \frac{1}{2} \\ (0, 1) & \text{if } p = \frac{1}{2} \\ 1 & \text{if } p > \frac{1}{2} \end{cases} \quad \left\{ \begin{array}{ll} \text{median} = 0 & \text{if } q > p \\ = 0.5 & \text{if } q = p \\ = 1 & \text{if } q < p \end{array} \right.$$

## Mode

$p > q \Rightarrow p$  will be the mode

~~or~~ or will be the mode

## \* Variance

(K = 0, 1)

$$\begin{aligned}
 E[X^2] &= P(X=1) \cdot K + P(X=0) \cdot 0 \\
 &= P(X=1) \cdot 1^2 + P(X=0) \cdot 0^2 \\
 &= p \cdot 1^2 + q \cdot 0^2 \\
 &= p = E[X]
 \end{aligned}$$

$$\text{Var}[X] = E[X^2] - E[X]^2$$

$$\begin{aligned}
 &= E[X] - E[X]^2 \\
 &= p - p^2 \\
 &= p(1-p) = pq
 \end{aligned}$$

$$\begin{aligned}
 \sigma^2 &= 0.40 \times (0-0.6)^2 + 0.6(1-0.6)^2 \\
 &= 0.40 \times 0.36 + 0.6(0.16) \\
 &= 0.24
 \end{aligned}$$

$$= pq$$

$$\boxed{
 \begin{aligned}
 \sigma^2 &= pq \\
 \sigma &= \sqrt{pq}
 \end{aligned}
 }$$

*Binomial test of statistical significance*

## Binomial Distribution

→ It is basis for

Binomial

test of statistical significance

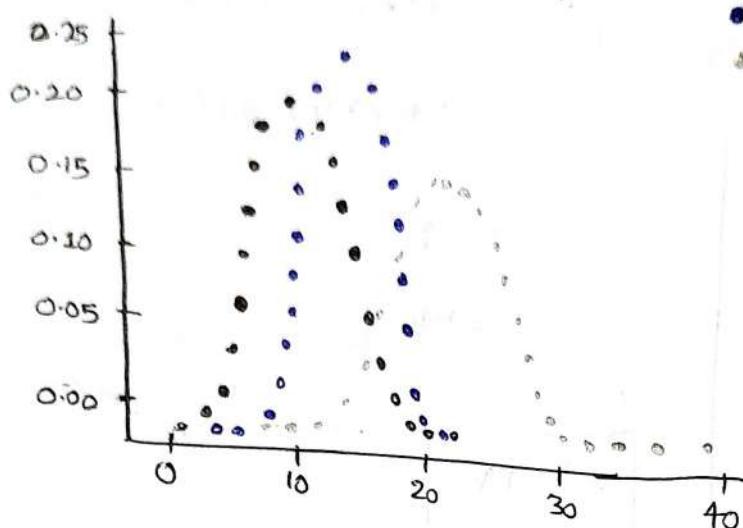
In probability theory & statistics, the binomial distribution with parameter  $n$  &  $p$  is the discrete probability distribution of no of success in a sequence of  $n$  independent experiments, each asking Yes/No questions and each with its own Boolean-valued outcomes: success ( $p$ ) or failure ( $q$ )

⇒ A single Success/failure experiment is also called a Bernoulli trial or Bernoulli Experiment & a sequence of outcomes are called a Bernoulli process

⇒ for single trial Binomial distribution is a Bernoulli distribution.

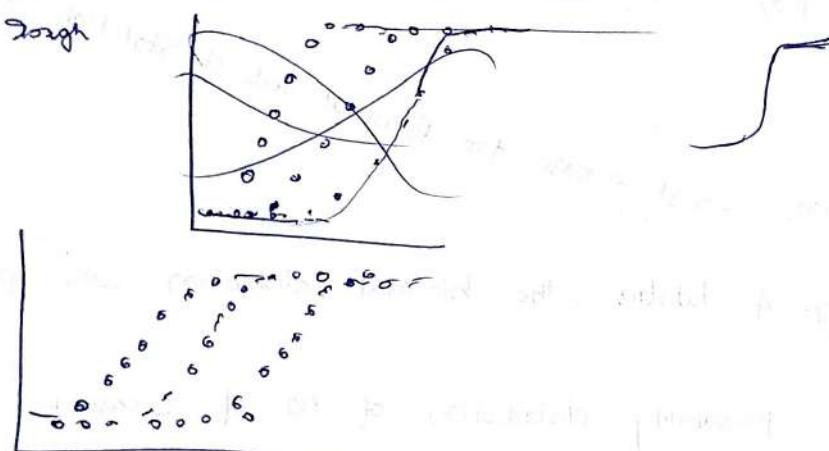
## Binomial distribution

### Probability mass function



- $p=0.5 \text{ & } n=20$
- $p=0.7 \text{ & } n=20$
- $p=0.5 \text{ & } n=40$

### Cumulative distribution function



### \* Discrete random variable

- Every outcome of the experiment is binary
- These experiments are performed for n trials

Ex: Tossing a coin 10 times

$$n=10$$

$$\{H, T\}$$

## Notation

$B(n,p)$

parameters :  $n \in \{0, 1, 2, \dots\} \Rightarrow$  no of trials or experiment

$p \in [0, 1] \rightarrow$  success probability for each trial

$$q = 1-p$$

Support :  $k \in \{0, 1, 2, 3, \dots, n\} \Rightarrow$  no of successes

PMF  $P(k, n, p) = {}^n C_k p^k (1-p)^{n-k}$

For  $k = 0, 1, 2, \dots, n$  where

$$\boxed{{}^n C_k = \frac{n!}{k!(n-k)!}} \Rightarrow \text{Binomial Coefficient}$$

Mean :  $n \cdot p$

Variance :  $n p q$

$$\sigma = \sqrt{npq}$$

Ex: ① Coin flip

No. of trials ( $n$ ) = 5

Probability of success ( $p$ ) = 0.5

No. of success ( $k$ ) = varies from 0 to 5

Q. What is the probability of getting exactly 3 heads in 5 flips?

$$n=5 \quad k=3$$

$$P(X=3) = {}^5C_3 (0.5)^3 (1-0.5)^{5-3}$$

$$= 0.3125$$

### Example:

Scenario: Inspecting 10 items in a factory where each item has a 10% chance of being defective

\* No. of Trials ( $n$ ) = 10

\* Probability of Success ( $p$ ) = 0.1 (defective item)

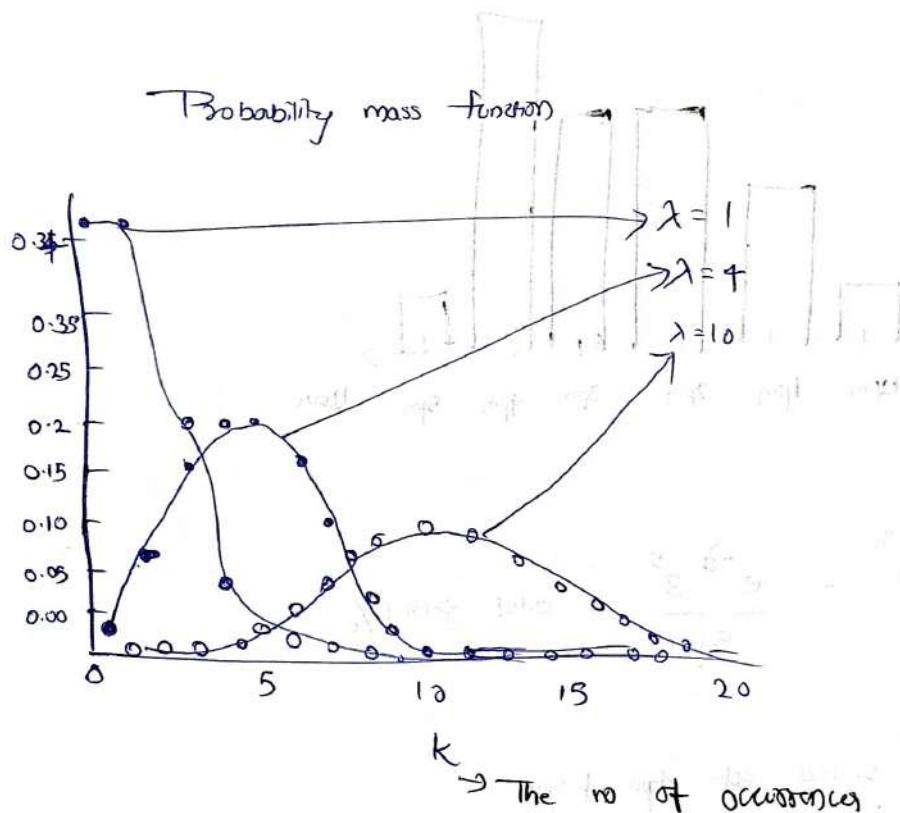
\* No. of Successes ( $k$ ) = varies from 0 to 10

Q. What is the probability of finding exactly 2 defective items in a sample of 10?

$$P(X=2) = {}^{10}C_2 (0.1)^2 (1-0.1)^{10-2} \approx 0.1937$$

## Poisson distribution :-

In probability theory & statistics, the poisson distribution is a discrete probability distribution that expresses the probability of a given no of events occurring in a fixed interval of time if those events occurs with a known constant mean rate & independently of the time since the last event.



$\lambda$  is the expected rate of occurrences

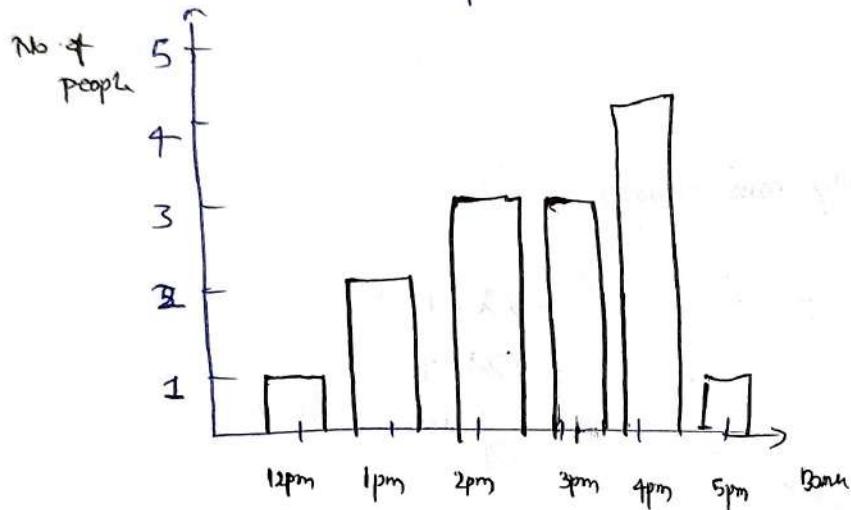
y-axis is probability of  $k$  occurrences given  $\lambda$ .

The function is defined only at integer values of  $k$ .

- ① Discrete random variable (PMF)
- ② Describe the no of events occurring in a fixed time intervals

Eg: No of people visiting hospital every hour

No of people visiting banks every hour



$$P(X=5) = \frac{e^{-\lambda} \lambda^5}{5!} = \frac{e^{-3} 3^5}{5!} = 0.101 \Rightarrow 10.1\%$$

— how many people visited at 4pm & 5pm

$$P(X=4) + P(X=5)$$

Mean of poisson Distribution:

$$\text{mean Absolute deviation} = E[X-\lambda] = \frac{\lambda^{(N+1)} e^{-\lambda}}{(N!)}$$

$$\text{mean} \neq \text{variance} = \lambda$$

$$\lambda \times t$$

?

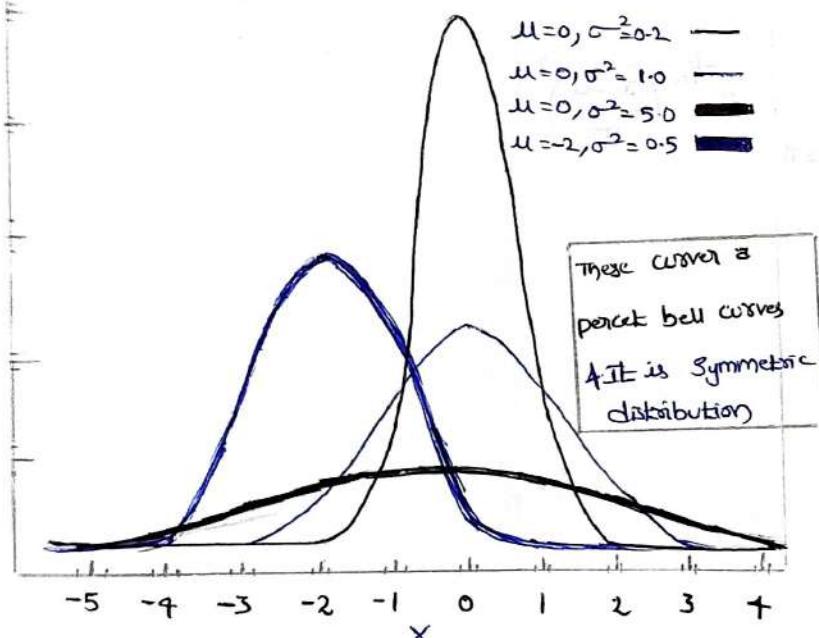
# Normal / Gaussian Distribution

In probability theory and statistics, a normal distribution or Gaussian Distribution is a type of continuous probability distribution for a real-valued random variable. The general form of its probability density function is

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

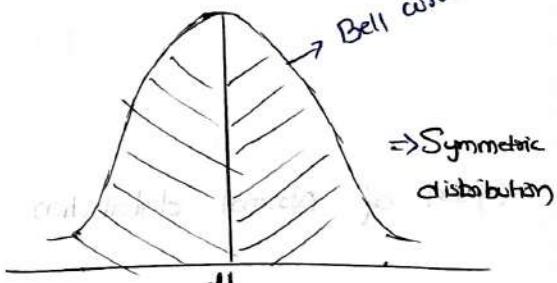
→ Continuous Random values (PDF)

probability density function



→ Symmetric distribution

→ Bell curve



mean ( $\mu$ ) = median = mode

Ex: IRIS Dataset

It usually follows normal distribution

↑ imp feature it has → Petal Length, Sepal length

Petal width, Sepal width

Based on this ↑ features, what flower it is

$\text{Ex } ②$  Weight of student in a class

$\text{Ex } ③$  Height of Students in a class

Notations :-  $N(\mu, \sigma^2)$

parameters:-  $\mu \in \mathbb{R} = \text{mean}$

$\sigma^2 \in \mathbb{R} > 0 = \text{Variance}$

$x \in \mathbb{R}$

Probability Density function =  $\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2}$

Mean of Normal distribution

$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$

Variance :-

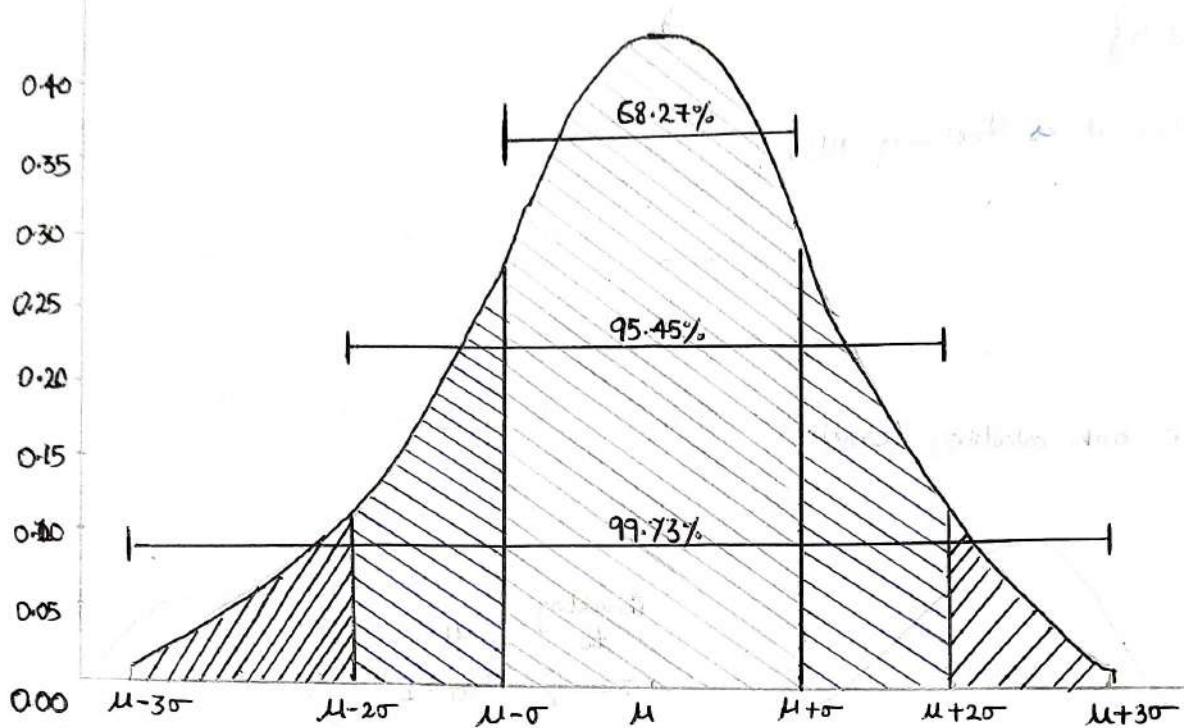
$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \mu)^2}{n}$$

Standard deviation ( $\sigma$ ) =  $\sqrt{\text{Variance}}$

$$= \sqrt{\sigma^2}$$

$$= \sigma$$

# Empirical Rule of Normal / Gaussian Distribution



$X = \{x_1, x_2, x_3, \dots, x_n\}$  Assume it follows Gaussian distribution

\* QQ plot is used to determine whether a random variable  $X$  will follow a Normal / Gaussian distribution or not.

Probability:-

$$P(\mu - \sigma \leq x \leq \mu + \sigma) \approx 68\%$$

$$P(\mu - 2\sigma \leq x \leq \mu + 2\sigma) \approx 95\%$$

$$P(\mu - 3\sigma \leq x \leq \mu + 3\sigma) \approx 99.7\%$$

1 Standard normal Distribution { Any distribution having  $\mu=0$  &  $\sigma=1$  }

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

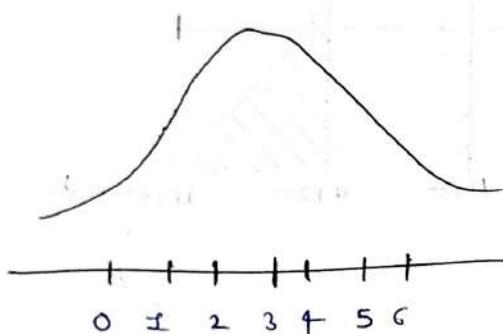
↳

consider it as following N(μ, σ²)

$$\mu = 3$$

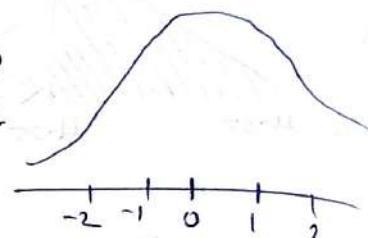
$$\sigma = 1.414$$

$\approx 1$  (To make calculation easy)



Converting  
to

$$\begin{array}{l} \mu = 0 \\ \sigma = 1 \end{array}$$



\* Normal Distribution is converted such that

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

$\mu = 0$  &  $\sigma = 1$  then the distribution is

Called Standard normal distribution

→ We uses Z-score =  $\frac{x_i - \mu}{\sigma}$  for converting

↓

$$\textcircled{1} \frac{1-3}{1} = -2$$

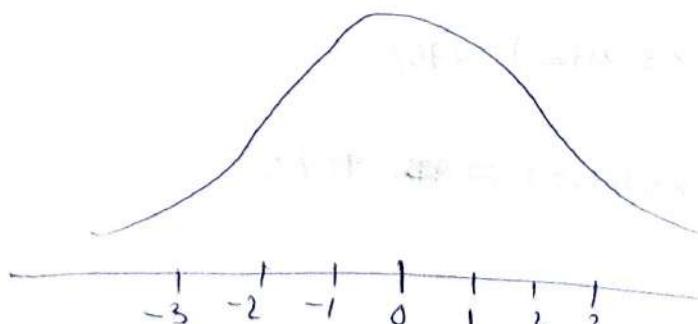
$$Y = \{-2, -1, 0, 1, 2\}$$

$$\textcircled{2} \frac{2-3}{1} = -1$$

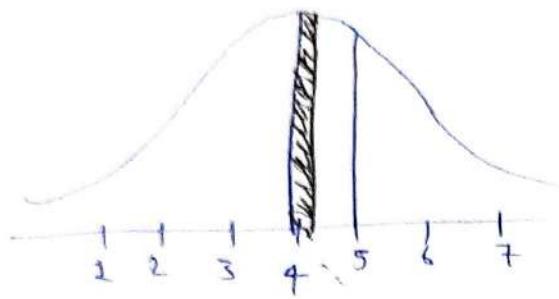
$$\textcircled{3} \frac{3-3}{1} = 0$$

$$\textcircled{4} \frac{4-3}{1} = 1$$

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



$$X \approx \text{SNIS} \quad (\mu=0, \sigma=1)$$



How many standard deviations 4.25 is away from the mean

$$x_i = 4.25$$

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

Eg: Dataset

Age	Years kg weight	Height cm's	INR Salary
24	70	175	40k
25	60	160	50k
25	70	170	35k
26	35	163	60k
20	70	175	45k
31	90	180	65k

- In feature units will be different, we try to bring all the feature in a same unit
- we use Standardization

$$Z\text{-score} = \frac{x_i - \mu_{\text{age}}}{\sigma}, \quad \frac{x_i - \mu_{\text{weight}}}{\sigma}$$

# Uniform Distribution

① Continuous uniform Distribution

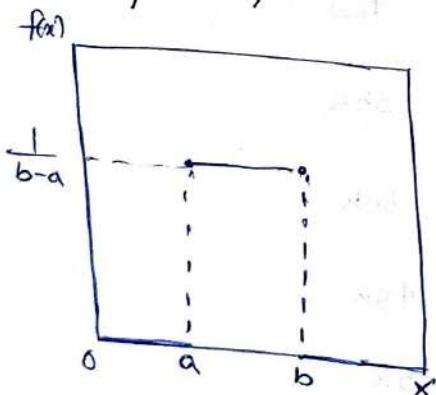
② Discrete uniform Distribution

① Continuous uniform Distribution (crv)

In probability theory and distribution statistics, the continuous uniform distributions or rectangular distributions are a family of symmetric probability distributions. Such a distribution describes an experiment where there is an arbitrary outcome that lies between certain bounds. The bounds are defined by the parameters,  $a$  and  $b$  which are the minimum and maximum values.

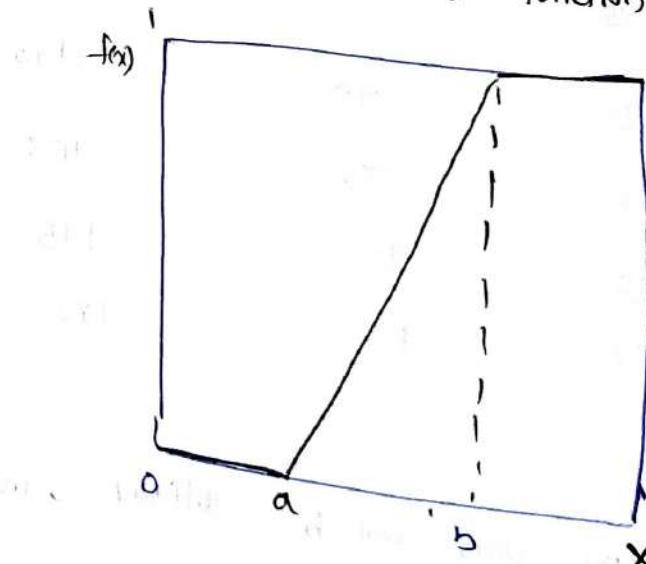
Continuous uniform:

Probability density function



using maximum constant

Cumulative distribution function



Notation :-  $U(a,b)$

parameters :-  $-\infty < a < b < \infty$

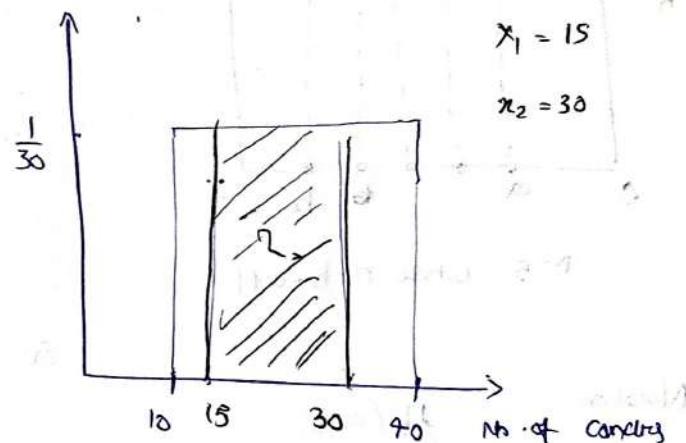
$$\text{pdf} = \begin{cases} \frac{1}{b-a} & x \in [a,b] \\ 0 & \text{otherwise} \end{cases}$$

$$\text{cdf} = \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a} & \text{for } x \in [a,b] \\ 1 & \text{for } x > b \end{cases}$$

$$\text{Variance} = \frac{1}{12} (b-a)^2$$

$$\text{Mean} = \frac{1}{2} (a+b)$$

$$\text{Median} = \frac{1}{2} (a+b)$$



- ① The number of candies sold daily at a shop is uniformly distributed with a maximum of 40 candies and a minimum of 10.

- ② Probability of daily sales to fall between 15 and 30?

$$P(15 \leq x \leq 30) = (x_2 - x_1) \times \frac{1}{b-a} = (30-15) \times \frac{1}{30} = 0.5 //$$

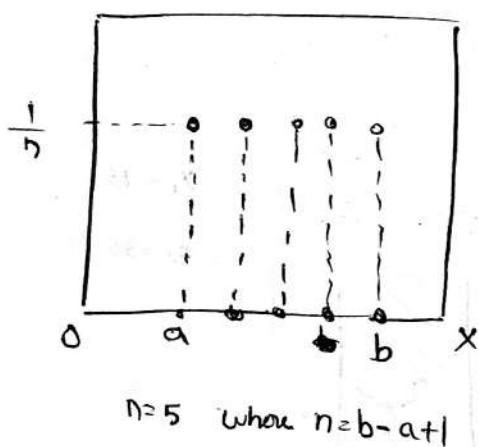
$$P(x \geq 20) = (40-20) \times \frac{1}{30} = 0.66 = 66\%$$

## ② Discrete uniform Distribution :-

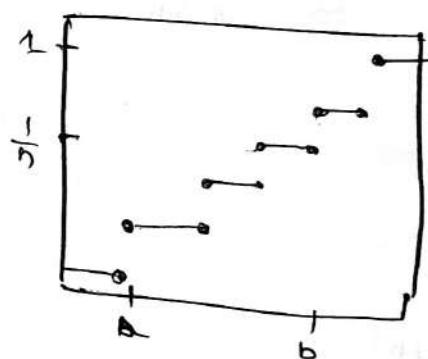
In probability theory & statistics, the discrete uniform distribution is a symmetric probability distribution wherein all finite no. of values are equally likely to be observed; every one of  $n$  values has equal probability equally likely to happen.

Discrete uniform

PMF



Cumulative distribution function



Notation

$$\text{U}(a, b)$$

Ex: Rolling a dice {1, 2, 3, 4, 5, 6}  $P(1)=\frac{1}{6}$

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

parameters

:  $a, b$  where  $b \geq a$

PMF

$$=\frac{1}{n}$$

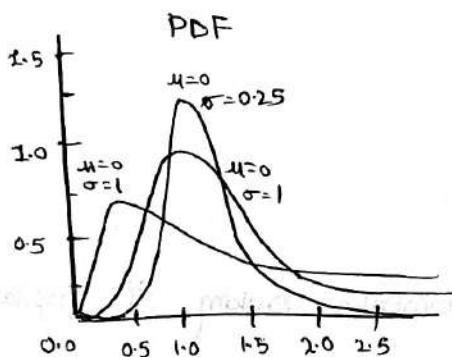
mean

$$\text{median} \rightarrow \frac{a+b}{2}$$

## Log Normal Distribution:

In probability theory, a log-normal (or lognormal) distribution is a continuous probability distribution of a random variable  $x$  as log-normally distributed, if  $y = \ln(x)$  has a normal distribution. Equivalently, if  $y$  has a normal distribution, then the exponential function of  $y$ ,  $x = \exp(y)$ , has a log-normal distribution.

Log-normal D \* it is a right skewed distribution



$x \approx \log \text{Normal distribution } (\mu, \sigma)$

$y \approx \ln(x) \rightarrow \text{Normal distribution}$

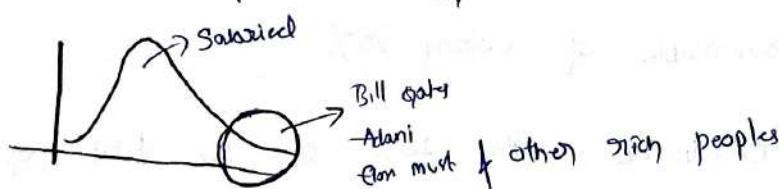
Natural log  
[ $\ln(x)$ ]

checked using Q&P plot

$x \approx \exp(y) \Rightarrow \log \text{normal distribution}$



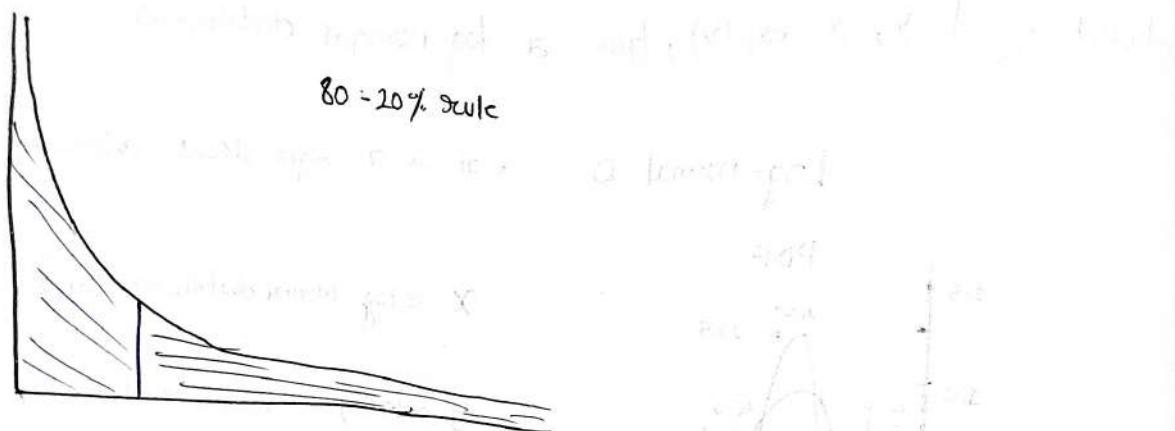
Ex ① Wealth distribution of the world



② Discussion forum → length of the comments

## Power law Distribution:

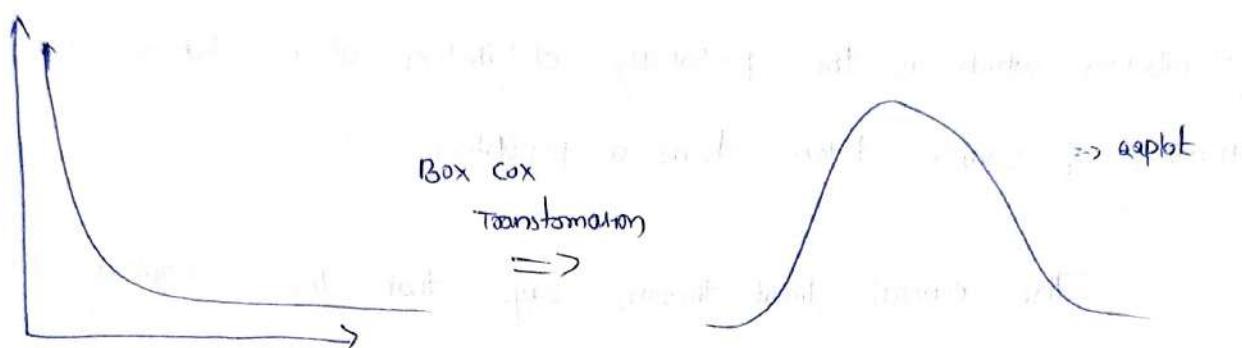
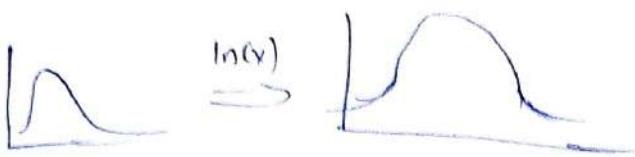
In statistics, a power law is a functional relationship between two quantities, where a relative change in one quantity results in a proportional relative change in the other quantity, independent of the initial size of those quantities: one quantity varies as a power of another.



An example power-law graph that demonstrates ranking of popularity. To the right is the long tail, and to the left are the few that dominate (also known as 80-20 rule).

### Eg:- IPL

- ① 20% of Team is responsible of winning 80%.
- ② 80% of wealth are distributed with 20% of the total population
- ③ 80% of the total oil is with 20% of the nation



The data which follows power law distribution we call it as Pareto distribution.

Project Distribution :-

e.g. 80% of the entire project is done by 20% of the team.

\* In statistics, a power law is a

The Pareto distribution, named after the Italian civil engineer, economist & sociologist Vilfredo Pareto, is a power-law PD that is used in description of social, quality, control, scientific, geophysical, actuarial & many other types.

The principle originally applied to describing the distribution of wealth in a society, fitness the trend that a large portion of wealth is held by a small fraction of the population.

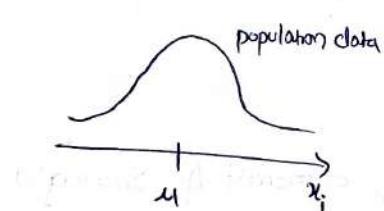
## ① Central Limit Theorem:

The central limit theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large number of samples taken from a population.

The central limit theorem says that the sampling distribution of the mean will always be normally distributed as long as the sample size is large enough. Regardless of whether the population has a normal, poisson, binomial or any other distribution, the sampling distribution of mean will be normal.

$$① X \approx N(\mu, \sigma)$$

$n$  = sample size  $\Rightarrow$  any value



$$\begin{aligned} S_1 &= \{x_1, x_2, \dots, x_n\} = \bar{x}_1 \\ S_2 &= \{x_2, x_3, \dots, x_n\} = \bar{x}_2 \\ &\vdots \\ S_n &= \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n\} = \bar{x}_n \end{aligned} \quad \left. \begin{array}{l} \text{Sample means} \\ \vdots \end{array} \right\}$$

Sample distribution of means

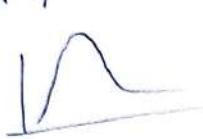


$\Rightarrow$  Gaussian / Normal

Distribution

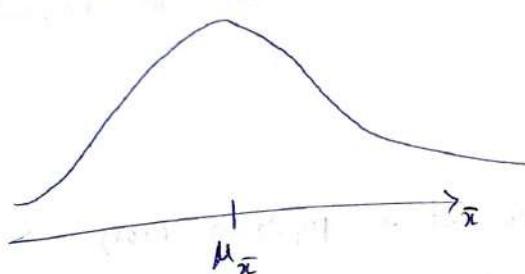
$\mu_x$

$x \neq N(\mu, \sigma)$



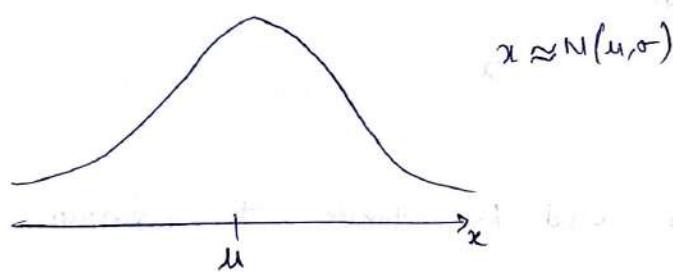
$$S_1 \rightarrow \bar{x}_1 \\ S_2 \rightarrow \bar{x}_2 \\ \vdots \\ S_n \rightarrow \bar{x}_n$$

$n \geq 30$  sample size



Central Limit Theorem

## ① Normal Distribution



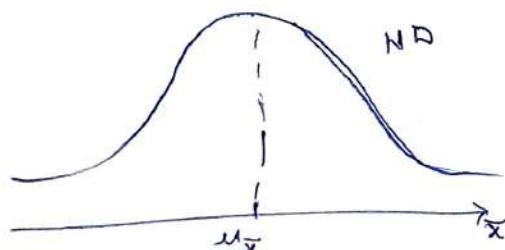
$$x \approx N(\mu, \sigma)$$

$\sigma$  = population std

$\mu$  = population mean

$n$  = Sample size

Sampling Distribution of mean



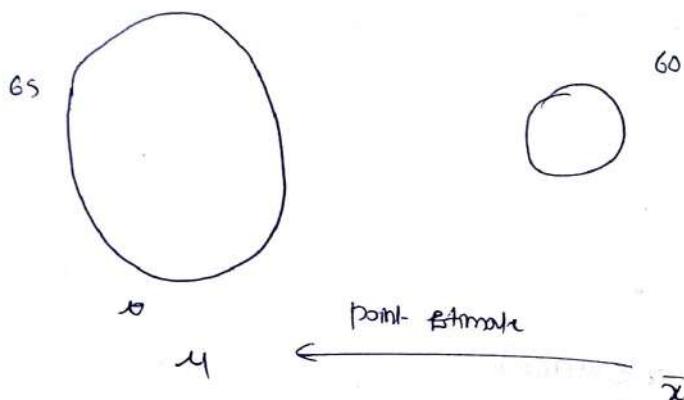
$$x = N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

Estimate :- It is a Specified observed numerical value used to estimate an unknown population parameter

Types:-

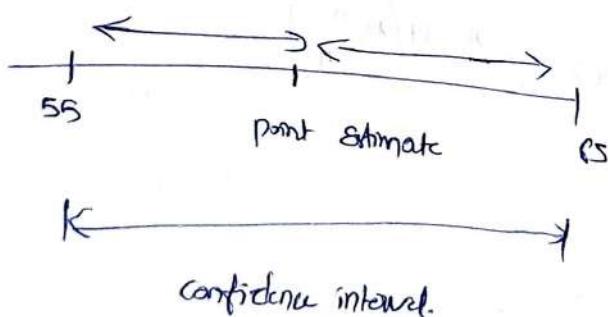
① Point Estimate :- Single numerical value used to estimate an unknown population parameter

Eg :- Sample mean is a point estimate of a population mean



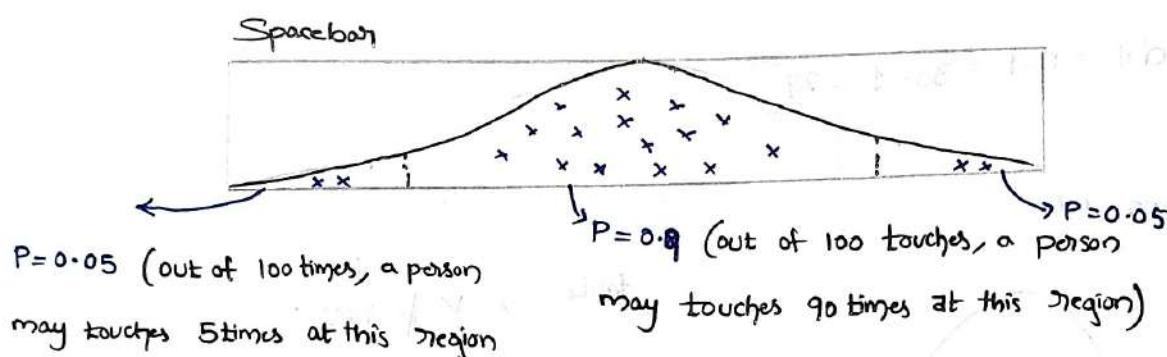
② Interval Estimate :- Range of values used to estimate the unknown population parameter

$$[55 - 65] \Rightarrow \text{Sample mean}$$



P-value: The P-value is a number, calculated from a statistical test that describes how likely you are to have found a particular set of observations if null hypothesis were true. P values are used in hypothesis testing to help decide whether to reject the null hypothesis.

Example 1 :- P-value giving the probability of Spacebar touch at Particular Region



Example 2 :- Coin is fair or not

- ① Null hypothesis :  $H_0$ : coin is fair
- ② Alternate hypothesis :  $H_1$ : coin is not fair
- ③ Experiment :- 100 times
- ④ Significance value :-  $\alpha = 0.05$

$$\text{Confidence interval} = 1 - \alpha$$

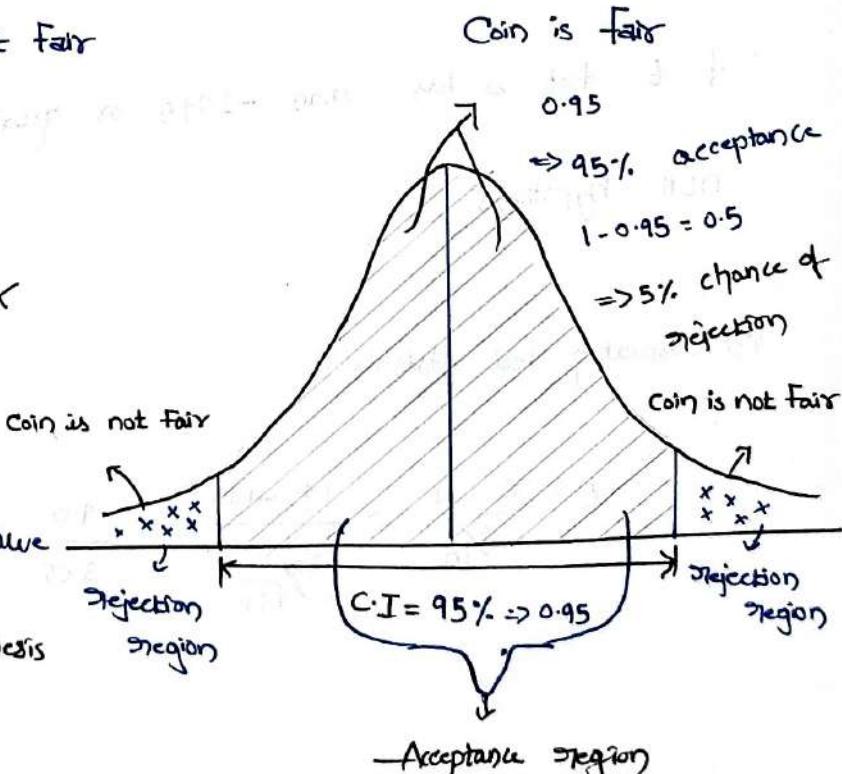
$$\begin{aligned} C.I &= 1 - 0.05 \\ &= 0.95 \end{aligned}$$

- ⑤ Conclusion :- If  $P <$  Significance value

Reject the null hypothesis

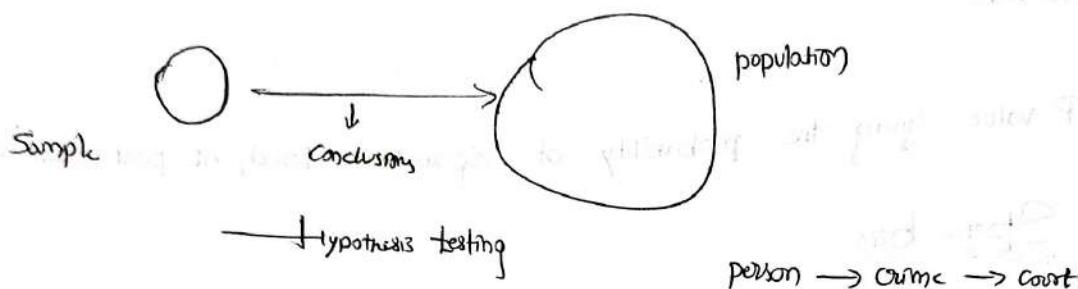
else

Fail to Reject the null hypothesis



## Hypothesis and Hypothesis Testing Mechanism:

Inferential Stats :- Conclusion or Inference



### Hypothesis Testing Mechanism:

① Null hypothesis ( $H_0$ ) - person is not guilty

- The assumption you are beginning with

② Alternative hypothesis ( $H_1$ ) - The person is guilty

- Opposite of Null hypothesis

③ Experiments → Statistical Analysis

→ Collect proof (DNA, finger print)

④ Accept the Null hypothesis or Reject the Null hypothesis

Ex:- Colleges at District A states its average passed percentage of students is 85%. A new college opened in the district and it was found that a sample of student has a pass percentage of 90% with a standard deviation of 4%. Does college have a different passed percentage?

Null hypothesis ( $H_0$ ) =  $\mu = 85\%$

Alternate hypothesis ( $H_1$ ) =  $\mu \neq 85\%$

## Hypothesis And Statistical Analysis:-

① Z-test

}  $\Rightarrow$  data dealing with average

Ex:- sample average

Z-table

② t-test

t-table

③ CHI SQUARE  $\Rightarrow$  Statistical analysis of categorical data

④ ANNOVA  $\Rightarrow$  Variance of data

Z-test :-

① When to use Z-test?

If i) population std & ii)  $n \geq 30$   
is known

② The average heights of all residents in a city is 168 cm with a  $\sigma = 3.9$ . A doctor believes the mean to be different. He measured the height of 36 individuals & found the average height to be 169.5 cm.

③ State null & alternate hypothesis

④ At a 95% confidence level, is there enough evidence to reject the null hypothesis

Ans

$$\mu = 168 \text{ cm}$$

$$\sigma = 3.9$$

$$n = 30$$

$$\bar{x} = 169.5$$

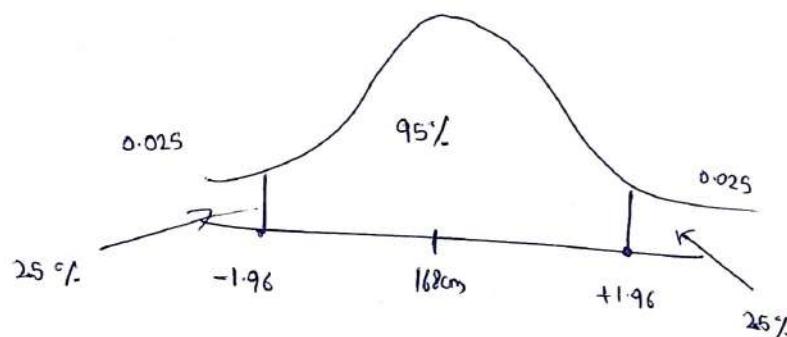
$$CI = 0.95$$

$$\alpha = 1 - 0.95 = 0.05 //$$

① Null hypothesis ( $H_0$ ) =  $\mu = 168 \text{ cm}$

②  $H_1 = \mu \neq 168 \text{ cm}$

③ Based on CI we will draw the Decision boundary



$$1 - 0.025 = 0.9750$$

$$Z\text{-Score} = 0.9750$$

$$\text{Area} \Rightarrow +1.96$$

If  $Z$  is less than  $-1.96$  or greater than  $+1.96$ , reject the Null hypothesis.

$Z$ -test

$$Z_K = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{169.5 - 168}{3.9/\sqrt{36}}$$

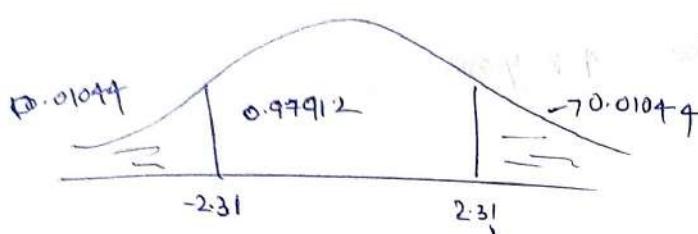
$$Z\text{-Score} = \frac{\bar{x} - \mu}{\sigma}$$

$$Z_d = \frac{1.5}{0.65} = \boxed{2.31}$$

Conclusion:

$2.31 > 1.96$  Reject the null hypothesis

$$P < 0.05$$



$$z\text{-table} \rightarrow 0.98956$$

$$1 - 0.98956 = 0.01044$$

$$\Rightarrow P = 0.01044 \times 2$$

$$1 - 0.02088$$

$$P \text{ value} = 0.01044 + 0.01044$$

$$= 0.02088$$

$$P < 0.05$$

$0.02088 < 0.05 \Rightarrow$  Reject the null hypothesis

Final conclusion:

the average height is not equal to 168cm

The average height seems to be increasing based on sample data

② A factory manufactures bulbs with an average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 years. He tests a sample of 40 bulbs and find the average time to be 4.8 years.

③ State null & alternate hypothesis

④ At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised?

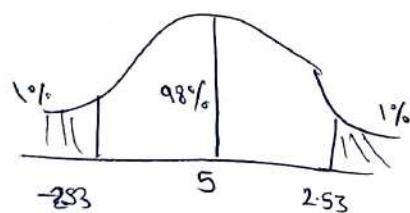
Ans

$$\mu = 5, \sigma = 0.50, n = 40, \bar{x} = 4.8$$

$$① H_0 : \mu = 5$$

$$② H_1 : \mu < 5$$

③ Decision boundary



$$\alpha = 0.02 \Rightarrow$$

$$z\text{-score} = 1 - 0.02$$

$$= 0.98$$

$\hookrightarrow$   
from

2.53

$$z\text{-test} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{4.8 - 5}{0.5/\sqrt{40}} = \frac{-0.2}{0.079} = -2.53$$

$$= -2.53 = 0.0570$$

Area under curve with z-score = 0.0570  
P-value = 0.0570

Compare p-value with  $\alpha$

$$0.0570 > 0.02 \Rightarrow \text{fail to reject}$$

We accept the Null hypothesis

We fail to reject Null hypothesis

## Student's t distribution:

In Z stats when we perform any analysis using Z-score

we require σ (population standard deviation) → is already known

→ How do we perform any analysis when we don't know the population standard deviation?

## Student's t distribution:

$$z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}}$$

z-test

t-table

s = sample standard deviation

$$t = \frac{\bar{x} - \mu}{s / \sqrt{n}}$$

t-table ⇒ t-test

## Degree of freedom

$$df = n - 1 = 3 - 1 = 2$$

3 people

T-stats : T-test → one sample test

- ① In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence?

Ans

$$\mu = 100 \quad n = 30 \quad \bar{x} = 140 \quad s = 20 \quad C.I = 95\% \quad \alpha = 0.05$$

(i) Null hypothesis  $H_0: \mu = 100$

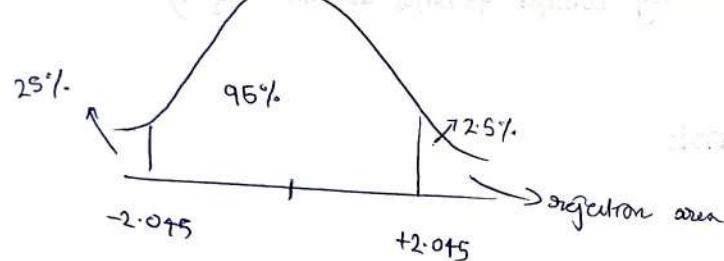
(ii) Alternate hypothesis  $H_1: \mu \neq 100$  {2 tail Test}

②  $\alpha = 0.05$

③ Degree of freedom :-

$$dof = n - 1 = 30 - 1 = 29$$

④ Decision Rule



if t test is less than  $-2.045$  or greater than  $2.045$ , Reject the null hypothesis

⑤ Calculating test statistic

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.45} = 10.96$$

$$t = 10.96$$

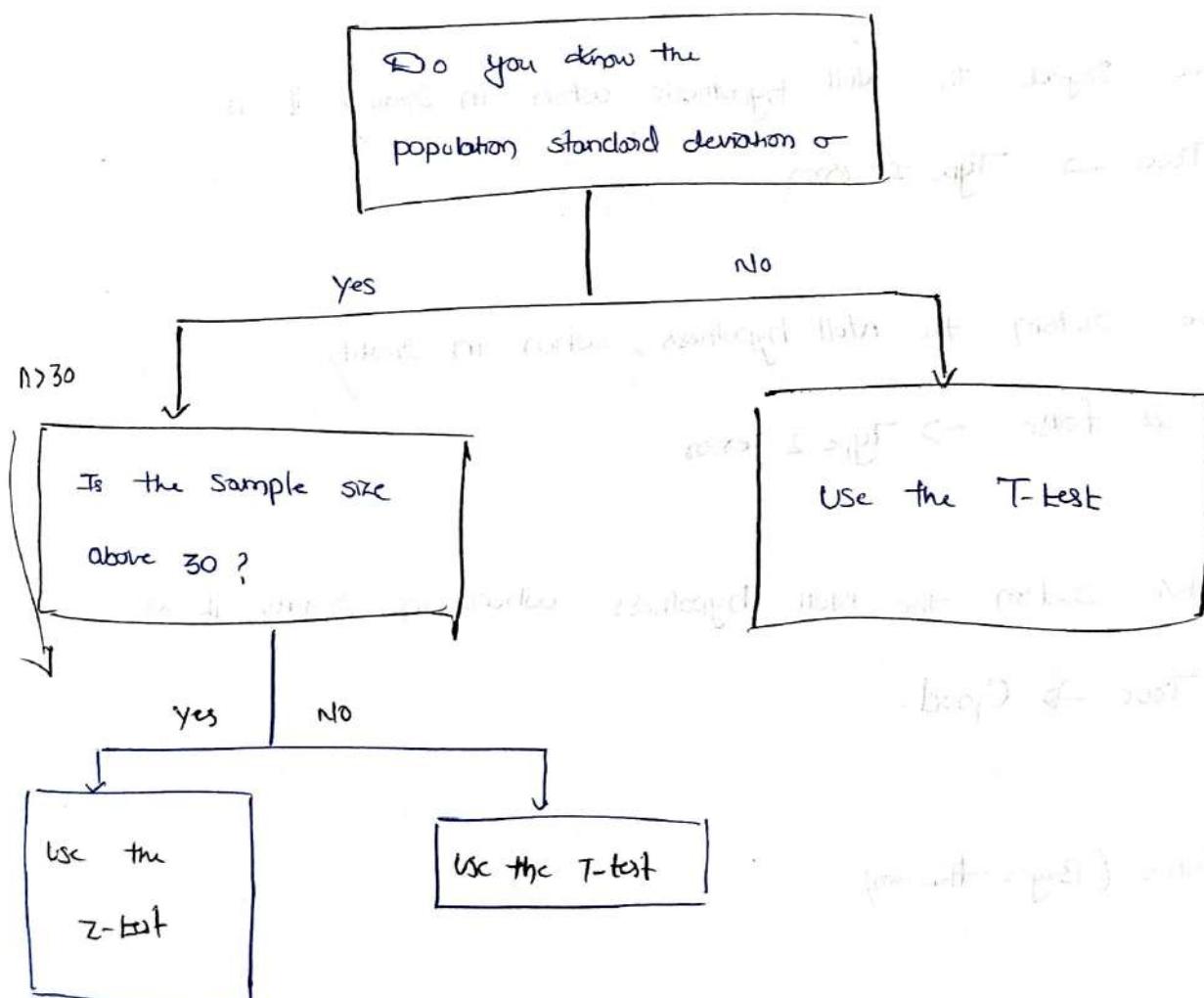
Since

$$t = 10.96 > 2.075 \quad \{ \text{reject the null hypothesis}$$

Conclusion: Medication used has affected the intelligence

Medication has increased

### When To Use T-test vs Z-test



## Type 1 and Type 2 Errors

Reality :- Null hypothesis is True or Null hypothesis is False

Decision :- Null hypothesis is True or Null hypothesis is False

Outcome 1 :- We reject the Null hypothesis when in Reality it is false  $\rightarrow$  Good

Outcome 2 :- We reject the Null hypothesis when in Reality it is True  $\rightarrow$  Type I error

Outcome 3 :- We retain the Null hypothesis, when in Reality it is false  $\rightarrow$  Type II error

Outcome 4 :- We retain the Null hypothesis when in Reality it is True  $\rightarrow$  Good.

## Baye's Statistics (Bayes Theorem)

Bayesian statistic is an approach to data analysis & parameter estimation based on Baye's theorem

Probability - [ Independent Events  $\rightarrow$  Rolling a dice  
 Dependent Events  
 ↓  
 [ Red ]  $P(R) = \frac{3}{5} \rightarrow \frac{2}{4}$   
 yellow  
 [ Yellow ]  
 Dials ]

$$P(R \text{ and } Y) = P(R) * P(Y|R) \rightarrow \text{Conditional probability}$$

$$= \frac{3}{5} * \frac{2}{4}$$

$$= \frac{6}{20}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\Rightarrow P(A \cap B) = P(Y_B) * P(B) \quad \text{--- (1)}$$

$$P(Y_A) = \frac{P(A \cap B)}{P(A)}$$

$$\Rightarrow P(A \cap B) = P(B/A) * P(A) \quad \text{--- (2)}$$

① & ②, we get

$$P(Y_B) * P(B) = P(B/A) * P(A)$$

$P(A|B)$  = probability of A given B is true

$$P(Y_B) = \frac{P(Y_A) + P(A)}{P(B)}$$

Likelihood

prior

Posterior

Marginal

$P(B/A)$  = probability of B given A is true

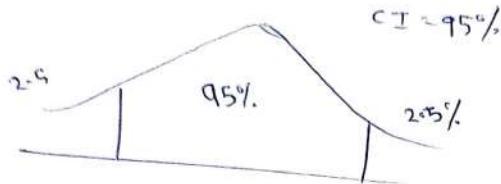
Data

Size of House	No of Rooms	Location	price
------------------	----------------	----------	-------

$$P(Y/x_1 x_2 x_3) = \frac{P(y) * p(x_1 x_2 x_3 | y)}{P(x_1 x_2 x_3)}$$

Bayes' theorem

## Confidence Intervals and margin of error



Point Estimate

$$\begin{array}{c} \boxed{\bar{x}} \longrightarrow \boxed{\mu} \\ \bar{x} = 2.5 \\ N = 3 \end{array}$$

$\bar{x}$  can less or greater than 3, this is a problem of point estimation

We should define a range of confidence Interval

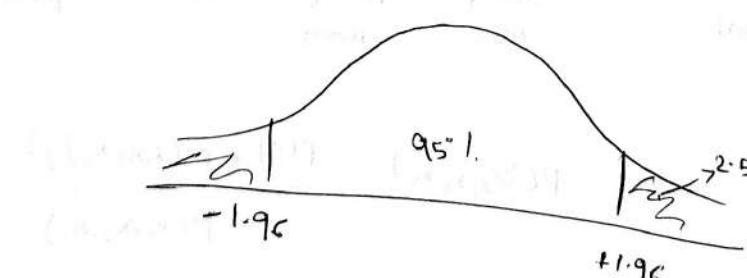
like PE lies in 2-4

Confidence Interval

point estimate  $\pm$  Margin of error

$$\boxed{\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}}} \rightarrow Z\text{-test}$$

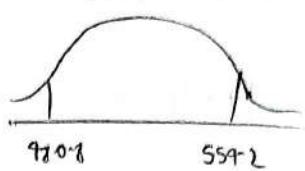
$$\boxed{\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}} \rightarrow t\text{-test}$$



$$1 - 0.95 = 0.9750$$

$+1.96 \Rightarrow 5\text{ table}$

$$\text{Lower CI} = 520 - (1.96) \cdot \left( \frac{100}{\sqrt{25}} \right) = 480.8$$



$$\text{Higher CI} = 520 + (1.96) \cdot \frac{100}{\sqrt{25}} = 559.2$$

Conclusion: I am 95% confident that about the mean CAT score is between 480.8 to 559.2.

### (HI-SQUARE ( $\chi^2$ ) Test):

When a fair coin is tossed 100 times, the theoretical considerations leads us to expect 50 heads & 50 tails. But in practice, these results are rarely achieved.

i.e. The result obtained in a experiment do not agree exactly with the theoretical results.

→ The magnitude of discrepancy between the theory and observation is given by the quantity  $\chi^2$  (a greek letter, pronounced as "chi-square").

→ If  $\chi^2=0$ , the observed & expected frequencies completely coincide.

→ As the value  $\chi^2$  increases, it affords as a measure of the corresponding discrepancy between the observed & theoretical frequencies increases.

→ Thus,  $\chi^2$  affords a measure of the correspondence between theory & observation.

Definition: If a set of events  $A_1, A_2, A_3, \dots, A_n$  are observed to occur with frequencies  $o_1, o_2, o_3, \dots, o_n$  respectively and according to probability rules  $A_1, A_2, \dots, A_n$  are expected to occur with frequencies  $E_1, E_2, \dots, E_n$  respectively with  $o_1, o_2, \dots, o_n$  are called observed frequencies and  $E_1, E_2, \dots, E_n$  respectively with  $o_1, o_2, \dots, o_n$  are called expected frequencies.

Def If  $o_i$  ( $i=1, 2, \dots, n$ ) is set of observed (experimental) frequencies &  $E_i$  ( $i=1, 2, \dots, n$ ) is a set of ~~observed~~ (~~experimental~~ theoretical) frequencies, corresponding expected

then  $\chi^2$  is defined as  $\chi^2 = \sum_{i=1}^n \frac{(o_i - E_i)^2}{E_i}$  with  $(n-1)$  degrees of freedom

It is used to test whether the differences b/w of  $\chi^2$  are significant

Note: If the data is given in a series of ' $n$ '-numbers the degrees of freedom  $= n-1$

In case of Binomial distribution, d.f.  $= n-1$

Poisson distribution, d.f.  $= n-2$

Normal distribution, d.f.  $= n-3$

$$\chi^2 = \sum_{i=1}^n \frac{(o_i - E_i)^2}{E_i}$$

$$= \frac{(o_1 - E_1)^2}{E_1} + \frac{(o_2 - E_2)^2}{E_2} + \dots + \frac{(o_n - E_n)^2}{E_n}$$

CHI Square for Goodness of Fit

In 2010 Census of the city, the weight of the individuals in a small city were found to be the following.

$<50\text{kg}$	$50-75$	$>75$
20%	30%	50%

In 2020, weight of  $n=500$  individual were sampled. Below are the result.

$<50$	$50-75$	$>75$
140	160	200

Using  $\alpha=0.05$ , would you conclude the population differences of weights has change in the last 10 years?

Ans

2010	$<50\text{kg}$	$50-75$	$>75$
Expected	20%	30%	50%

2020 observed $n=500$	$<50$	$50-75$	$>75$
	140	160	200

Expected	$<50$	$50-75$	$>75$
	$0.2 \times 500$	$0.3 \times 500$	$0.5 \times 500$
	=100	=150	=250

① Null hypothesis  $H_0$ : The data meets the expectation

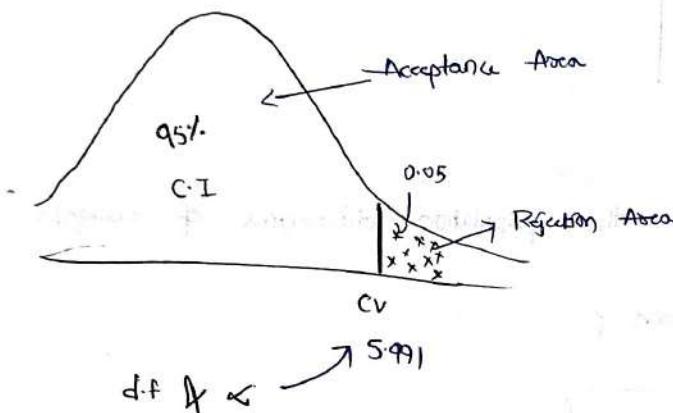
② Alternate hypothesis  $H_1$ : The data does not meet the expectation

③  $\alpha = 0.05$  C.I = 95%

④ Degrees of freedom

$$d.f = k-1 = 3-1 = 2$$

⑤ Decision boundary



If  $\chi^2$  is greater than 5.99, reject  $H_0$  else

We fail to reject the Null hypothesis

⑥ Calculation:

$$\chi^2 = \sum \frac{(O-E)^2}{E}$$

$$= \frac{(140-100)^2}{100} + \frac{(160-150)^2}{150} + \frac{(200-250)^2}{250}$$

$$= 16 + 0.66 + 10 \quad 26.66 > 5.99, \text{ Reject } H_0$$

$$\chi^2 = 26.66$$

The weights of 2020 populations are different than those experiment in 2010 population

## Analysis of variance (ANOVA)

Super imp statistical method

Definition: ANOVA is a statistical method used to compare the means of 2 or more groups.

### ANOVA

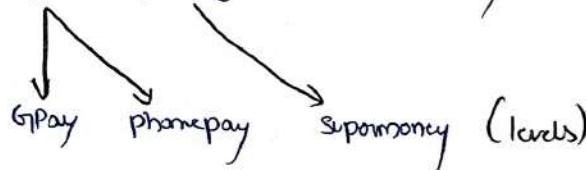
- ① Factors (variable)

- ② levels

Eg:- Medicine (Factor)

[Dosage] 5mg 10mg 15mg → level

Eg:- Mode of payment (Factors)



### Assumptions of ANOVA (imp of interview)

→ centre limit theorem

- ① Normality of Sampling Distribution of mean

The distribution of Sample mean is normally distributed

- ② Absence of outliers

outlying score need to be removed from the dataset

- ③ Homogeneity of variance [ $\sigma_1^2 = \sigma_2^2 = \sigma_3^2$ ]

Population variance in different levels of each independent variables are equal

- ④ Samples are independent & randomly selected

## Types of ANNOVA (3 types)

① One Way ANNOVA :- one factor with atleast 2 levels, these levels are independent

Eg:- Doctor wants to test a new medication to decrease headache  
They split the participants in 3 conditions [10mg, 20mg, 30mg] Doctor ask the participants to state the headache [1-10]

Medication  $\rightarrow$  Factor

10mg	20mg	30mg
5	7	2
3	9	1
-	-	-
-	-	-

(subject) running for study  
(subject) participant running for study

② Repeated Measures ANNOVA :- one factor with atleast 2 levels, levels are dependent

Running  $\rightarrow$  Factor

levels	Day 1	Day 2	Days
8	5	7	
7	4	9	
-	-	-	

③ Factorial ANOVA: Two or more factors (each of which with at least) 2 levels, levels can be independent and dependent.

		Running → Factors		
		Day 1	Day 2	Day 3
Levels				
Gender factor	male	8	5	4
	female	9	7	3
		2	4	6
		7	8	3

Hypothesis Testing In ANOVA (Partitioning of variance in the Anova)

Null Hypothesis:  $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_n$ , all population means are equal.

Alternate hypothesis:  $H_1$ : At least one of the sample mean is now unequal.

$$\boxed{\mu_1 \neq \mu_2 \neq \mu_3 \neq \dots \neq \mu_n} \rightarrow \text{wrong}$$

Test Statistics

$$F = \frac{\text{Variance between Sample}}{\text{Variance within Sample}}$$

Variance between sample			
	$x_1$	$x_2$	$x_3$
Variance within Samples	1	4	2
	5	6	3
	3	7	4
	9	8	8
	5	9	6

$$\bar{x}_1 = 3 \quad \bar{x}_2 = 19/5 \quad \bar{x}_3 = 4$$

### Ex:- One way ANOVA

One factor with atleast 2 levels, levels are independent

- ① Doctors want to test a medication which reduces headache. They splits the patient into 3 condition [15mg, 30mg, 45mg]. Later on the doctor ask the patient to rate the headache between [1-10]. Are there any difference biffereance between the 3 condition using  $\alpha = 0.05$ ?

	15mg	30mg	45mg
9	7	4	
8	6	3	
7	6	2	
8	7	3	
2	3	7	
9	7	3	
8	6	2	

① Define null & alternate hypothesis?

$$H_0 : \mu_{15} = \mu_{30} = \mu_{45}$$

$H_1$  : at least one mean is not equal

② Significance

$$\alpha = 0.05, C.I = 1 - 0.05 = 0.95$$

③ Calculate Degree of freedom

$$N = 21 \quad a = 3 \quad n = 7$$

↓  
Categories

$$df_{\text{between}} = a - 1 = 3 - 1 = 2$$

$$\left. \begin{array}{l} df_1, df_2 \\ (2, 18) \end{array} \right\}$$

$$df_{\text{within}} = N - a = 21 - 3 = 18$$

↓ with  $\alpha = 0.05$

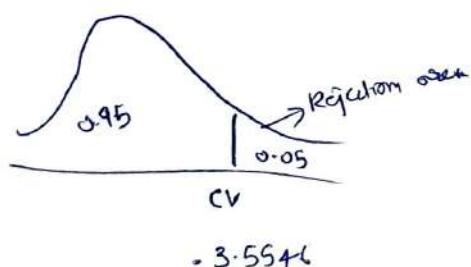
F-table

↓ finds

Critical Value

$$df_{\text{total}} = N - 1 = 20$$

④ Decision Boundary



Decision rule

If  $F$  is greater than  $-3.5546$ , reject the NULL Hypothesis

⑤ Calculate F test statistic

	Sum of Square	df	mean Square	F
Between between within	98.67	2	49.34	
	10.29	18	0.57	
Total	108.96	20		

$$\textcircled{1} \quad SS_{\text{between}} = \sum \frac{(\bar{x}_{ai})^2}{n} - \frac{\bar{T}^2}{N}$$

$$15 \text{mg} = 4+8+7+8+7+5+9+8 = 57$$

$$30 \text{mg} = 7+6+7+8+7+4 = 47$$

$$45 \text{mg} = 7+3+2+3+4+2+2 = 21$$

$$= \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57^2 + 47^2 + 21^2]}{21}$$

$$= 98.67$$

$$\textcircled{2} \quad SS_{\text{within}} = \sum y_i^2 - \sum (\bar{x}_{ai})^2$$

$$\sum y_i^2 = 9^2 + 8^2 + 7^2 + 8^2 + 1^2 - \dots$$

$$= 853$$

$$= 853 - \frac{(57^2 + 47^2 + 21^2)}{21}$$

$$\approx 10.29$$

$$F\text{-test} = \frac{MS_{\text{Between}}}{MS_{\text{Within}}}$$

$$F\text{-test}_{(1)} = \frac{99.34}{0.54}$$

$$= 86.56$$

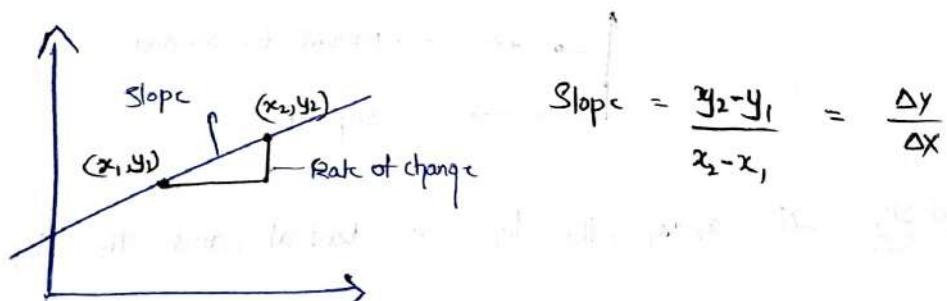
If  $f$  is greater than  $3.5546$ , Reject the  $H_0$ .

$86.56 > 3.5546$  Reject the  $H_0$ .

Slope → Derivative as a concept

The Slope of a line is a measure of how steep the line is, and it represents the rate of change of one variable with respect to another.

In the context of a two-dimensional cartesian co-ordinate System, the slope indicates the ratio of the vertical change (rise) to the horizontal change (run) between two points on a line.



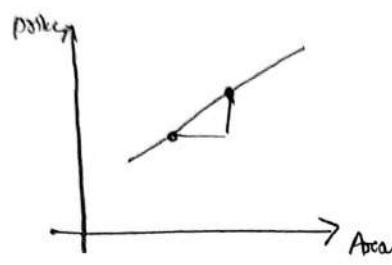
where  $y_2 - y_1$  is the vertical (rise)

$x_2 - x_1$  is the horizontal (run)

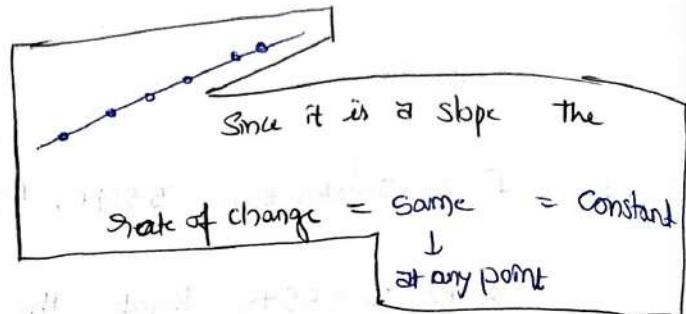
Example Data Set:



2 features : Area & force



as area increase force decreases



Interpretation of Slope:

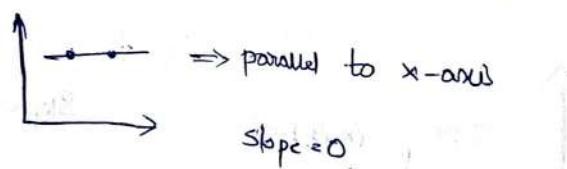
\* Positive slope: If slope  $> 0$ , the line rises as it moves left to right, the larger the slope, the steeper the line.



\* Negative slope: If slope  $< 0$ , the line falls as it moves from left to right, the more negative the slope, the steeper the line in downwards direction.

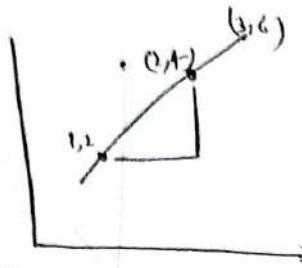


\* Zero slope: If slope  $= 0$  the line is horizontal, meaning there is no vertical change as the line moves from left to right.



\* undefined slope: If  $x_2 = x_1$ , the line is vertical, and the slope is undefined because you cannot divide by 0.

Example 1: Positive Slope:



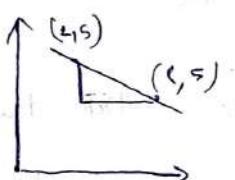
$$\text{Slope} = \frac{6-2}{3-1} = \frac{4}{2} = 2$$

This means for every 1 unit you move horizontally from left to right the line moves 2 units vertically up

Example 2: Negative slope:

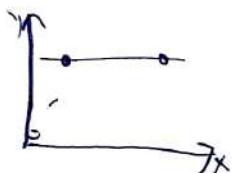
(2, 5) (6, 3)

$$\text{Slope} = \frac{3-5}{6-2} = \frac{-2}{4} = -\frac{1}{2}$$



This means for every 2 units you move horizontally to the right, the line moves 1 unit vertically down

Example 3: Zero slope: Consider the point (4, 4), (5, 4)



$$\text{Slope} = \frac{4-4}{5-4} = \frac{0}{1} = 0 \Rightarrow \frac{0}{4} = 0$$

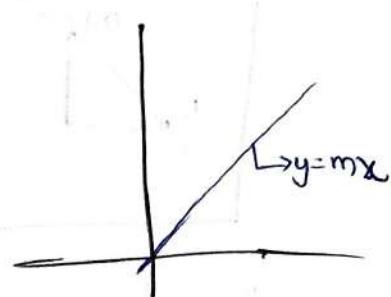
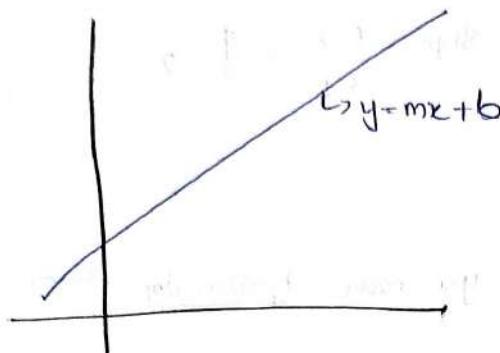
This means the line is horizontal

Example 4: Undefined slope

(3, 2) (3, 7)

$$\text{Slope} = \frac{7-2}{3-3} = \frac{5}{0}$$

## Slope in a equation of a line



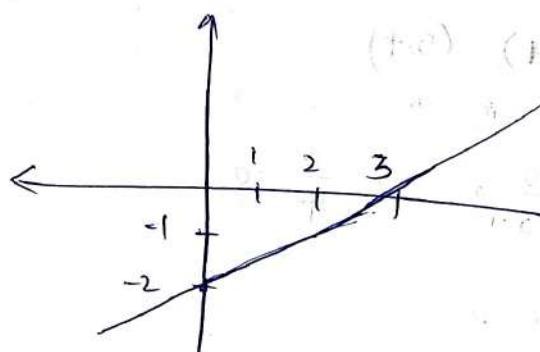
$$y = mx + b$$

$m$  is the slope of the line

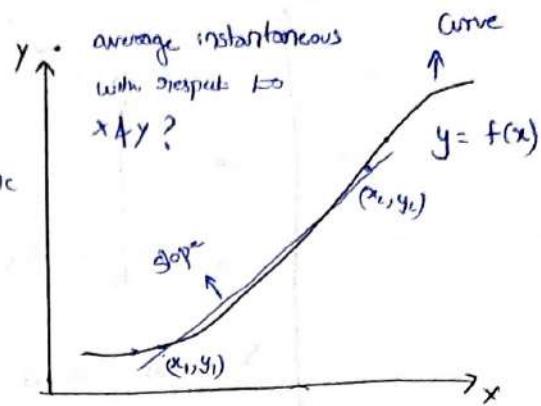
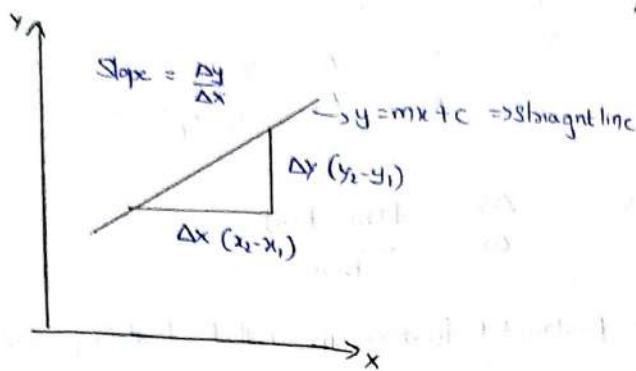
$b$  is  $\approx$  the y intercept where the line crosses  $x$  axis

$$\boxed{m=3} \quad \boxed{b=-2}$$

$$\Rightarrow y = 3x - 2$$

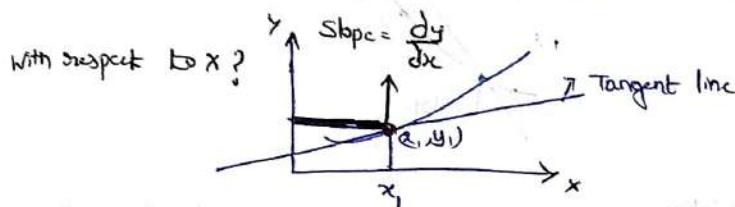


Derivative:



In this figure the Slope is change as it is a curve

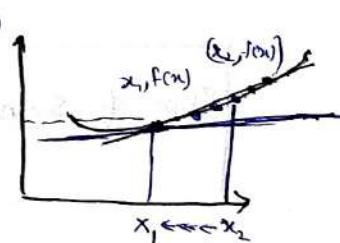
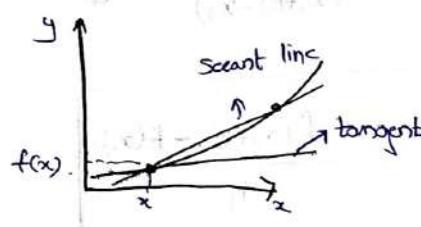
⇒ To get average instantaneous we draw a Secant line ↗



Rate of change of  $x_1$  with respect to  $y$ ? cannot draw secant ↗ so we calculate by drawing tangent

Slope of Tangent = Instantaneous rate of change of  $x$  with respect to  $y$

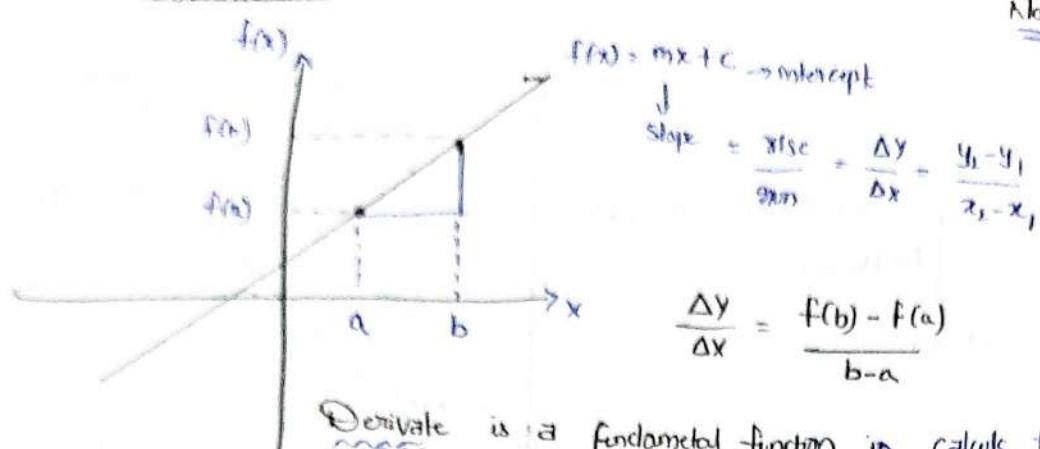
$$\Rightarrow \text{slope} = \frac{dy}{dx}$$



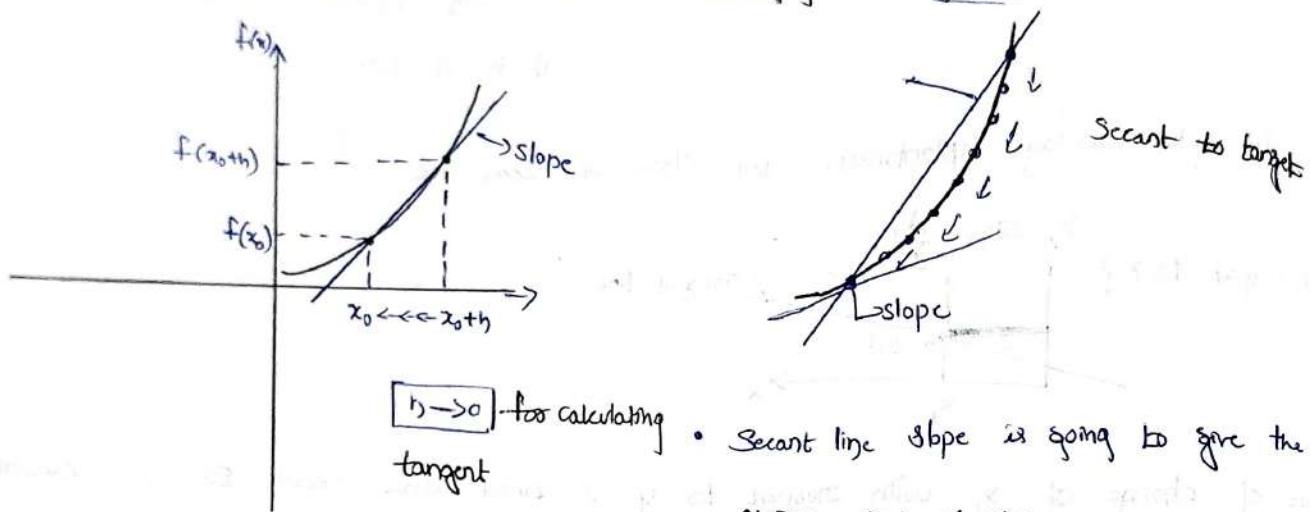
Mathematics notation of Derivative with limits ↗

Secant to Tangent then slope ↗

The derivative is a fundamental concept in calculus that represents the rate of change of a function in changing at any given point. It is essentially the slope of the tangent line to the function's graph at that point. The derivatives is used to understand how a function behaves as its input changes & it is a key tool for analyzing the dynamics of systems in mathematics, physics, economics, engineering & many other fields.



Derivative is a fundamental function in calculus that represents the rate of which a function is changing at a given point.



- Secant line slope is going to give the average rate of change

- To get exact slope at a point you need to draw a tangent line & then find a slope

$$\text{Slope} = \frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{f(x_0+h) - f(x_0)}{(x_0+h) - x_0} = \frac{\Delta y}{\Delta x}$$

Slope of Secant line =  $\frac{f(x_0+h) - f(x_0)}{h}$

To calculate rate of change at specific point we need to convert Secant line to tangent by moving from  $x_0+h$  to  $x_0$  [we use limit to move]

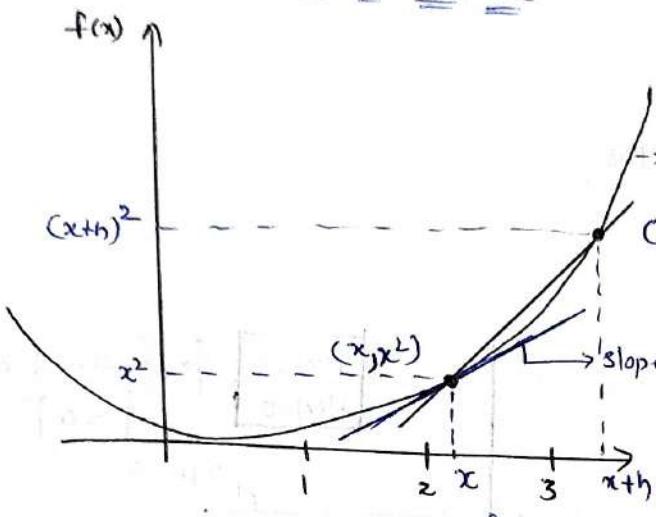
$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x_0+h) - f(x_0)}{h}$$

$\leftarrow$  derivative of  $f(x)$

$$f'(x) = \frac{dy}{dx} = \frac{d(f(x))}{dx}$$

# Finding = Derivative at a point :

\* Proving derivative of  $x^2$  is  $2x$



$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$\Rightarrow f(1) = 1$$

$$f(2) = 4$$

$$= 9$$

$\vdots$

$n^2$

limit  
 $h \rightarrow 0$

for calculating slope at target

$$f'(x) = \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h}$$

$$f'(x) = \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - x^2}{h}$$

$$= \lim_{h \rightarrow 0} \frac{h(2x+h)}{h}$$

$$= 2x + 0$$

$$f'(x) = 2x \quad \rightarrow \text{find a derivative at a point } x$$

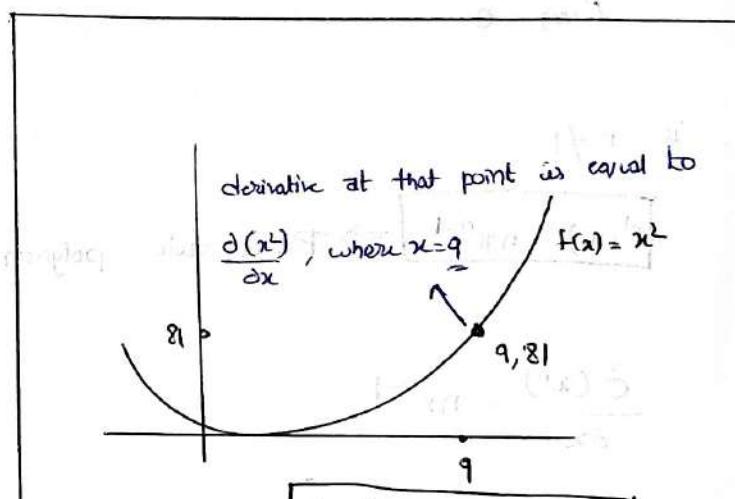
where  $f(x) = x^2$

Hence proved that  $\frac{\partial (f(x))}{\partial x} \Rightarrow \boxed{\frac{\partial (x^2)}{\partial x} = 2x}$

$$f'(x^2) = 2x =$$

$$\boxed{\frac{\partial (x^n)}{\partial x} = nx^{n-1}}$$

$$\Rightarrow \frac{\partial (x^2)}{\partial x} = 2x^{2-1} = 2x =$$



$$\boxed{\frac{\partial (x^2)}{\partial x} = 2x = 18}$$

# Power Rule In Derivatives

→ JI as applied to polynomials

$$f(x) = x^2$$

$$f'(x) = 2x$$

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Different polynomial equations have different graphs  
 ↓

$$x^2+3, x^3, x^4+2x+1 \dots \text{etc}$$

$$f(x) = x^n, n \neq 0$$

If  $n=0$   $f(x) = x^0 = 1 \Rightarrow$  constant value

$$f'(x) = 0$$

4  $n \neq 1$

$$f'(x) = nx^{n-1} \Rightarrow \text{Power rule polynomial expression.}$$

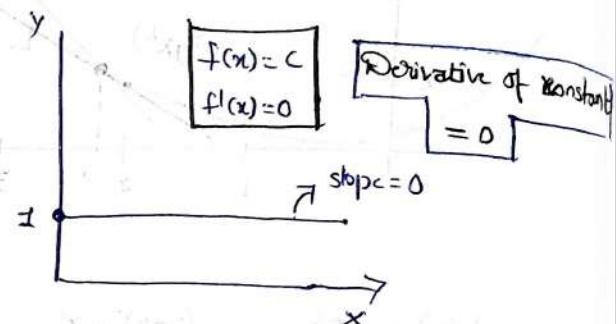
$$\frac{\partial (x^n)}{\partial x} = nx^{n-1}$$

$$\Rightarrow f(x) = x^3, x=2$$

$$\frac{\partial (x^3)}{\partial x} = 3x^2 = 3 \cdot 2^2 = 12$$

$$\Rightarrow \frac{\partial (3x^2)}{\partial x} = 3 \cdot \frac{\partial x^2}{\partial x} = 3 \cdot 2x = 6x$$

$$\Rightarrow \frac{\partial (\frac{1}{x})}{\partial x} = \frac{\partial (x^{-1})}{\partial x} = -x^{-2} = -\frac{1}{x^2}$$



Assignment :-

$$\textcircled{1} \ f(x) = x^8 \quad \underline{\text{Ans}} \quad f'(x) = \frac{d x^8}{d x} = 8x^{8-1} = 8x^7$$

$$\textcircled{2} \ f(x) = x^{-2} \quad \Rightarrow \ f'(x) \text{ at } x = \underline{-1}$$

$$\underline{\text{Ans}} \quad \frac{d x^{-2}}{d x} = -\frac{1}{x^2} = \frac{-1}{(-1)^2} = -1$$

\* Derivative Rules :- Constant, sum, difference and Constant Multiple

$$\frac{d(x^n)}{dx} = nx^{n-1}, n \neq 0 \quad \left\{ \text{power rule} \right\} \Rightarrow \text{polynomial}$$

$$\frac{d(x^0)}{dx} = \frac{d(1)}{dx} = 0 \quad \Rightarrow \text{Derivation of a constant is 0}$$

$$\frac{d(c)}{dx} = 0 \quad \begin{matrix} \nearrow \text{constant} \\ c \end{matrix}$$

$$\Rightarrow \frac{d(cf(x))}{dx} = c \frac{d(f(x))}{dx} = c f'(x)$$

$$\frac{d(3x^4)}{dx} = 3 \frac{d(x^4)}{dx} = 3 \times 4x^{4-1} = 12x^3$$

$$\text{at } x=2 \quad \Rightarrow 12 \times 2^3 = 12 \times 8 = 96$$

Assignment

$$\frac{d[2f(x)]}{dx} = \frac{d[2x^5]}{dx} = 2 \frac{d(x^5)}{dx} = 2 \times 5x^{5-1} = 10x^4$$

\* Sum of 2 functions  $f(x), g(x)$

$$\frac{\partial [f(x) + g(x)]}{\partial x} = \frac{\partial (f(x))}{\partial x} + \frac{\partial (g(x))}{\partial x}$$

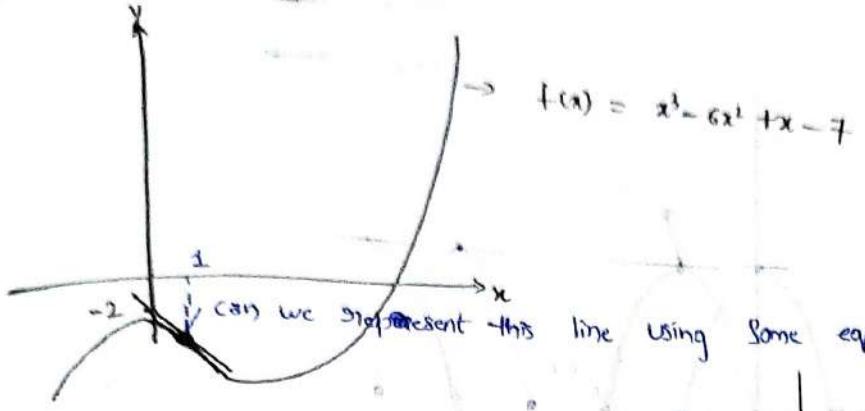
$$\begin{aligned}\frac{\partial [x^4 + x^{-2}]}{\partial x} &= \frac{\partial (x^4)}{\partial x} + \frac{\partial (x^{-2})}{\partial x} \\ &= 4x^3 - 2x^{-3}\end{aligned}$$

Assignment :-

$$\begin{aligned}① \frac{\partial (x^2+2)}{\partial x} &= \frac{\partial (x^2)}{\partial x} + \frac{\partial (2)}{\partial x} \\ &= 2x\end{aligned}$$

$$\begin{aligned}② \frac{\partial [4x^3 - 6x^2 + 2x + 100]}{\partial x} &= \frac{\partial [4x^3]}{\partial x} - \frac{\partial [6x^2]}{\partial x} + 2 \frac{\partial (x)}{\partial x} + \frac{\partial (100)}{\partial x} \\ &= 12x^2 - 12x + 2 + 0 \\ &= 12x^2 - 12x + 2\end{aligned}$$

## Yield of polynomial



$$f'(x) = \frac{d(x^3)}{dx} - \frac{d(6x^2)}{dx} + \frac{d(7x)}{dx} - 7$$

$$= 3x^2 - 12x + 1 - 0$$

$$= 3x^2 - 12x + 1$$

$$\text{at } x=1$$

$$= 1 - 12 + 1 - 7 = -11 //$$

$$f'(1) = 3 - 12 + 1 = -9 \Rightarrow \text{slope}$$

$$y = mx + c$$

$$y = -9x + c$$

$$-11 = -9x + c$$

$$c = -9(1) - 11$$

$$\boxed{c = -2}$$

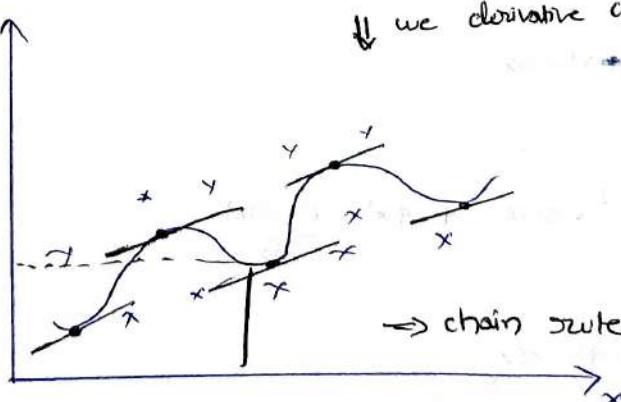
$$\Rightarrow y = mx + c$$

$$\boxed{y = -9x - 2}$$

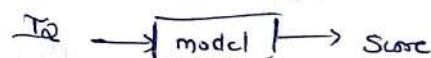
↓      ↓  
 Slope    intercept

*(x) to look at that  
 equation that particular target at that  
 particular point at x=1*

y



Ex:	ID	Dataset	ID	Score
			100	91
			180	89
			90	90
			85	88



we try to find right slope & right intercept

we will satisfied this line such that  
for any data point model should predict the  
output

*(x) to look at that  
 equation that particular target at that*

*particular point at x=1*

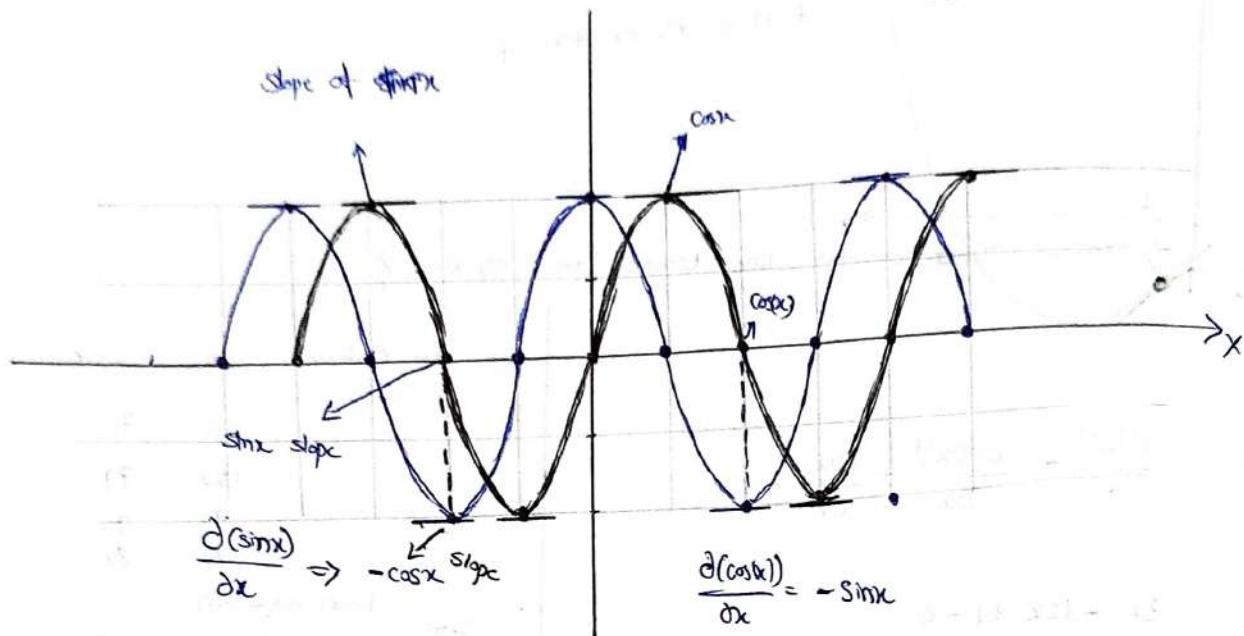
if we derivative concepts to get slope.

⇒ with respect any value predict  
the output

# Trigonometric Functions

$\sin x$

$\cos x$



Logarithmic function:

Exponential function

constant

power rule

$$f(x) = \ln(x) \text{ then}$$

$$f'(x) = \frac{1}{x}$$

$$f(x) = e^x \quad \text{if } f(x) = c$$

$$f'(x) = e^x$$

$$\text{then } f'(x) = 0$$

$$\text{if } f(x) = x^n$$

$$f'(x) = nx^{n-1}$$

## Product rule in Derivatives:-

$$\boxed{\frac{d [h(x) \cdot f(x)]}{dx} = h'(x) \cdot f(x) + h(x) \cdot f'(x)}$$

$$\textcircled{1} \quad \frac{d(x^2 \cos x)}{dx} = \frac{d(x^2)}{dx} \cdot \cos x + x^2 \cdot \frac{d(\cos x)}{dx}$$

by product rule

$$= 2x \cdot \cos x + x^2 \cdot (-\sin x)$$

Assignment

$$\textcircled{2} \quad \frac{d(4x^2 \sin x)}{dx} = \frac{d(4x^2)}{dx} \cdot \sin x + 4x^2 \cdot \frac{d(\sin x)}{dx}$$

$$= 8x \sin x - 4x^2 \cos x$$

## Chain Rule :- [of Derivation]

The Chain Rule is a fundamental theorem in calculus that is used to find the derivatives of a Composite function. When a function is composed of other function, the chain rule allows us to differentiate it with respect to the innermost variable.

Formal definition:

If  $y = f(g(x))$  where  $y = f(u)$  &  $u = g(x)$ , then the derivative of  $y$  with respect to  $x$

$$\frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

In simpler term, this can be expressed as

$$\frac{dy}{dx} = f'(g(x)) \cdot g'(x)$$

This means that to differentiate a composite function, you first differentiate the outer function with respect to its inner function and then multiply by the derivative of the inner function.

Example 1.

$$y = \underbrace{(3x^2 + 2x + 1)}_u^5$$

Step 1: Identify the outer and inner function

Outer function  $f(u) = u^5$

Inner function  $u = g(x) = 3x^2 + 2x + 1$

Step 2: Differentiate the outer function

$$f'(u) = \frac{df}{du} = \frac{d}{du} u^5 = 5u^4$$

Step 3: Differentiate the inner function

$$g'(x) = \frac{dy}{dx} = \frac{d(g(x))}{dx} = \frac{d(3x^2 + 2x + 1)}{dx} = 6x + 2$$

Step 4:  $\frac{dy}{dx} = f'(g(x)) \cdot g'(x)$

$$= 5u^4 \cdot (6x + 2)$$

$$\frac{dy}{dx} = 5(3x^2 + 2x + 1)^4 \cdot (6x + 2)$$

## ② Trigonometric

$$y = \sin(4x^3 + x)$$

$$u = g(x) = 4x^3 + x$$

$$= \sin(u)$$

$$\Rightarrow y = f(g(x))$$

Step 1: Identify the outer and inner functions.

$$\text{Outer function: } f(u) = \sin(u)$$

$$\text{inner function: } u = g(x) = 4x^3 + x$$

$$\text{Step 2: } f'(u) = \frac{\partial f}{\partial u} = -\cos(u)$$

$$\text{Step 3: } g'(x) = \frac{\partial u}{\partial x} = \frac{\partial(4x^3 + x)}{\partial x} = 12x^2 + 1$$

Step 4: Chain Rule

$$\frac{\partial u}{\partial x} = f'(g(x)) \cdot g'(x) = -\cos(4x^3 + x) \cdot (12x^2 + 1)$$

$$\boxed{\frac{\partial u}{\partial x} = -\cos(4x^3 + x) \cdot (12x^2 + 1)}$$

Example of ~~two~~ composition of three function

$$y = \sqrt{\sin(3x)}$$

$\underbrace{\hspace{1cm}}$        $\underbrace{\hspace{1cm}}$        $\left\{ \begin{array}{l} \text{3 composite function} \\ \text{outer, middle, inner} \end{array} \right.$

Step 1: Identify the function

$$\left. \begin{array}{l} \text{outer function: } f(u) = \sqrt{u} = u^{1/2} \\ \text{middle function: } g(v) = \sin(v) \\ \text{inner function: } h(x) = 3x \end{array} \right\}$$

$$u = \sin(3x)$$

$$v = 3x$$

$$\text{Step 1: } f(u) = \frac{1}{2} u^{-1/2} = \frac{1}{2\sqrt{u}}$$

$$g(v) = -\cos v$$

$$h(x) = 3$$

Step 2: Apply chain rule

$$\frac{\partial u}{\partial x} = \frac{\partial f}{\partial u} \cdot \frac{\partial u}{\partial v} \cdot \frac{\partial v}{\partial x} \quad u = \sin(3x)$$

$$= f'(g(h(x))) \cdot g'(h(x)) \cdot h'(x)$$

$$= \frac{1}{2\sqrt{\sin(3x)}} \cdot \cos(3x) \cdot 3$$

$$\boxed{\frac{\partial u}{\partial x} = \frac{3\cos(3x)}{2\sqrt{\sin(3x)}}}$$

The chain rule is a powerful tool in calculus for differentiating composite functions. By breaking down the differentiation process into manageable steps, it allows us to compute derivatives of complex expressions efficiently. Understanding and applying the chain rule is essential for solving a wide range of problems in mathematics, physics, engineering, economics, and data science.

# Applications of chain rule in Data Science

3) Backpropagation: Training Deep learning models

chain rule for calculating the gradients of loss function with respect to weight that are initialized.

4) Gradient Descent Optimization:

Linear Regression - the goal is to minimize a cost function

↳ update slopes - Gradient Descent optimization - chain rule

5) Chain rule in f.e

$$z = \log(x^2 + y^2)$$

$$x \quad z \quad \boxed{\frac{\partial z}{\partial x} = \frac{1}{x^2+y^2} \cdot 2x} \Rightarrow \text{How small change in } x \text{ & } y \text{ affect } z, \text{ chain rule}$$

$$y \quad z \quad \boxed{\frac{\partial z}{\partial y} = \frac{1}{x^2+y^2} \cdot 2y}$$

4) Regularization Techniques:

Overfitting And underfitting.

# Application of Linear Algebra, stats and Differential Calculus

My Aim:

1) Linear Algebra

2) Statistics

3) Differential Calculus



=> Application in Data Science

① Simple linear regression

② Multiple linear regression

③ Dimensionality Reduction → PCA

④ Artificial Neural Network



Deep learning

=> Linear Algebra + Differential calculus + Stats

=> Model Trained



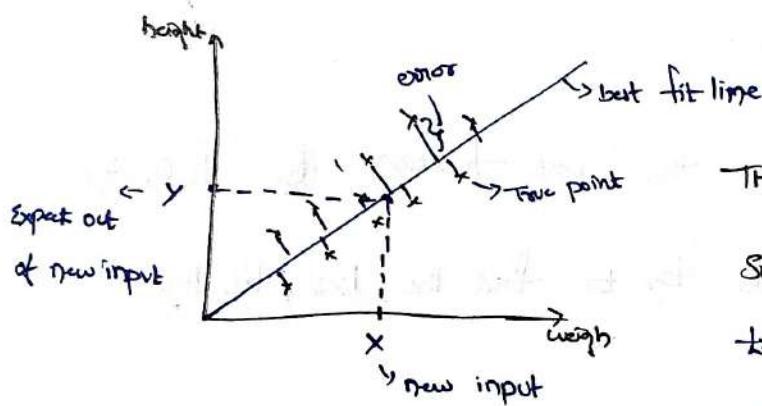
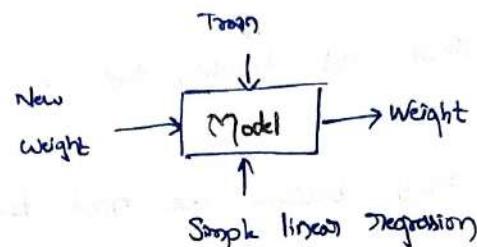
## Simple linear Regression: (Having single input/independent feature):

In Supervised ML the regression problems are solved by using Simple linear regression.

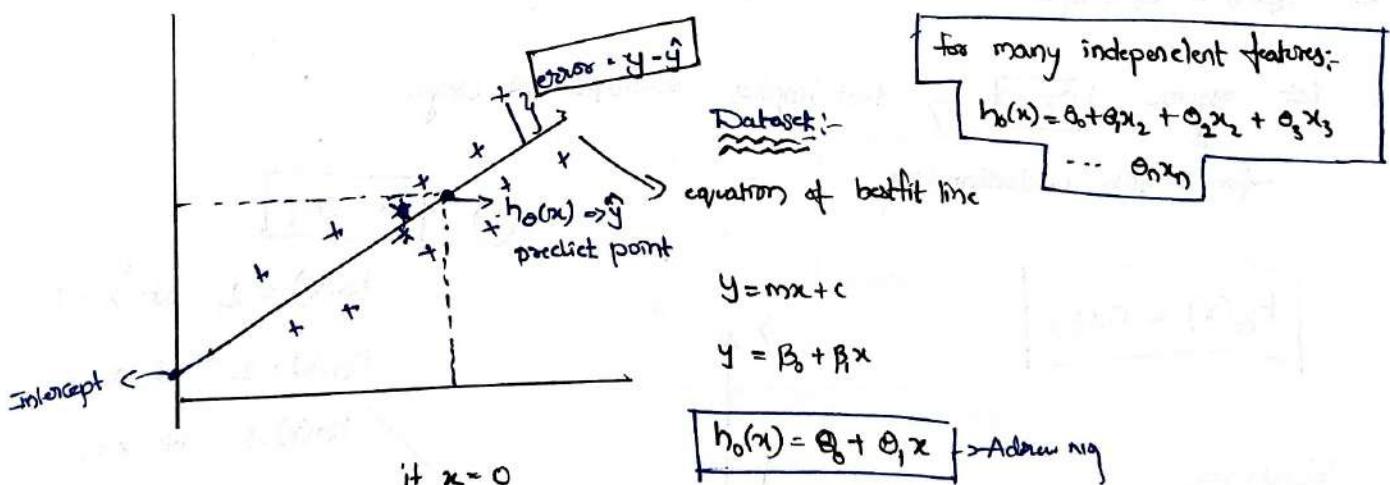
Ex. Dataset:-

	Independent Weight	Dependent or o/p Height	
74	170 cm		
80	180 cm		
75	175.5 cm		
-	-		

} True output



The best fit line should be created in such a way that the sum of all distances between the all the true point and expected point should be minimal.



$$y = mx + c$$

$$y = \theta_0 + \theta_1 x$$

$$h_0(x) = \theta_0 + \theta_1 x \rightarrow \text{Adressing}$$

$\theta_0$  = Intercept

$h_0(x) = \theta_0$  intercept

$\theta_1$  = Slope of Co-efficient { if unit movement in x axis what is the movement with respect to y-axis }

$$\text{Error} = y - \hat{y}$$

Cost function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

↑                              ↑                              ↑  
Slope                          predicted      True O/P  
↓                              ↓                              ↓  
Intercept                           Error

⇒ Mean Square Error

We also have other cost function that are mean absolute error, root MSE etc.

We are using MSE because we need to minimize it such a way that by changing the theta zero & theta one

Main aim? What we need to solve.

The main aim is to minimize the cost function  $J(\theta_0, \theta_1)$

By continuously change  $\theta_0$  &  $\theta_1$ , we try to find the best fit line

①  $h_\theta(x) = \theta_0 + \theta_1 x$

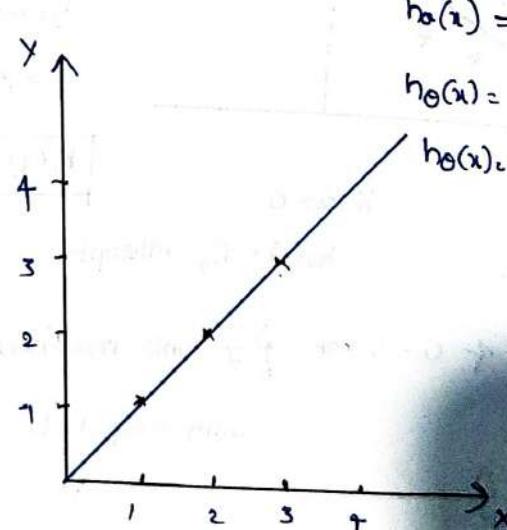
Let assume  $\boxed{\theta_0 = 0}$   $\Rightarrow$  that implies intercept at origin

for simple understanding

$$\boxed{h_\theta(x) = \theta_1 x}$$

Dataset

x	y
1	1
2	2
3	3



i)  $\boxed{\theta_1 = 1}$

$$h_\theta(x) = 1 \text{ at } x = 1$$

$$h_\theta(x) = 2 \text{ at } x = 2$$

$$h_\theta(x) = 3 \text{ at } x = 3$$

$$J(\theta_1) = \frac{1}{2 \times 3} \sum_{i=1}^3 (h_{\theta}(x^i) - y^{(i)})^2$$

$$= \frac{1}{2 \times 3} \left[ (1-1)^2 + (2-2)^2 + (3-3)^2 \right]$$

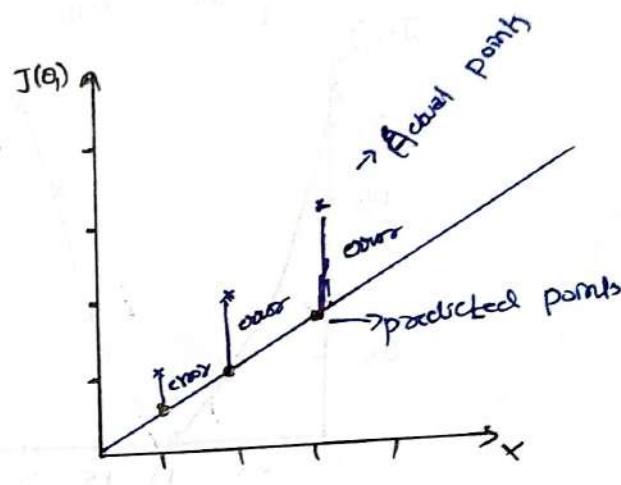
$$J(\theta_1) = 0$$

(ii) Let  $\theta_1 = 0.5$

$$x=1, h_{\theta}(x) = (0.5) \times 1 = 0.5$$

$$x=2, h_{\theta}(x) = 1$$

$$x=3, h_{\theta}(x) = 1.5$$



$$J(\theta_1) = \frac{1}{2 \times 3} ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2)$$

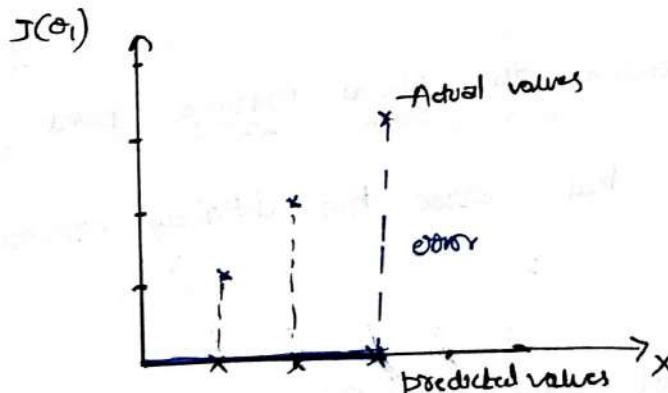
$$J(\theta_1) \approx 0.53$$

(iii) Let  $\theta_1 = 0$

$$h_{\theta}(x) = 0 \quad x=1$$

$$h_{\theta}(x) = 0 \quad x=2$$

$$h_{\theta}(x) = 0 \quad x=3$$



$$J(\theta_1) = \frac{1}{2 \times 3} [(0-1)^2 + (0-2)^2 + (0-3)^2]$$

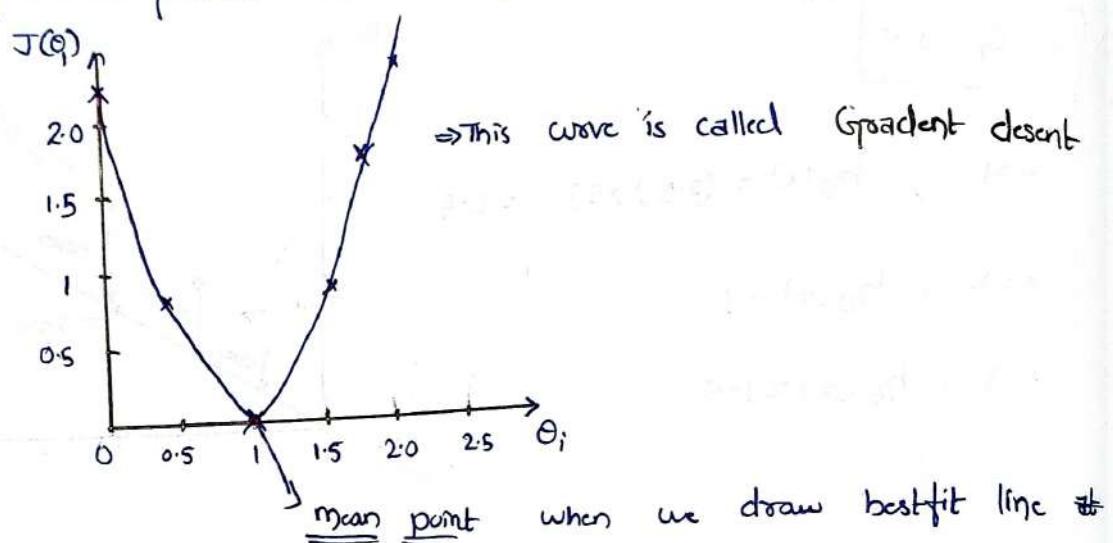
$$\approx 2.3$$

$$\Rightarrow \theta_1 = 0 \quad J(\theta_1) = 2.3$$

$$\theta_1 = 0.5 \quad J(\theta_1) = 0.58$$

$$\theta_1 = 1 \quad J(\theta_1) = 0$$

After continuing calculating  $J(\theta_1)$  value we get the following curve



using mean point error was very less. So, we can say this point is Global minimum.

$\Rightarrow$  we need achieve the global minima point or atleast near to this point to say that error has definitely minimized.

Note: We cannot keep on selecting different - different theta one ( $\theta_1$ ) values all the time. So we should apply Convergence algorithm. We should select one  $\theta_1$  value & try to find out mechanism of changing  $\theta_1$  values.

Convergence Algorithm: The main aim is to optimize the change of  $\theta_j$  value

Repeat until Convergence [means repeat until global minima]

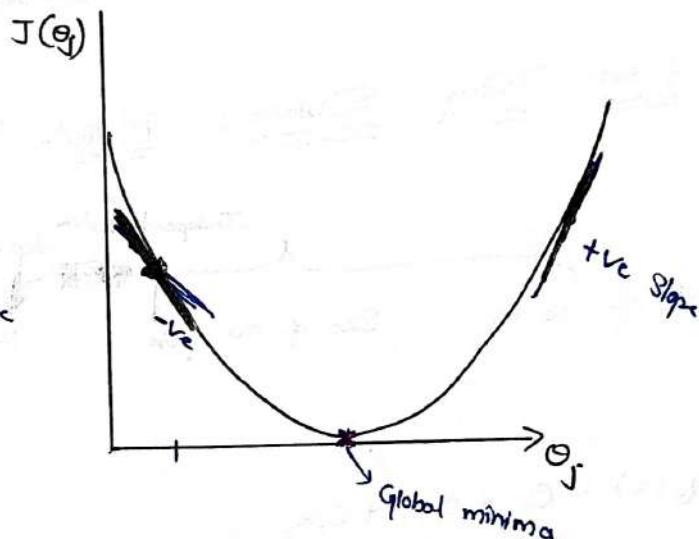
{

$$\theta_j = \theta_j - \alpha \frac{\partial J(\theta_j)}{\partial \theta_j}$$

slope

}

If Slope is -ve



$$\Rightarrow \frac{\partial J(\theta_j)}{\partial \theta_j} = -\nu_c$$

$$\Rightarrow \theta_j = \theta_j - \alpha (-\nu_c) = \theta_j + (\nu_c) \rightarrow \text{it add some value to } \theta_j$$

∴  $\theta_j$  is increased until it get at or near global minima

$$\text{If Slope is +ve} \Rightarrow \theta_j = \theta_j - \alpha (+\nu_c) = \theta_j - (+\nu_c)$$

It will decrease  $\theta_j$  until it get at or near global minima

$\alpha$  - alpha is a learning rate  $\alpha = 0.001$  → very good practice

It controls the speed at which convergence should happen

If it is more small it will take more time to converge

If it is very-very big it will jump continuously and it may never converge

Select a small value but not very-very small

## Linear Regression Example

Dataset

I/p  
Weight

O/p  
height

$$h_{\theta}(x) = \theta_0 + \theta_1 x \quad \text{I/p or independent}$$

$\theta_0$  = Intercept

$\theta_1$  = Slope

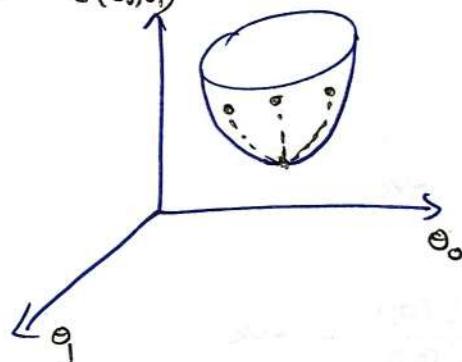
  Dataset : Multiple Regression :-

No. of rooms      Size of the house      Independent feature      dependent feature

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$\theta_0$  = intercept

$\theta_1, \theta_2$  = slope



Multiple regression formula:

$$\boxed{h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n}$$

## Performance matrix used in Linear Regression:

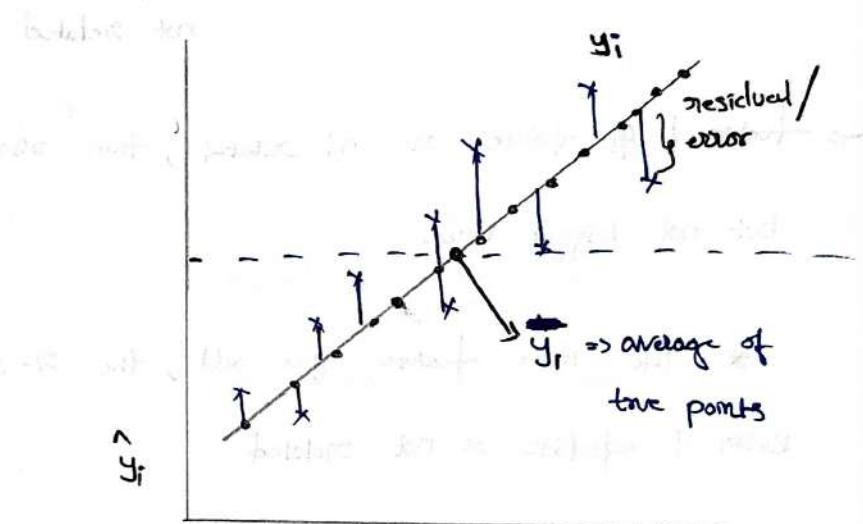
It is used to determine whether the model is good or not for a specific problem statement.

⇒ The metrics that basically used as

- ① R Squared
- ② Adjusted Squared

① R-Squared:-

$$R\text{-squared} = 1 - \frac{SS_{Res}}{SS_{Total}}$$



$$= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y}_i)^2} \Rightarrow \text{Smaller compare to denominator}$$

$$= 1 - \frac{\text{Small number}}{\text{big number}}$$

$$= 1 - \text{Small number}$$

$$\approx 1$$

$0.70 \Rightarrow 70\% \text{ Accuracy}$

$1 \Rightarrow 100\% \text{ Accuracy}$

$0.80 \Rightarrow 80\%$

$0.94 \Rightarrow 94\%$

## 2) Adjust $\hat{R}$ squared :

### Dataset

- Size of the house      price       $\Rightarrow R\text{-squared} = 75\%$
  - Size of the house      Location      price       $\Rightarrow R\text{-squared} = 80\%$
  - Size of the house      Location      gender      price       $\Rightarrow R\text{-squared} = 83\% \Rightarrow 0.83$   
 not related
- $\rightarrow$  feature + gp feature are not related, then also R-squared value will increase but not bigger value
- \* The more feature you add, the R-squared value is going to increase even if feature is not related

To overcome this we use  $\text{Adjusted R-squared}$

$$\text{Adjusted R-squared} = 1 - \frac{(1-R^2)(N-1)}{N-P-1}$$

N = No of data points

P = No. of independent features

If the feature is related  $\Rightarrow$  value increase

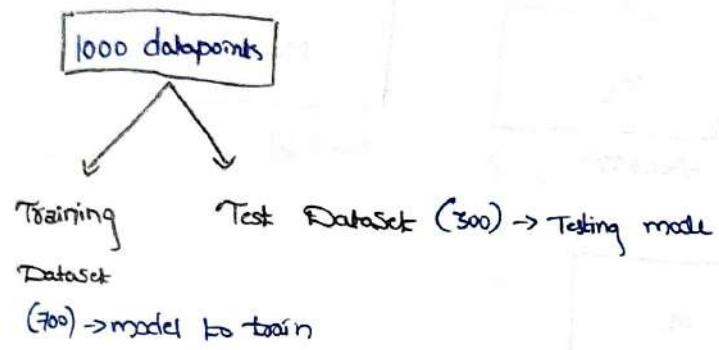
If the feature is not related  $\Rightarrow$  value decrease

$\text{Ex: } P=2 \quad R^2 = 90\% \quad R^2\text{-adjusted} = 81\%$ one no-related feature is added $\Rightarrow P=3 \quad R^2 = 92\% \quad R^2\text{-adjusted} = 82\%$
---

## Overshooting And undershooting (Bias And Variance)

- ① Training data
  - ② test data
  - ③ validation dataset

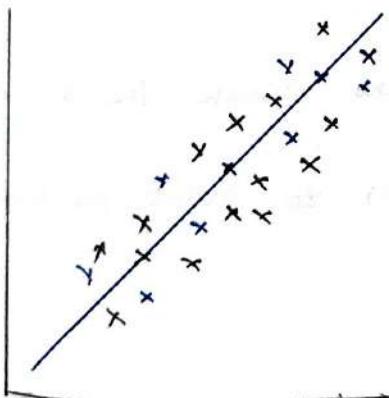
Draft :-



Train	Very	Good Accuracy [low bias]	90%
Test	Very	Good Accuracy [low variance]	85%

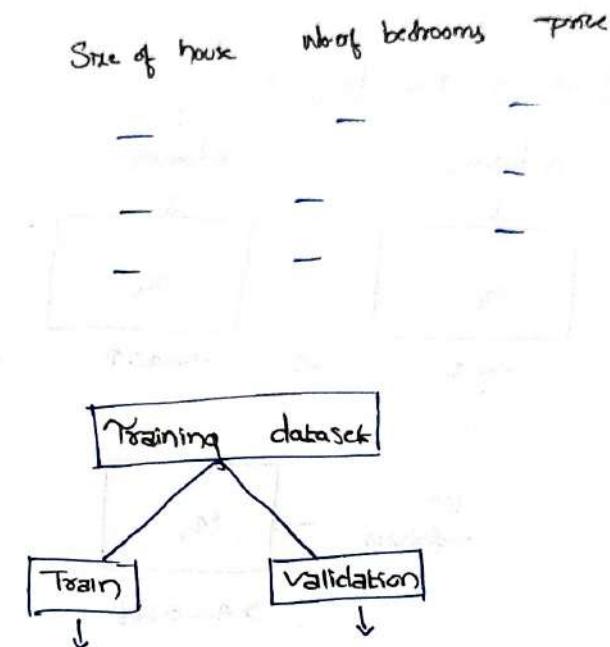
$\Rightarrow$  Generalized model

Train | Very Good (90%)  
 Test | Bad Accuracy (50%)  
 -> Model is overfitting



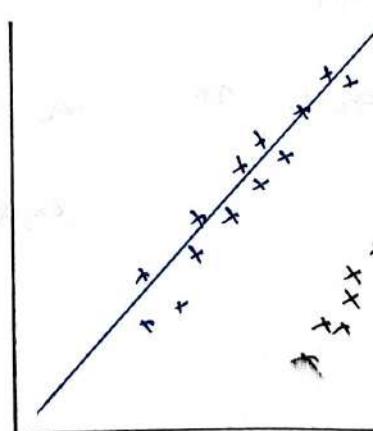
## Generalized model

low bias, low variance

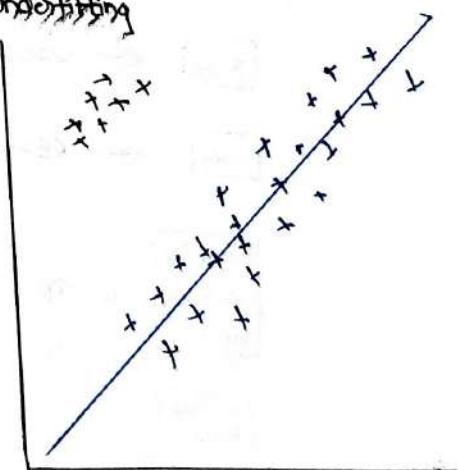


Training | Model Accuracy is low [High Bias]  
Test | Model Accuracy is low [High Variance]

Model is underfitting



Overfitting  
Low Bias, High Variance



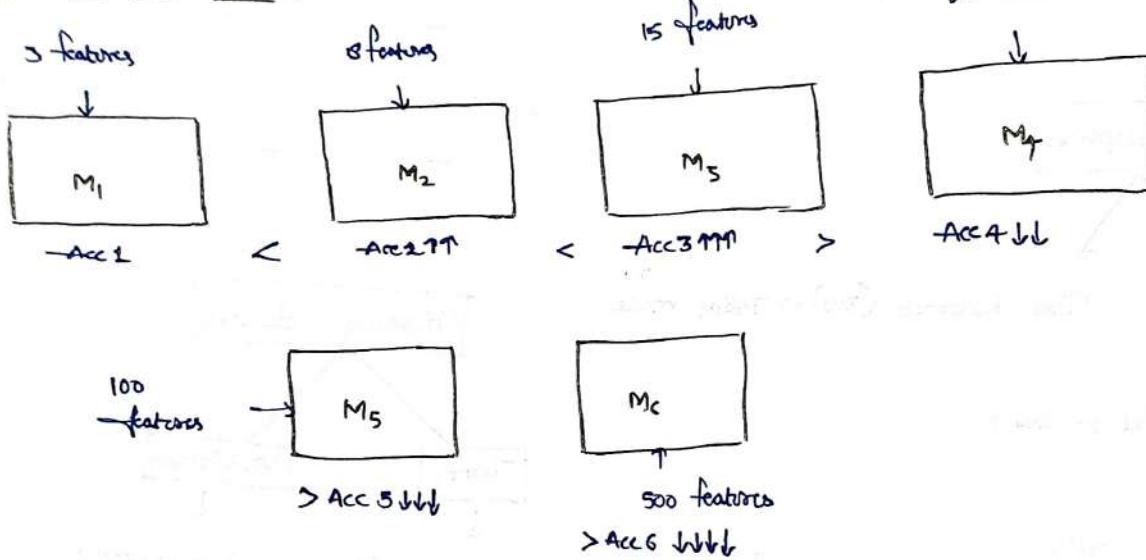
underfitting

# Principal Component Analysis (PCA) [Dimensionality Reduction]

① Cause of Dimensionality :-

Dimensions {features}

=> Price of the house :



② Model performance degrades :-

A person asking the domain expert about the house price

$f_1$   
loc<sup>+</sup>  $\leftarrow$  500k - 550k

3bhk  $\leftarrow$  600k ↑↑

beach  $\leftarrow$  150k - 700k ↑↑

Near to  
celebrity  
house  $\leftarrow$  ↑↑

Grocery  
shop

As the features increase the domain

Expert confused to predict the price

Two different ways to remove curse of Dimensionality

1 Feature Selection



The important features  
are Selected

② Dimensionality Reduction (PCA)



Feature Extraction

original

$f_1, f_2, f_3$  o/p

↓ Feature Extraction

$D_1, D_2$  o/p

Feature Selection

vs Feature Extraction

↳ Dimensionality Reduction

③ Why Dimensionality Reduction?

\* Prevent → Curse of Dimensionality

\* Improve the performance of the model

\* Visualize the data → understand the data

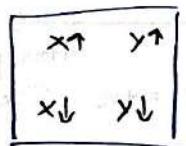
100 feature are difficult to visualize

⇒ 3-5 feature are easy to visualize

Feature Selection:

I/p o/p

+vc



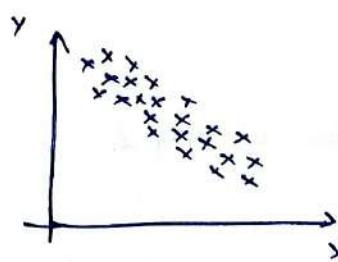
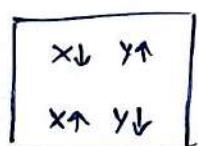
x y

-

- -

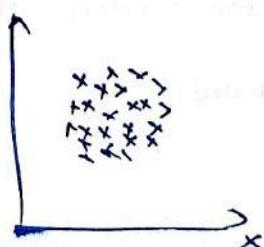
- -

-vc



No linear relationship

b/w x & y



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

↑                      ↓                      ↗

it may be +ve / -ve /  $\approx 0$

↳ Selection

Pearson Correlation =  $\frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = [-1 \rightarrow 1]$

The more towards the value of +1 the more the  
+ve correlated  $x \& y$

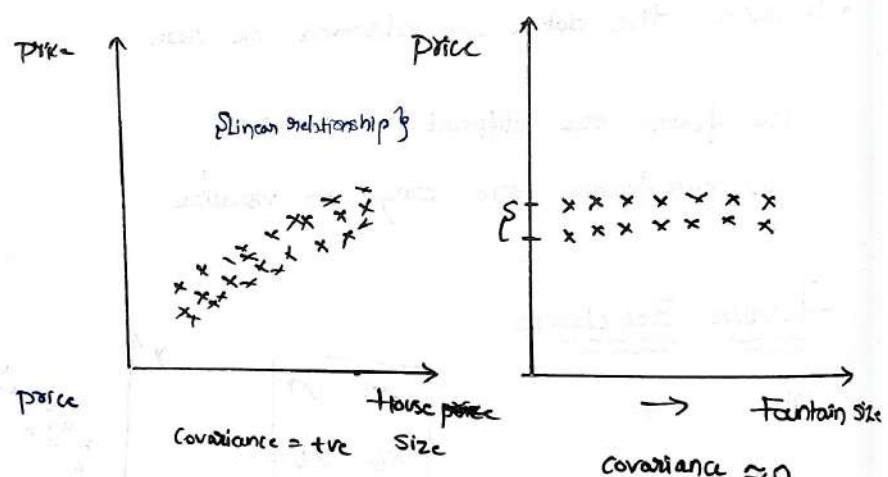
Dataset      Housing :-  
Example

① Feature Selection :-  
Independent features

House size	Fountain size
-	-
-	-

Fountain size feature is dropped

so, basically in Feature Selection the important  
features are selected



even though fountain size  
increase the house price  
is stagnant in a  
specific region

$\Rightarrow$  covariance is very low  
or zero

If it is not that important  
feature

Example :- Feature Extraction

~~ROOM~~ Size | No. of. Rooms | price Dimensionality Deduction 2f → 1f

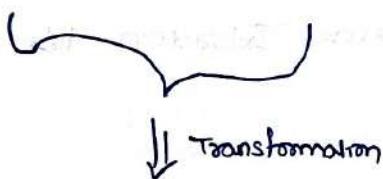
We cannot use feature selection to perform dimensionality reduction because both the features are very correlated or -very correlated

\* We cannot drop one feature

→ So use feature extraction by taking the features & applying transformations to extract new features

$\Rightarrow$  feature extraction

Room Size      No. of Rooms



House price | price

{ Dimensionality reduction }

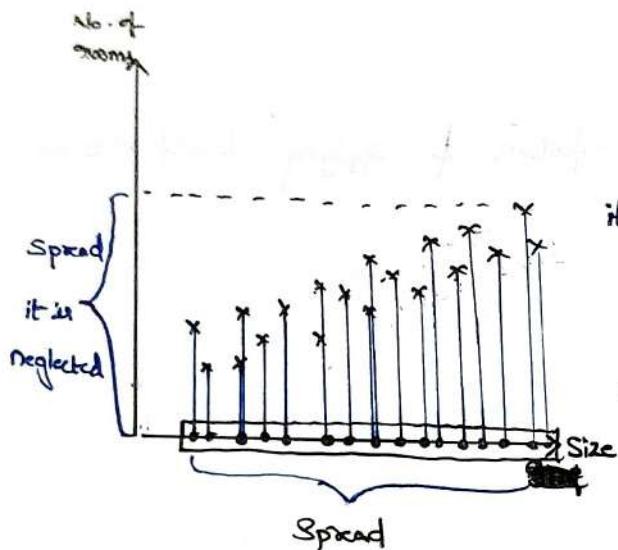
→ Housing dataset

PCA

Size of house | No. of rooms | price

$2D \rightarrow 1D$

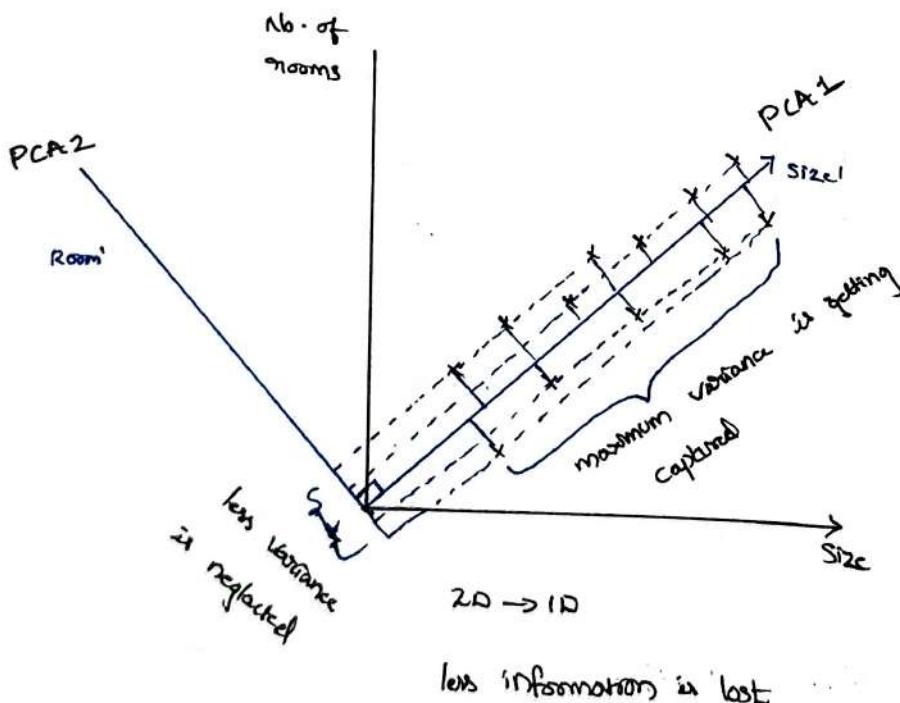
↓  
our aim



if project of x-axis a huge spread is neglected  
 ⇒ loss of information (No of Rooms)  
 ⇒ In this feature extraction lots of information is lost

Spread increases → Variance increase

★ To prevent huge loss of information PCA is used.



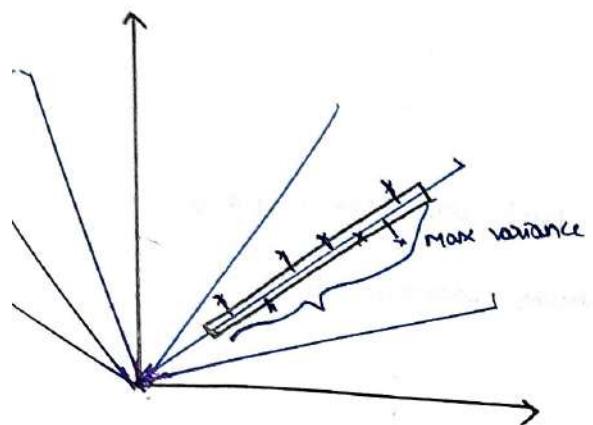
We will apply Eigen decomposition  
 on matrix If it is transformed  
 To get new axis

3 Dimension

⇒ PC1, PC2, PCA

Variance (PC1) > Variance (PC2) > Variance (PC3)

⇒  $\sigma^2(\text{PC1}) > \sigma^2(\text{PC2}) > \sigma^2(\text{PC3})$



**2D - 1D**

We Select 2 Best PCA line (i.e. the lines capturing high variances)

3 Dimensions: To get the best principle component which captures maximum variance

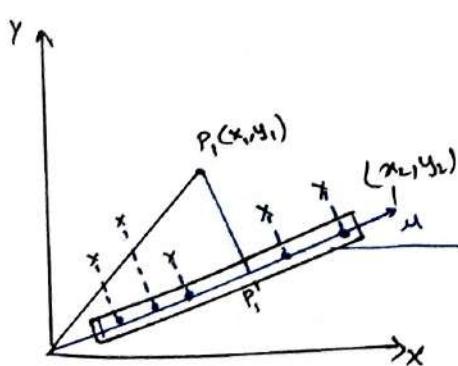
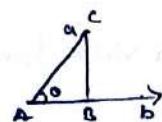
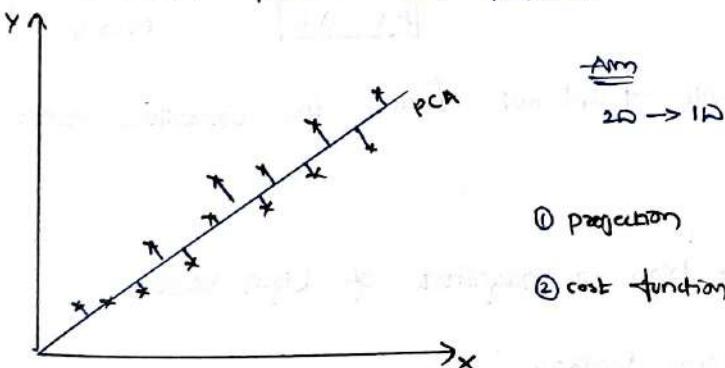
$$\text{Var}(\text{pc}_1) > \text{Var}(\text{pc}_2) > \text{Var}(\text{pc}_3)$$

$$3D \Rightarrow \text{pc}_1, \text{pc}_2, \text{pc}_3$$

**3D  $\Rightarrow$  1D**

**3D  $\Rightarrow$  2D**

Maths Intuition behind PCA Algorithm:



$$\text{Proj}_{P_i} u = \frac{P_i \cdot u}{\|u\|} \Rightarrow \boxed{\text{Proj}_{P_i} u = P_i \cdot u}$$

$$P_0^1, P_1^1, P_2^1, P_3^1, P_4^1 \dots P_n^1$$

scalar values

$$p_0^1, p_1^1, p_2^1, \dots, p_n^1$$

$$x_0^1, x_1^1, x_2^1, x_3^1, \dots, x_n^1$$

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

↓  
Cost function.

} goal : find the best unit vector which captures maximum variance }

Obviously we ~~does~~ cannot select different-different unit vectors

The technique eigen decomposition which we specifically say as Eigen vector & Eigen values

Eigen vector & Eigen values

To find the best unit vector following steps are required.

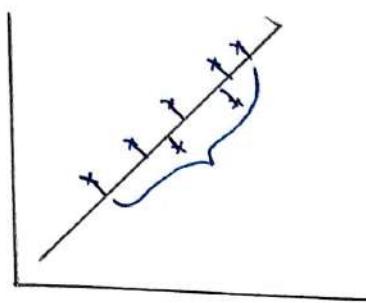
- ① Covariance matrix between features
- ② Eigen vectors & Eigen values will be found out from the covariance matrix
- ③ Largest Eigen vector  $\Rightarrow$

$$AV = \lambda V$$

linear transformation of Matrix

Eigen vector having Eigen value high  $\rightarrow$  magnitude of Eigen vector  
It will capture the maximum variance

Eigen decomposition of covariance matrix:-



→ we will find Eigen vectors & Eigen values for the matrix.

$$\begin{bmatrix} \text{covariance Matrix } A \\ \end{bmatrix} \times \begin{bmatrix} v \end{bmatrix} = \lambda \times v$$

$A \cdot v = \lambda \cdot v$

Linear transformation on A

$$\# \xrightarrow{LT} \#$$

- \* Eigen vector having maximum magnitude that will be use as principal component and it will calculate maximum variance

$$A \cdot v = \lambda \cdot v$$

Eigen vector  $\Rightarrow$  Max magnitude  $\Rightarrow$  Max Eigen value  $\Rightarrow$  Best PC  $\rightarrow$  PC1

Steps to calculate Eigen value and vectors:

### ① Covariance of feature

$$\begin{bmatrix} x & y \end{bmatrix} \quad z$$

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

If it is a  $2 \times 2$  matrix

$$\begin{array}{|c|c|} \hline x & y \\ \hline \text{cov}(x,x) & \text{cov}(x,y) \\ \hline \text{cov}(y,x) & \text{cov}(y,y) \\ \hline \end{array}$$

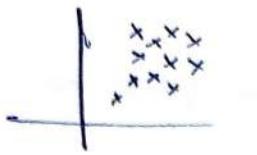
$\Rightarrow$  covariance matrix = A

$$A \cdot v = \lambda \cdot v$$

x	y	z	
x	$\text{cov}(x)$	$\text{cov}(x,y)$	$\text{cov}(x,z)$
y	$\text{cov}(y,x)$	$\text{var}(y)$	$\text{cov}(y,z)$
z	$\text{cov}(z,x)$	$\text{cov}(z,y)$	$\text{var}(z)$

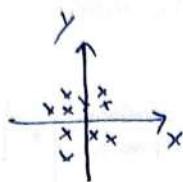
Since we are having 2 features x,y we will get 2  $\lambda$  values  $\lambda_1$  &  $\lambda_2 \Rightarrow$  eigen values.

Magnitude  
↓  
PC1      PC2



$2D \rightarrow 1D$

① Standardize the data



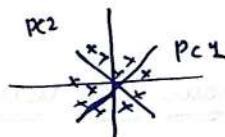
② Covariance matrix of  $x$  &  $y$

$A_2$	$\begin{array}{ c c } \hline \text{Var}(x) & \text{Cov}(x,y) \\ \hline \text{Cov}(y,x) & \text{Var}(y) \\ \hline \end{array}$
-------	---

③ Find out eigen vectors & Eigen values

$$A \cdot V = \lambda V$$

$$\begin{bmatrix} \lambda_1, \lambda_2 \\ \downarrow \quad \downarrow \\ PC_1 \quad PC_2 \end{bmatrix} \rightarrow \text{Eigen values}$$



$3D \rightarrow 1D$

$$\begin{array}{c} \lambda_1, \lambda_2, \lambda_3 \\ \downarrow \quad \downarrow \quad \downarrow \\ PC_1 \quad PC_2 \quad PC_3 \\ \checkmark \quad \quad \quad | \\ [1D] \quad [1D] \rightarrow [2D] \end{array}$$

$3D \rightarrow 1D$

$$\begin{array}{c} \lambda_1, \lambda_2, \lambda_3 \\ \downarrow \quad \downarrow \quad \downarrow \\ PC_1 \quad PC_2 \quad PC_3 \\ \checkmark \quad \quad \quad | \\ [1D] \end{array}$$

$2D \rightarrow 1D$

$$\begin{array}{c} \lambda_1, \lambda_2 \\ \downarrow \quad \downarrow \\ PC_1 \quad PC_2 \\ \checkmark \rightarrow \text{projection} \\ [1D] \end{array}$$



We are trying to find best principle component that can fit for any data capturing maximum variance.

③ Perception :- It is a single layer neural network & it is a type of artificial neuron or the simplest form of a neural network

=> It performs binary classification that maps input features to an output features / decision, usually classifying data into one of two categories, such as 0 or 1

① Input layer

② hidden layer

③ Weights

④ Activation function

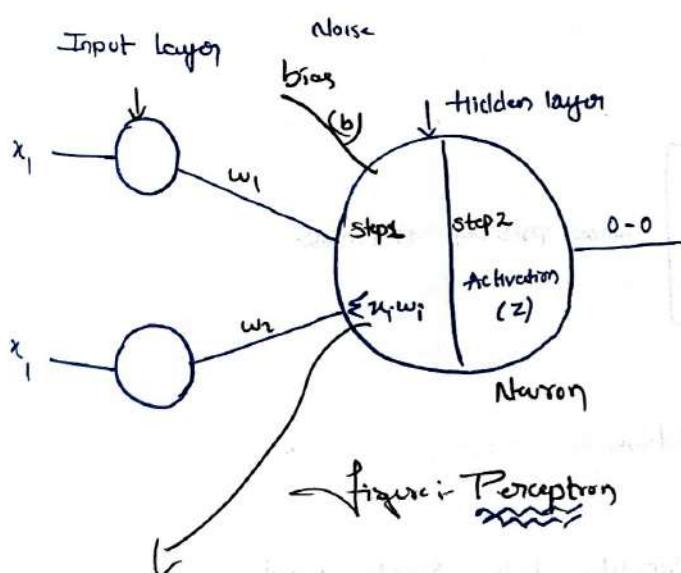
Dataset :-

IQ	No. of study hours	Pass/Fail
95	3	0
110	4	1
100	5	1

Activation function:

The aim is to transform the output

b/w the value 0 to 1 {use step function} -1 to 1



If weights are initialize zero then no processing happens so this bias should add 3 so that entire processing will not become zero

Noise :- A data that may completely different & our model is able to handle it

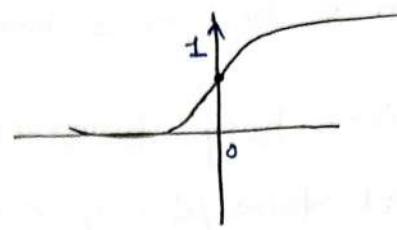
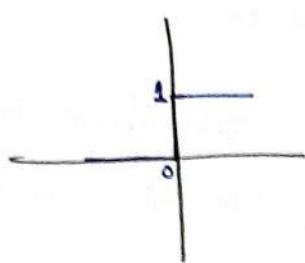
$$Z = x_1 w_1 + x_2 w_2 + b$$

$$Z = \sum_{i=1}^n w_i x_i + b$$

The output & real output are same error is zero

If not same change weight

### Step function



$$\begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

$$\begin{cases} 1 & z > 0.5 \\ 0 & z \leq 0.5 \end{cases}$$

Threshold value is change

$$0 \rightarrow 0.5$$

$\Rightarrow$  If output & real output are not same then change weights in order to get correct output

### Step 1

$$z = \sum_{i=1}^n w_i x_i + b$$

$$z = b + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

$$y = mx + c$$

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_n x_n$$

update

change

weights

in order to

get correct output

step 2

step 3

step 4

step 5

step 6

step 7

step 8

step 9

step 10

step 11

step 12

step 13

step 14

step 15

step 16

step 17

step 18

step 19

step 20

step 21

step 22

step 23

step 24

step 25

step 26

step 27

step 28

step 29

step 30

step 31

step 32

step 33

step 34

step 35

step 36

step 37

step 38

step 39

step 40

step 41

step 42

step 43

step 44

step 45

step 46

step 47

step 48

step 49

step 50

step 51

step 52

step 53

step 54

step 55

step 56

step 57

step 58

step 59

step 60

step 61

step 62

step 63

step 64

step 65

step 66

step 67

step 68

step 69

step 70

step 71

step 72

step 73

step 74

step 75

step 76

step 77

step 78

step 79

step 80

step 81

step 82

step 83

step 84

step 85

step 86

step 87

step 88

step 89

step 90

step 91

step 92

step 93

step 94

step 95

step 96

step 97

step 98

step 99

step 100

step 101

step 102

step 103

step 104

step 105

step 106

step 107

step 108

step 109

step 110

step 111

step 112

step 113

step 114

step 115

step 116

step 117

step 118

step 119

step 120

step 121

step 122

step 123

step 124

step 125

step 126

step 127

step 128

step 129

step 130

step 131

step 132

step 133

step 134

step 135

step 136

step 137

step 138

step 139

step 140

step 141

step 142

step 143

step 144

step 145

step 146

step 147

step 148

step 149

step 150

step 151

step 152

step 153

step 154

step 156

step 158

step 160

step 162

step 164

step 166

step 168

step 170

step 172

step 174

step 176

step 178

step 180

step 182

step 184

step 186

step 188

step 190

step 192

step 194

step 196

step 198

step 200

step 202

step 204

step 206

step 208

step 210

step 212

step 214

step 216

step 218

step 220

step 222

step 224

step 226

step 228

step 230

step 232

step 234

step 236

step 238

step 240

step 242

step 244

step 246

step 248

step 250

step 252

step 254

step 256

step 258

step 260

step 262

step 264

step 266

step 268

step 270

step 272

step 274

step 276

step 278

step 280

step 282

step 284

step 286

step 288

step 290

step 292

step 294

step 296

step 298

step 300

step 302

step 304

step 306

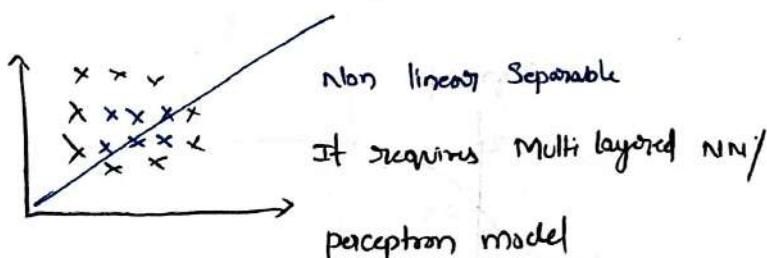
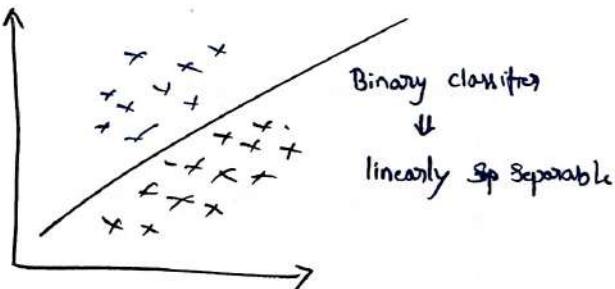
step 308

step 310

## Perception models

### Single Layer perception model

- It is a feed forward neural N/w
- The real output is 0. example we get 1 as output then again we go back & update weights & perform feed forward Neural N/w until it matches
- This not a very efficient technique except linearly separable problems
- problem



### Multi Layered perception model

- \* Forward propagation
- \* Backward propagation
- \* Loss functions
- \* Activation functions
- \* Optimizer

# Artificial Neural Network

Dataset 1.

① Forward propagation

$x_1$        $x_2$        $x_3$

② Backward propagation

IQ      study hours      play hours

Pass / Fail

③ Loss function

95      4      4

④ Activation function

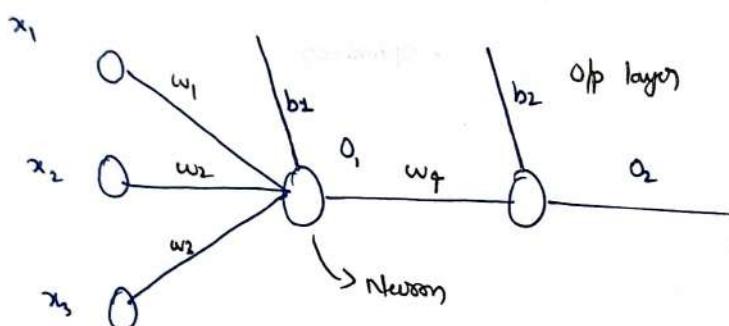
100      5      2

⑤ Optimizers

95      2      7

0

Forward propagation



The hidden layer can have more than one Neuron

① Forward propagation:

Hidden layer 1

$$\text{Step 1: } z = \sum_{i=1}^n w_i x_i + b$$

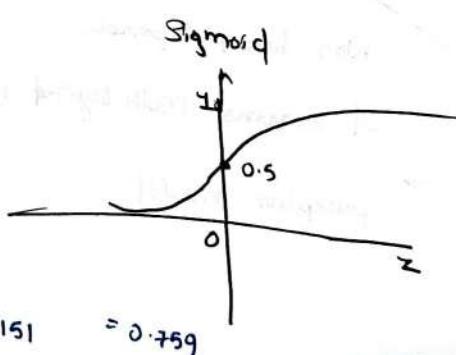
$$= 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 1 \times 0.01$$

$$= 1.151$$

Step 2: Activation ( $z$ )

$$f(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^{-1.151}} = 0.759$$

$$o_1 = 0.759$$



hidden layer 2

Step 1 :  $Z = O_1 \cdot w_2 + b_2$

$$= 0.759 \cdot 0.02 + 0.03$$
$$= 0.04518$$

Step 2 : Activation ( $z$ )

$$Z = \frac{1}{1+e^{-Z}} = \frac{1}{1+e^{-(0.04518)}} = 0.51129$$

$$O_2 = 0.51129 = \hat{y} \text{ (predicted)}$$

Loss function =  $(\hat{y} - O_2)^2$  real output

$$\approx 0.49 \text{ (error)}$$

Error is high our main aim is to reduce the error

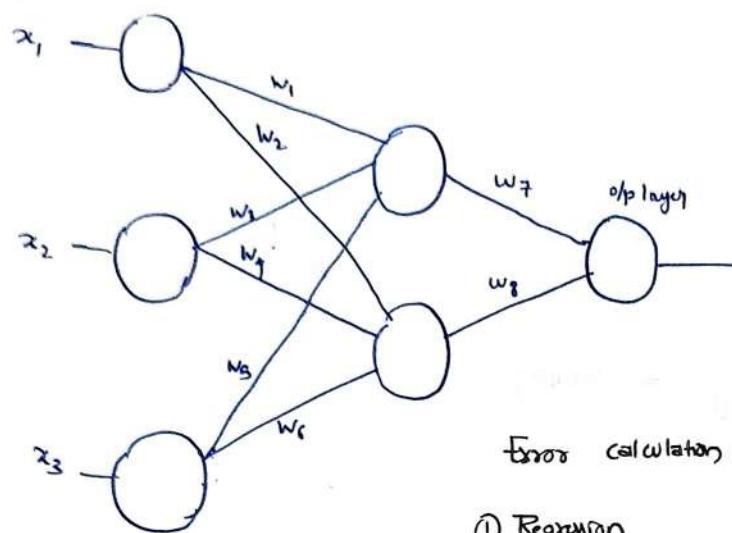
To minimize the error the weights should be updated.

Updated process will going from backward propagation & update the weight again all the steps repeated and loss function is used if error the again repeat.

Loss function vs cost function

$$\rightarrow (Y - \hat{Y})^2 \quad \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

## Backward propagation And Weight update formula



Loss function  
 $(y - \hat{y})^2$

	$x_1$	$x_2$	$x_3$	Pass / fail
IQ	study hours	play hours	O/p	
95	4	4	1	
100	5	2	1	
95	2	7	0	

Error calculation functions

① Regression

- MSE
- MAE
- huber loss

② Classification

- Binary class entropy
- Categorical class entropy

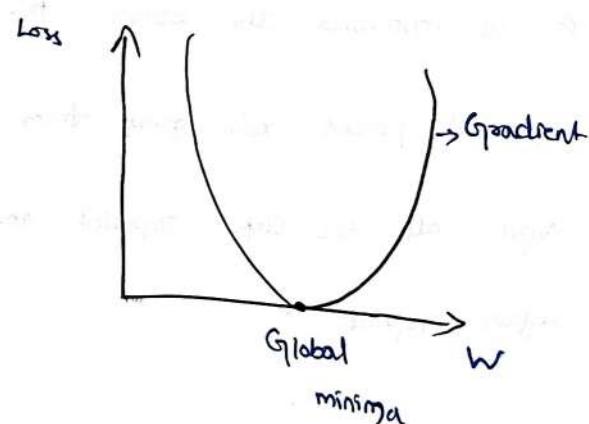
\* If we get very big loss function we need to update the weight it will happen in backward propagation

Weight update formula:

$$w_{\text{new}} = w_{\text{old}} - \eta \left[ \frac{\partial L}{\partial w_{\text{old}}} \right]$$

slope

$\eta$  ⇒ It is learning rate



If it is more small it will take more time to converge

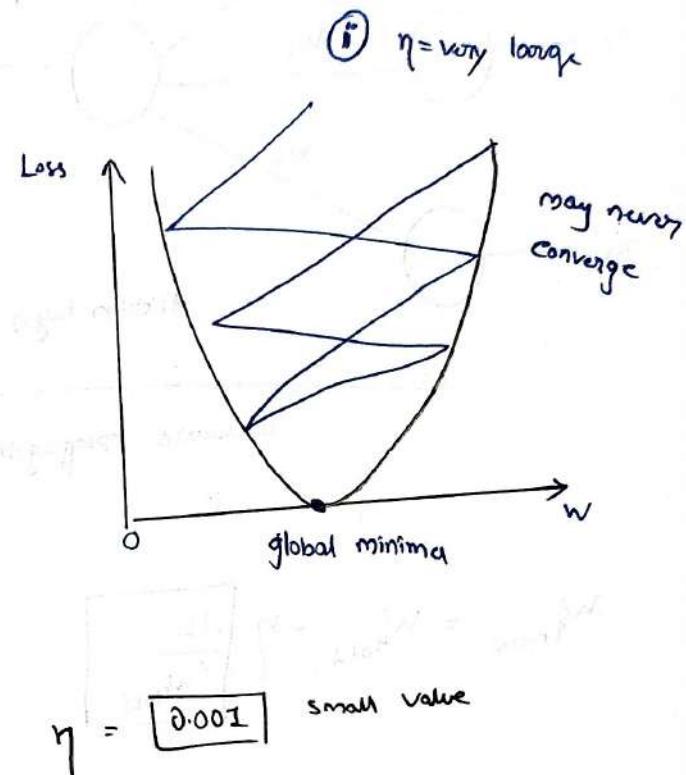
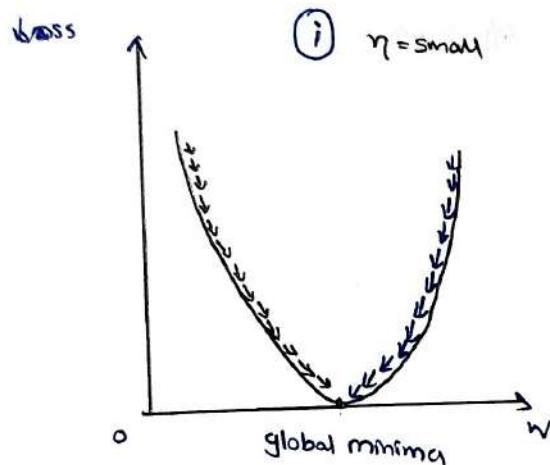
If it is very big value will jump continuous & may never converge

\* Select a small value but not very-very small

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

optimizes  $\rightarrow$  To reduce the loss function

Gradient Descent optimizer:



$$w_{\text{new}} = w_{\text{old}} - \eta (-v_c)$$

$$= w_{\text{old}} + \eta (+v_c)$$

↳ learning rate

$$\boxed{w_{\text{new}} \gg w_{\text{old}}}$$

$$w_{\text{new}} = w_{\text{old}} - \eta (+v_c)$$

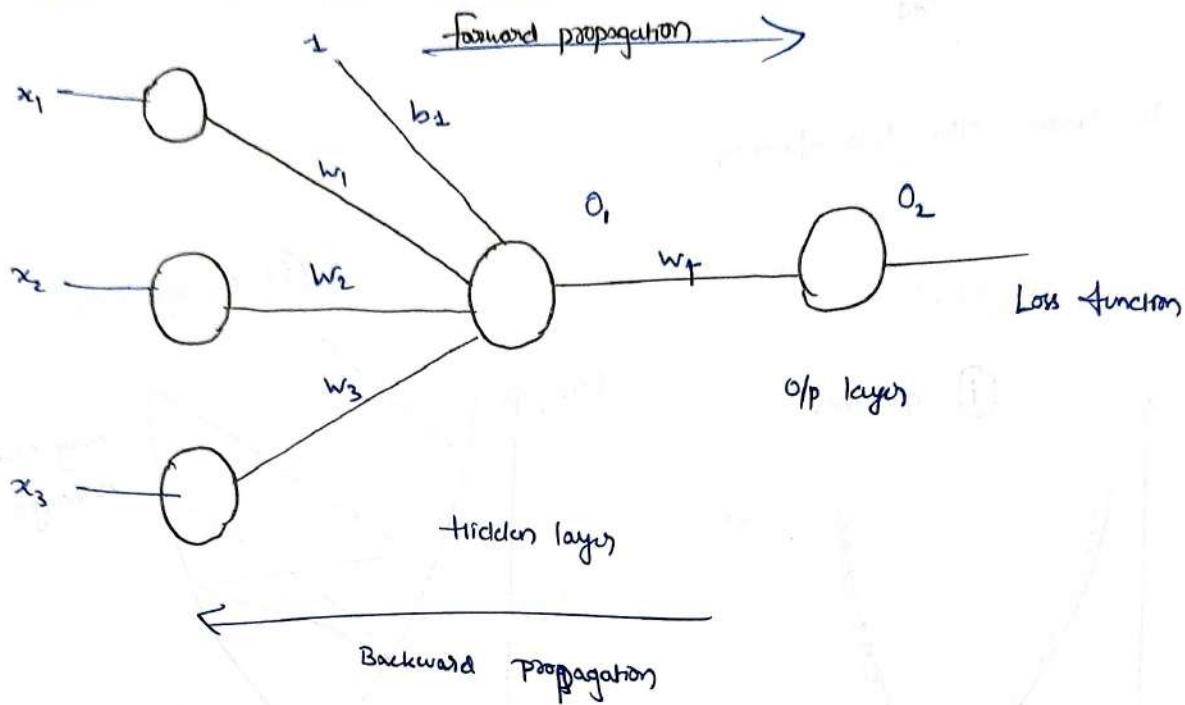
$$= w_{\text{old}} - \eta (+v_c)$$

$$\boxed{w_{\text{new}} \ll w_{\text{old}}}$$

It will stop when  $w$  reaches global Minima

$$\boxed{w_{\text{new}} = w_{\text{old}}}$$

## Chain Rule of Derivation :-



$$w_{t \text{ new}} = w_{t \text{ old}} - \eta \boxed{\frac{\partial L}{\partial w_{t \text{ old}}}}$$

$$\frac{\partial L}{\partial w_{t \text{ old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial w_t}$$

This splitting of the derivative is basically called as chain rule of derivative

$$\textcircled{1} \quad w_{t \text{ new}} = w_{t \text{ old}} - \eta \frac{\partial L}{\partial w_{t \text{ old}}}$$

$$\boxed{\frac{\partial L}{\partial w_{t \text{ old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_1}}$$

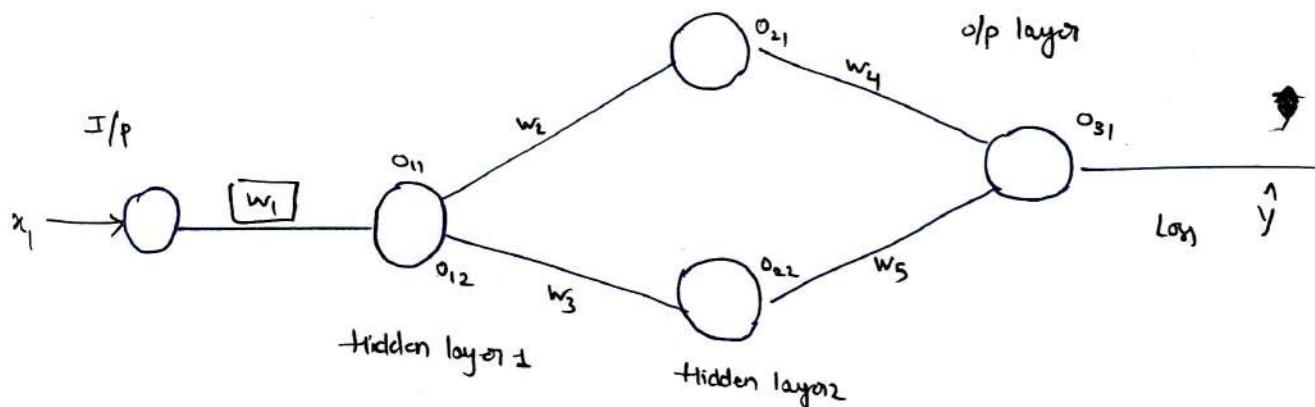
$$③ w_2 \text{new} = w_{2 \text{old}} - \eta \frac{\partial L}{\partial w_{2 \text{old}}}$$

$$\boxed{\frac{\partial L}{\partial w_{2 \text{old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_2}}$$

$$③ w_3 \text{new} = w_{3 \text{old}} - \eta \frac{\partial L}{\partial w_{3 \text{old}}}$$

$$\boxed{\frac{\partial L}{\partial w_{3 \text{old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_3}}$$

Assignment :- How do you update  $w_{1 \text{new}}$ ?



$$w_{1 \text{new}} = w_{1 \text{old}} - \eta \frac{\partial L}{\partial w_{1 \text{old}}}$$

$$\frac{\partial L}{\partial w_{1 \text{old}}} = \left[ \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1 \text{old}}} \right]$$

+

$$\left[ \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{12}} * \frac{\partial o_{12}}{\partial w_{1 \text{old}}} \right]$$