

# Mathematics - Basics to Advanced

## for Data Science & GenAI

### INTRODUCTION TO LINEAR ALGEBRA

Linear Algebra is a branch of maths that focuses on the study of vectors, vector/linear spaces, linear transformations, and system of linear eqns.

- Applications: → Data Representation and Manipulation : {Linear Algebra works with higher dimension data}
- ML & AI : → Model Train → Dimensionality Reduction
- Neural Network → Computer Graphics → Optimization  
(Fw & Bw propagation) → (gradient descent)

Scalar: A scalar is a single numerical value. It represents a magnitude (or) quantity & has no direction.

Ex: Car Speed = 45 kmph → (magnitude), Simple Linear Regression :  $y = mx + c$

(slope) ↓  
↓ (intercept)  
↓ (scalar value)

Vector: A vector is an ordered list of no.'s. It can represent a point in space (or) quantify with both magnitude & direction.

Ex: Speed of the car is 45 kmph & is moving towards East direction :  $\rightarrow$  E,  
a vector representing person's weight over time : [70, 72, 75, 78]

•  $\hat{i}, \hat{j}$  → unit vector towards 'x' & 'y' axis

Addition of 2 vectors:  $A = \begin{bmatrix} x_1 \\ y_1 \end{bmatrix}, B = \begin{bmatrix} x_2 \\ y_2 \end{bmatrix} \rightarrow A+B = \begin{bmatrix} x_1+x_2 \\ y_1+y_2 \end{bmatrix} = \begin{bmatrix} x_3 \\ y_3 \end{bmatrix}$

(Point in space) → ex:  $(-4, 3) \rightarrow \begin{bmatrix} -4 \\ 3 \end{bmatrix} \rightarrow$  represented as  $\begin{bmatrix} -4 \\ 3 \end{bmatrix}$

Multiplication of Vectors: 3 types :

- Dot Product (Inner Product)
- Element wise multiplication
- Scalar multiplication

Dot Product: The dot prod. of 2 vectors results in a scalar & is calculated as the sum of the products of their corresponding components.

Ex:  $A = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, B = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$  :  $\begin{array}{c} \uparrow \quad \uparrow \\ \text{B}(2 \times 2) \\ \downarrow \quad \downarrow \\ \text{length of projected } B \end{array}$  :  $\begin{array}{l} A \cdot B = 5 \times 2 + 0 \times 2 = 10 \\ \text{To: } (\text{len. of projected } B) \cdot (\text{len. of vector } A) \end{array}$  (scalar val.)

Ex:  $A \cdot B = \begin{bmatrix} 2 \\ 3 \end{bmatrix} \begin{bmatrix} 4 & 5 \\ 3 & 5 \end{bmatrix} = 2 \times 4 + 3 \times 5 = 23$

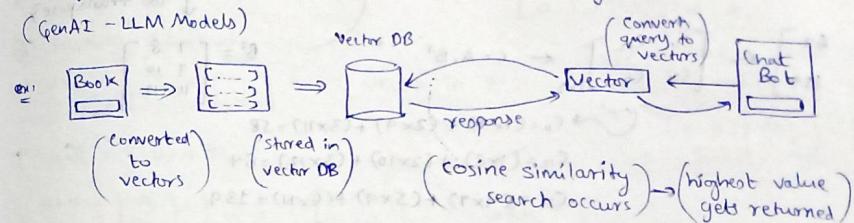
$\Rightarrow \vec{A} \cdot \vec{B} = 0 \rightarrow$  project the vector to the origin

Applications: 1) Cosine Similarity: It's a measure used to determine how similar 2 vectors are. It calculates the cosine of the angle b/w 2 vectors, providing a similarity score that ranges from -1 (dissimilar) to 1 (complete similar)

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

} Dot product

2) Vector Databases: When we talk about Vector database, we basically design a RAG (Retrieval Augmented Generation).



Element Wise Multiplication: In this, corresponding elements of 2 vectors are multiplied to form a new vector of the same dimension.

Ex:  $A = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, B = \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} \Rightarrow A \otimes B = \begin{bmatrix} 3 \\ 8 \\ 15 \end{bmatrix}$

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

Ex:  $A = \begin{bmatrix} 3 \\ 5 \\ 7 \end{bmatrix}, C = 4 \rightarrow CA = \begin{bmatrix} 12 \\ 20 \\ 28 \end{bmatrix}$

Applications: Normalization, Standardization.

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

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{bmatrix}_{m \times n}; a_{ij} \rightarrow i=\text{row}, j=\text{col.}$$

Ex: Data representation, images in computer vision, confusion matrix, nlp, neural net

Matrices Operations:

→ Matrix addition & subtraction

→ Scalar Matrix Multiplication

→ Matrix Multiplication

$$\begin{bmatrix} 50 & 10 \\ 5 & 35 \end{bmatrix}_{2 \times 2}$$

: 50 = true +  
10 = false -  
5 = false -  
35 = true -

Matrix Add. & Subtract.: Add or subtract corresponding elems of 2 matrices of the same dimension.

Ex:  $A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{bmatrix}, B = \begin{bmatrix} \dots & \dots & \dots \\ 7 & 8 & 9 \\ 1 & 2 & 3 \end{bmatrix} \rightarrow A+B = \begin{bmatrix} 1+4 & 2+5 & 3+6 \\ 4+7 & 5+8 & 6+9 \\ 7+1 & 8+2 & 9+3 \end{bmatrix} = \begin{bmatrix} 5 & 7 & 9 \\ 11 & 13 & 15 \\ 8 & 10 & 12 \end{bmatrix}$

Scalar Multiplication: Scalar multiplication involves multiplying every element of a matrix by a scalar value. ( $B = cA$ )

Matrix Multiplication: It involves the dot product of rows of the first matrix with columns of the 2nd matrix.

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

$$\text{Ex: } [1 \ 2 \ 3] \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix} = 1 \cdot 2 + 2 \cdot 3 + 3 \cdot 4 ; A = \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{2 \times 3} ; B = \begin{bmatrix} 7 & 9 & 11 \\ 8 & 10 & 12 \end{bmatrix}_{2 \times 3}$$

$$\begin{bmatrix} 58 & 64 \\ 139 & 154 \end{bmatrix} \leftarrow \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}_{2 \times 2} \leftarrow C = A \cdot B^T$$

↓ do transpose

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

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

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

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

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

### INTRODUCTION TO FUNCTIONS & TRANSFORMATIONS

Functions: A func. is a mathematical relationship that uniquely associates elements of one set (~domain) with elements of another set (~codomain).

Notation: A func. 'f' mapping ele. 'x' from set X to set Y is denoted by  $f: X \rightarrow Y$ . If  $x'$  is an element of X, then  $f(x')$  is the corresponding ele. in Y.

$$\text{Ex: } f(x) = 2x + 3$$

Vector Transformations: It refers to operations that map 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 computer graphics & data science.

Ex:

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

$$\text{Ex: } V' = 2V = 2 \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} 4 \\ 6 \end{bmatrix}$$

2) Rotation: transformation that turns vectors around the origin.

$$\text{Ex: } V = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \in \mathbb{R}^2 \xrightarrow{\text{Tr(2D, real no)}} V' = \begin{bmatrix} 0 \\ 1 \end{bmatrix} \xrightarrow{\text{rotation}} \begin{array}{c} \uparrow \\ \downarrow \end{array} \Rightarrow \begin{array}{c} \uparrow \\ \downarrow \end{array}$$

3) Reflection: transformation that flips vectors over a specified axis (or) plane.

$$\text{Ex: } V = \begin{bmatrix} 3 \\ 4 \end{bmatrix} \xrightarrow{\text{across the y-axis}} \begin{bmatrix} 3 \\ 4 \end{bmatrix} \cdot \begin{bmatrix} -1 \\ 1 \end{bmatrix} = \begin{bmatrix} -3 \\ 4 \end{bmatrix} \xrightarrow{\text{Ans}}$$

4) Shearing: It refers to distortion of shape of an obj. by shifting some of its points in a particular direction.

Linear Transformation: It's a func. b/w 2 vector spaces that preserves the operations of vector addition & 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  $\mathbf{u}, \mathbf{v} \in V$  & scalar  $c \in \mathbb{R}$ :

- Additivity:  $T(\mathbf{u} + \mathbf{v}) = T(\mathbf{u}) + T(\mathbf{v})$
- Homogeneity:  $T(c\mathbf{u}) = cT(\mathbf{u})$

where  $\mathbf{u}, \mathbf{v} \in V$  and 'c' is a scalar value.

Ex: reflection: The reflct. tran. T across Y-axis maps a vector  $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix} \in \mathbb{R}^2$  to  $T(\mathbf{x}) = \begin{bmatrix} -x \\ y \end{bmatrix}$

$\therefore$  Transformation can be expressed as:  $A = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \xrightarrow{\text{LT}} T(\mathbf{x}) = A \cdot \mathbf{x}$

④ Checking additivity: (we have to check whether it's LT or not)

$$\text{Let } \mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \text{ & } \mathbf{v} = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} \text{ be 2 vector in } \mathbb{R}^2.$$

$$T(\mathbf{u} + \mathbf{v}) = T(\mathbf{u}) + T(\mathbf{v})$$

$$\mathbf{u} + \mathbf{v} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} + \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} u_1 + v_1 \\ u_2 + v_2 \end{bmatrix} \xrightarrow{\text{(LHS)}}$$

$$T(\mathbf{u} + \mathbf{v}) = A(\mathbf{u} + \mathbf{v}) = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 + v_1 \\ u_2 + v_2 \end{bmatrix} = \begin{bmatrix} -u_1 - v_1 \\ u_2 + v_2 \end{bmatrix} \xrightarrow{\text{(LHS=RHS)}}$$

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

$$\text{RHS} \Rightarrow T(\mathbf{u}) + T(\mathbf{v}) = \begin{bmatrix} -u_1 \\ u_2 \end{bmatrix} + \begin{bmatrix} -v_1 \\ v_2 \end{bmatrix} = \begin{bmatrix} -u_1 - v_1 \\ u_2 + v_2 \end{bmatrix} \xrightarrow{\text{(RHS)}}$$

② Checking Homogeneity:

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

$$c\mathbf{u} = c \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} cu_1 \\ cu_2 \end{bmatrix}, T(c\mathbf{u}) = A(c\mathbf{u}) = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} cu_1 \\ cu_2 \end{bmatrix} = \begin{bmatrix} -cu_1 \\ cu_2 \end{bmatrix} \xrightarrow{\text{LHS=RHS}}$$

$$c \cdot T(\mathbf{u}) = c(A \cdot \mathbf{u}) = c \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = c \begin{bmatrix} -u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} -cu_1 \\ cu_2 \end{bmatrix}$$

Similarly there are examples which don't follow Linear Transformation.

Why? ~ used in :

- Dimensionality Reduction (PCA)
- Neural Networks
- Optimizing & solving sys. of eqn's
- Feature Engg.
- Image & Signal Processing (CNN)
- Data preprocessing (Normalization & Standardization)

Also, in LT → Origin must be fixed  
→ All lines must remain lines

(The word 'transformation' means the same thing as func.)

Magnitude & Unit Vector:

$$\text{Ex: } \mathbf{B} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \xrightarrow{\text{Magnitude}} \|\mathbf{B}\| = \sqrt{x_1^2 + x_2^2 + x_3^2} \quad \text{vector length}$$

$$\|\mathbf{B}\| = \sqrt{x_1^2 + x_2^2 + x_3^2}$$

Unit Vector:  $\|\mathbf{u}\| = 1$ , denoted by  $\hat{\mathbf{u}}$ .

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

scalar  
multiplication

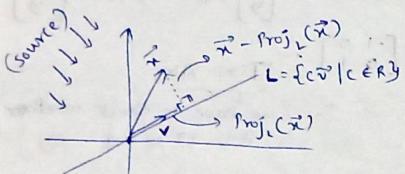
A point in space  $\vec{x}$  (2D) can also be written as:  $2\hat{i} + 3\hat{j}$ .

$$\text{ex: } \vec{v} = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix} \in \mathbb{R}^3 \rightarrow \|\vec{v}\| = \sqrt{1^2 + 2^2 + 0^2} = \sqrt{5}$$

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

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

## Introduction to Projections:



$\Rightarrow \text{Proj}_L(\vec{x})$  = projection of vector  $\vec{x}$  to  $L = \{c\vec{v} \mid c \in \mathbb{R}\}$

(some vector in line where  $\vec{x}$ - $\text{Proj}_L(\vec{x})$  is perpendicular to  $L$ )

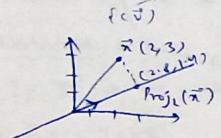
• Dot product of 2  $\perp$  vectors = 0.

$$\text{Proj}_L(\vec{x}) = (\vec{x} - \vec{v}) \cdot \vec{v} = 0 \rightarrow \vec{x} \cdot \vec{v} - \vec{v} \cdot \vec{v} = 0 \rightarrow \vec{x} \cdot \vec{v} = \frac{\vec{x} \cdot \vec{v}}{\vec{v} \cdot \vec{v}}$$

$$\therefore \text{Proj}_L(\vec{x}) = c\vec{v} = \left( \frac{\vec{x} \cdot \vec{v}}{\vec{v} \cdot \vec{v}} \right) \cdot \vec{v} \iff \text{Proj}_L(a) = \left( \frac{a \cdot b}{b \cdot b} \right) \cdot b$$

It gives us the point on  $\vec{v}$  where projection of  $\vec{x}$  is intersected.

$$\text{ex: } L = \{c \begin{bmatrix} 2 \\ 1 \end{bmatrix} \mid c \in \mathbb{R}\}, \vec{x} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$



$$\text{Proj}_L(\vec{x}) = \begin{pmatrix} [2] & [2] \\ [2] & [2] \end{pmatrix} \cdot \begin{bmatrix} \frac{2}{5} \\ \frac{1}{5} \end{bmatrix} = \frac{7}{5} \cdot \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{14}{5} \\ \frac{7}{5} \end{bmatrix}$$

$$\therefore \text{Proj}_L(\vec{x}) = \begin{bmatrix} 2.8 \\ 1.4 \end{bmatrix}$$

## INVERSE FUNCTIONS OR TRANSFORMATIONS

### Inverse of a function:

If func.  $f$  maps an ele.  $x$  from 'set  $X$ ' to an ele.  $y$  in 'set  $Y$ ', ( $f: X \rightarrow Y$ ), the inverse func.  $f^{-1}$  maps 'y' back to 'x' ( $f^{-1}: Y \rightarrow X$ ), for every  $y \in Y$ , there is a unique  $x \in X$  such that  $f(x)=y$ .

Conditions:

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

Identity function:  $I_x: X \rightarrow X$ . For a set  $X$ , the identity func.  $I_x$  is defined as:  $I_x(a) = a$ , for all  $a \in X$ .

(Identity func. on set 'X' & it maps every element 'x' in  $X$  to itself.)

- Properties:
- 1) Preservation: doesn't alter any ele.
  - 2) Linearity: Identity func. is a linear transformation
  - 3) Identity Matrix: all diag.-s = 1 & rest = 0's
  - 4) Inverse: Identity func. is its own inverse

ex:  $I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}_{3 \times 3}$

Existence & Uniqueness: A func. 'f' has an inverse if & only if it's bijective.

Injective (1-to-1): diff. ele.'s in the domain map to diff. ele.'s in the codomain.

surjective (onto): every ele. in codomain is the image of atleast 1 ele. in the domain.

Find the Inverse: ex:  $y = 2x + 3$  for 'x';  $y = 2x + 3 \Rightarrow x = \frac{y-3}{2}$

$$\therefore \text{Inverse func. } f^{-1}(y) = \frac{y-3}{2}$$

$$\text{Verification: } 1) f(f^{-1}(y)) = f\left(\frac{y-3}{2}\right) = 2\left(\frac{y-3}{2}\right) + 3 = 3y$$

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

## Applications of Inverse func.'s

(mean) (std.deviation)

### Normalization & Standardization

Use Case: After training model on standardized data, the predictions are often rescaled back to the org. scale to interpret the results.

Standard Normal Distribution:  $\mu=0$  &  $\sigma=1$

Standard transformation  $\rightarrow z = \frac{x-\mu}{\sigma}$

Inverse transformation  $\rightarrow x = z\sigma + \mu$

Min-Max Normalization: Org. Transformation:  $z = \frac{x - \min(x)}{\max(x) - \min(x)}$   $\Rightarrow z: x \rightarrow y$

Inverse Transformation:  $x = z(\max(x) - \min(x)) + \min(x)$

Distribution of Data:  $\rightarrow$  Normal distribution

### Data encryption & decryption

### Find inverse of a matrix

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

Determinant:  $\rightarrow$  It's a scalar value.

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \Rightarrow ad - bc \Rightarrow \text{scalar, determinant}$$

If  $\det(A) \neq 0$   $\rightarrow$  Inverse of the matrix exists.

## EIGEN VECTORS AND EIGEN VALUES

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

Eigen Vector ( $v$ ): A non-zero vector that only changes in scale (not direction).

For a square matrix ' $A$ ' =  $\begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix}$ , an eigen vector & its corresponding eigen value ' $\lambda$ ' satisfies :  $Av = \lambda v$

1) Find Eigen Values:  $\det(A - \lambda I) = 0$

$$\begin{bmatrix} 4-\lambda & 1 \\ 2 & 3-\lambda \end{bmatrix} \rightarrow \det. = (4-\lambda)(3-\lambda) - 2(1) = \lambda^2 - 7\lambda + 10$$

2) Solve the eqn.:  $\lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \rightarrow \lambda_1 = 5, \lambda_2 = 2$  eigen values of  $A$

3) Find eigen vectors:  $(A - \lambda I)v = 0$

$$\text{For } \lambda_1 = 5 \rightarrow A - 5I = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix} - \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix} = \begin{bmatrix} -1 & 1 \\ 2 & -2 \end{bmatrix} \rightarrow \begin{bmatrix} -1 & 1 \\ 2 & -2 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0 \rightarrow -x + y = 0 \Rightarrow x = y$$

$$\text{For } \lambda_2 = 2 \rightarrow A - 2I = \begin{bmatrix} 4-2 & 1 \\ 2 & 3-2 \end{bmatrix} \rightarrow \begin{bmatrix} 2 & 1 \\ 2 & 1 \end{bmatrix} \cdot \begin{bmatrix} x \\ y \end{bmatrix} = 0 \rightarrow 2x + y = 0 \rightarrow y = -2x$$

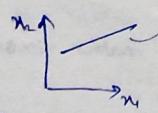
$$\therefore A = \begin{bmatrix} 4 & 1 \\ 2 & 3 \end{bmatrix}, \lambda_1 = 5, \lambda_2 = 2 \quad \underbrace{v_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}}_{\text{eigen vector}}, \quad \underbrace{v_2 = \begin{bmatrix} 1 \\ -2 \end{bmatrix}}_{\text{eigen value}} \quad (Av = \lambda v)$$

These eigen vectors & values describe how a matrix  $A$  scales & rotates vectors in its transformation. Eigenvalues indicate the factor by which the eigenvectors are stretched (or) compressed, and eigenvectors provide the directions in which this stretching (or) compression occurs.

Application: 'PCA' in dimensionality reduction.

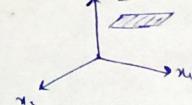
## EQN. OF A LINE, PLANE, HYPERPLANE

$$\text{Line eqn.: } y = mx + c \quad (\text{or}) \quad ax + by + c = 0 \quad (\text{or}) \quad w_1x_1 + w_2x_2 + b = 0 \quad \begin{array}{l} \text{(eqn. of} \\ \text{line)} \end{array}$$



• If line passes via Origin  $\rightarrow b=0$   $\therefore w^T x = 0$

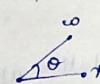
For 3D plane :



$$w_1x_1 + w_2x_2 + w_3x_3 + b = 0 \rightarrow w^T x + b = 0$$

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

for plane (line via origin):  $\Pi_n: w^T x = 0$



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

$w$  will always be  $\perp$  to plane if  $\theta = 90^\circ$

## INTRODUCTION TO STATISTICS

Statistics is the science of collecting, organizing & analyzing data.

Application: • Data Exploration & Summarize • Model building & Validation  
• Statistical Analysis • Hypothesis testing • Optimization & Efficiency  
• Reporting

Types of Statistics:

→ Descriptive Statistics

Involves summarizing & organizing data to make it understandable.

• Measures of Central Tendency:  $\begin{array}{l} \text{Mean} \\ \text{Median} \\ \text{Mode} \end{array}$

• Measure of dispersion (variance, std.deviation) • P-value

• Data Distributions: Histograms, Box plot, Pie chart

• Summary Statistics: Five No. Summary, Q1, Q2, Q3, Maximum

• Statistical Analysis Test:  $\begin{array}{l} Z\text{-test} \\ \text{ANOVA} \\ \chi^2\text{ test} \end{array}$

ex: Let say there are 20 students class in your college, & you have collected the height of students in the class.

→ Descriptive: "What is the avg. height of the entire classroom?"  
→ Inferential: "Are the heights of the sample students in classroom similar to what you expect in the entire college?"

→ Population

• It's the entire set of objects of interest in a particular study.

• Contains all the observation of interest

- parameter: a numerical value summarizing the entire population.

$$\begin{array}{l} \downarrow \sigma^2, \mu \rightarrow \text{mean} \\ \downarrow \text{variance} \end{array}$$

ex: all customers in a city

to understand the purchasing behaviour of all customers.

Inferential Statistics

Involves methods for making predictions or inferences about a population based on a sample of data.

• Hypothesis Testing

• P-value

• Confidence Interval

• Z-test  
ANOVA  
 $\chi^2$  test

Sample Data

• A sample is a subset of the population that is used to represent the entire group.

• Represents a portion of population.

- parameter: a numerical value summarizing the sample data.

$$\begin{array}{l} \downarrow \text{mean, variance} \end{array}$$

- samples should be randomly selected to avoid bias.

ex: 500 consumers from the city

Behaviour → population.

## Types of Sampling Techniques:

### 1) Probability Sampling:

a) Simple Random Sampling: every member of the population has an equal chance of being selected.

ex: selecting people randomly

b) Systematic Sampling: select every  $n^{\text{th}}$  member of the population from a random starting point.

ex: Feedback Survey → selected every 7<sup>th</sup> person's feedback.

c) Stratified Sampling: divide the population into strata (groups) based on specific characteristics & then randomly sampling from each strata.

d) Cluster Sampling: divide the population into clusters, randomly selecting clusters, then sampling all the members from the selected clusters.

e) Multi Stage Sampling: combining several sampling methods. Usually involves selecting clusters, then randomly sampling within those clusters.

### 2) Non Probability Sampling:

Select individuals who are easiest to reach. ex: surveying people at mall.

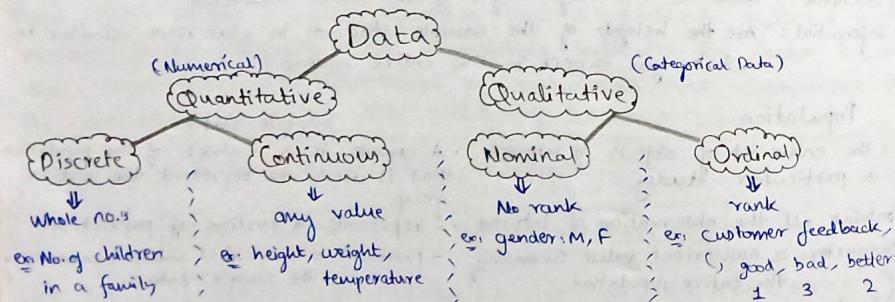
a) Convenience Sampling: selecting individual who are easiest to reach.

b) Judgemental Sampling: select individual based on the researcher's judgement.  
ex: choose experts in a field to participate.

c) Snowball Sampling: existing study subjects, recruit future subjects from among their acquaintances.

ex: survey members of a rare disease.

d) Quota Sampling: sampling based on age, group - gender, etc.



## Scales of Measurement of Data:

The scales of measurement describe the nature of info. within the values assigned to variables. We have 4 primary scales:

1) Nominal Scale

2) Ordinal Scale

3) Interval

4) Ratio.

### 1) Nominal Scale:

This scale classifies data into distinct categories that don't have an intrinsic order.  
Characteristics:  
→ Data is categorized based on labels, names or qualities.  
→ Categories are mutually exclusive.  
→ No logical order among categories [No rank].

ex: gender → M, F

2) Ordinal Scale: This scale classifies the data into categories that can be ranked in order.

Characteristics:  
→ Data is categorized & ranked in a specific order.  
→ The interval btw ranks are not necessarily equal.  
ex: customer feedback → Satisfied, Very Satisfied, Not satisfied  
(1) (2) (3)

3) Interval Scale: The interval scale not only categorizes & orders, but also specifies the exact difference btw intervals. It lacks a true zero point.

Characteristics:  
→ Data is ordered with consistent interval btw values.  
→ Allows for meaningful comparison of differences (Ratio can't be measured).  
→ No true zero point.

ex: IEL Scores: 90, 100, 110 → difference = 100 - 90 = 10. ( $10 \neq 0$ )

### 4) Ratio Scale:

→ The order matters.  
→ Differences are measurable.  
(Ratio can be measured).  
→ Contains a 0 starting point.

ex: student marks

(1) 0, 10, 20, 30

(2) 0, 30, 60, 90

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

## DESCRIPTIVE STATISTICS

### Measures of Central Tendency:

These are statistical measures which provides a single value that summarizes a set of data by identifying the central position within that dataset.

1) Mean: sum of all values divided by the no. of values.

$$\rightarrow \text{Population Mean } (\mu) \left\{ \begin{array}{l} \mu = \frac{\sum_{i=1}^N x_i}{N} \\ \text{Population size } (N) \end{array} \right.$$

$$\rightarrow \text{Sample mean } (\bar{x}) \left\{ \begin{array}{l} \bar{x} = \frac{\sum_{i=1}^n x_i}{n} \\ \text{sample size } (n) \\ \quad \downarrow n \leq N \end{array} \right.$$

Here 'X' = random variable.

Characteristics:

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

$$\text{ex: } X = \{1, 2, 3, 4, 5, 7\} \rightarrow \mu = \frac{1+2+3+4+5+7}{6} = 3$$
$$\text{ex: } X = \{1, 2, 3, 4, 5, 100\} \rightarrow \mu = \frac{1+2+3+4+5+100}{6} = 23$$

2) Median: It's the middle value in a dataset when the values are arranged in ↑ or ↓ order.

Characteristics:

- Not affected by extreme outliers.
- Used for ordinal, interval & ratio data.

$$\text{ex: } X = \{1, 2, 3, 4, 5, 6\} \rightarrow \text{No. of ele. } 5 \text{ (odd)} \rightarrow \text{Median} = 3$$
$$\text{ex: } X = \{1, 2, 3, 4, 5, 10\} \rightarrow \text{No. of ele. } 6 \text{ (even)} \rightarrow \text{Median} = \frac{3+4}{2} = 3.5$$

3) Mode: the value that appears most frequently in a dataset.  
 Chars: → Not affected by extreme values. → Used for nominal, ordinal, interval & ratio data.

### Choosing the appropriate measure:

1) Mean: Best used when data is symmetrically distributed without outliers. Provides a mathematical avg., which is useful for further statistical calculations.

2) Median: Best used when data is skewed or contain outliers. Provides middle value, which better represents the center of a skewed dataset.

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.

### Measure of Dispersion:

These describe the spread or variability of a dataset. They indicate how much the values in a dataset differ from the central tendency.

1) Range: It's the difference btwn the max. & min. value in a dataset.

Chars: → Simple to calculate  
 → Sensitive to outliers  
 → Rough measure of dispersion.

$$\text{Range} = \text{Max. Value} - \text{Min. Value}$$

2) Variance: It measures the avg. squared deviation of each value from the mean. It provides a sense of how much the values in a dataset

$$\rightarrow \text{Population Variance: } \sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

( $x_i$  = data points,  $\mu$  = population mean,  $N$  = population size)

$$\rightarrow \text{Sample Variance: } s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

( $x_i$  = data points,  $\bar{x}$  = sample mean,  $n$  = sample size)

### Chars:

- Provides a precise measure of variability.
- Units are squared of the original data units.
- More sensitive to outliers than the range.



Variance: determines the spread of the data.

3) Standard Deviation: It's the square root of the variance.

Chars: → Provides a clear measure of spread in the same units as the data.  
 → Sensitive to outliers.

If we have std. deviation, we can find out variance & vice-versa.

Variance = provides a measure of dispersion of data points, in squared units, which can be difficult to interpret directly.

Standard deviation = provides a measure of the dispersion in the same units as the original data, making it easier to interpret & understand.

### Why 'n-1' in sample variance?

When we calc. var. using just the sample data, we're only working with a small part of the whole population. This can lead to underestimating the true variability, which is the bias we're correcting. By using 'n-1', we adjust for this & get a more accurate estimate.

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1} \rightarrow \text{Bessel's Correction}$$

Random Variables:  $X$ : it's a func. which derives its values from process or experiments.

↓

### Discrete Random Variable

ex: tossing a coin,  
 rolling a dice

### Continuous Random Variable

ex: time how many inches it's going to rain  
 [0, 1.1, 5.5, 10.5, 10.75]

• Percentage: ex: {1, 2, 3, 4, 5, 6} → percentage of odd no.s } =  $\frac{3}{6} \times 100 = 50\%$  in this group

• Percentile: A percentile is a value below which certain % of observations lie

ex: {2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, 9, 9, 10},  $x=9$

$$\text{Percentile of value } x = \frac{\text{No. of values below } x}{n} \times 100 = \frac{11}{14} \times 100 = 78.57\% \text{ of value } 9$$

(means 78.57% of the entire distribution is less than 9.)

If percentile is given & we need the value;

$$\text{ex: } 25\% \rightarrow \text{Value} = \frac{\text{Percentile } x(n+1)}{100} \Rightarrow \text{Val.} = \frac{25}{100} \times 15 = 3.75$$

(position)

$$\therefore \frac{3+4}{2} \approx 3.5$$

∴ 25% is 3.75

(gives us position in dataset)

• Quartiles:

$$25\% = 1^{\text{st}} \text{ Quartile}$$

$$50\% = 2^{\text{nd}} \text{ Quartile}$$

$$75\% = 3^{\text{rd}} \text{ Quartile}$$

(every 25% of percentage is defined by quartile)

5 Number Summary: {to identify & remove the outliers}

1) Minimum    2)  $Q_1$  (25%) → 1<sup>st</sup> Quartile    3) Median    4)  $Q_3$  (75%)    5) Maximum  
 ex: 1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27

1<sup>st</sup> Step: (to remove outliers): define lower fence & higher fence value. Anything higher than higher fence value will be removed as an outlier & anything lower than lower fence value will be removed as an outlier.

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

$$\therefore Q_1 = \frac{\text{Percentile } x(n+1)}{100} = \frac{25}{100} \times 20 = 5^{\text{th}} \text{ position} = 3$$

$$Q_3 = \frac{75}{100} \times 20 = 15^{\text{th}} \text{ position} = 7 \quad \therefore IQR = 7 - 3 = 4$$

$$\therefore [(\text{Lower fence } \leftarrow \text{Higher fence})] = [-3 \longleftrightarrow 13]$$

→ Interquartile Range ( $Q_3 - Q_1$ )

∴ '27' eliminated as an outlier

These 5 values define 5 no. summary & based on this we can define box plot, used to visualize the outliers.

## Histogram & Skewness.

(graphical representation of the distribution of numerical data.)

ex: ages = {11, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 58} → histogram  
let no. of bins = 10 →  $\frac{58}{10} = 5$  bin size  
range: 0 - 50

(no. of intervals into which my data range is divided.)

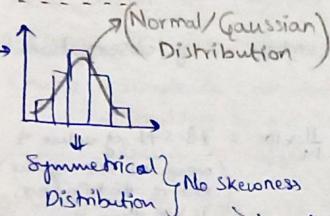
Bins → [0-5, 5-10, 10-15, ..., 45-50] :  
(>5, <20)

(width of these intervals)

Count(freq.)



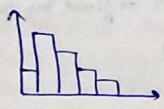
Skewness:



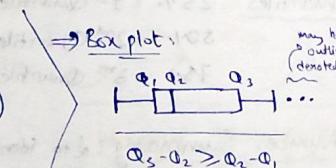
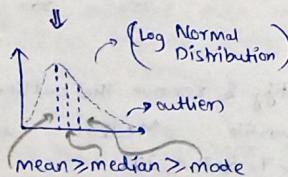
How it's done?  
Kernel density estimator

This smoothening of graph is known as: probability distribution of a continuous variable.

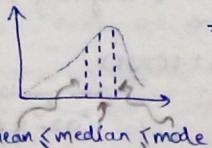
(Mean = Median = Mode) ↪ Right Skewed:



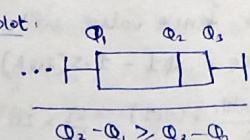
⇒ (+)ve Skewed



Left Skewed:



⇒ (-)ve Skewed → Box plot:



Covariance & Correlation are 2 statistical measures used to determine the relationship btw 2 variables. Both are used to understand how changes in 1 variable are associated with changes in another variable.

Covariance: It's a measure of how 2 random variables change together. If the variables tend to ↑ & ↓ together, the covariance is (+ve). If one tends to ↑ when the other ↓, the covariance is (-)ve.

$$\text{eg: } \left\{ \begin{array}{l} x_1 = y_1 \\ x_2 = y_2 \end{array} \right\} \rightarrow (+) \text{ covariance} \quad \left\{ \begin{array}{l} x_1 = y_1 \\ x_2 = y_2 \\ x_3 = y_1 \\ x_4 = y_2 \end{array} \right\} \rightarrow (-) \text{ covariance} \quad \left\{ \begin{array}{l} x_1 = y_1 \\ x_2 = y_2 \\ x_3 = y_1 \\ x_4 = y_2 \end{array} \right\}$$

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

Adv.: Quantify the relationship btw 'x' & 'y'

Disadv.: Co-var. doesn't have a specific limit value.  
 $\text{Cov}(x, y) \rightarrow -\infty \text{ to } +\infty$

Correlation: Pearson Correlation Coefficient, Spearman Rank Correlation

1) Pearson Correlation Coefficient:  $r_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \cdot \sigma_y}$  range: [-1 to 1]

→ The more the val. towards 1, the more (+)ve correlated 'x' & 'y' are!

→ The more the val. towards -1, the more (-)ve correlated 'x' & 'y' are!

Here we can see:  $x_1 \rightarrow y_1 \in x_2 \rightarrow y_2$ . So in this case pearson corr. must come as +1, but we are getting as 0.88. This is because of the slower growth in middle of graph.

∴ Pearson's Corr. cannot capture correlation for non-linear data.  $\rightarrow \{ \text{L} \rightarrow \text{U} \rightarrow \text{L} \}$

2) Spearman Rank Correlation:

$$r_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \cdot \sigma(R(y))}, R \rightarrow \text{rank}; \text{ ex: } \begin{array}{|c|c|c|c|} \hline x & y & R(x) & R(y) \\ \hline 1 & 2 & 2 & 1 \\ \hline 3 & 4 & 3 & 2 \\ \hline 5 & 6 & 4 & 3 \\ \hline 7 & 8 & 5 & 5 \\ \hline 0 & 7 & 1 & 4 \\ \hline \end{array}$$

Application: feature Selection.

## INTRODUCTION TO PROBABILITY

Mutual Exclusive Event: 2 events are mutually exclusive if they cannot occur at the same time. (it's about determining the likelihood of an event.)

eg: tossing a coin.  $\{H, T\} \quad P(H) = \frac{1}{2}, P(T) = \frac{1}{2}$

$$P(H \text{ or } T) = P(H) + P(T) = \frac{1}{2} + \frac{1}{2} = 1 \quad \{ \text{Additive rule for mutual exclusive events} \}$$

Non-mutual exclusive event: eg: taking a card from a deck.

$$K, Q \sim \{ \text{K, Q, J, A} \} \quad \Rightarrow P(K \text{ or } Q) = P(K) + P(Q) - P(K \text{ and } Q) = \frac{4}{52} + \frac{4}{52} - \frac{1}{52} = \frac{15}{52} \quad \{ \text{non-mutual exclusive event} \}$$

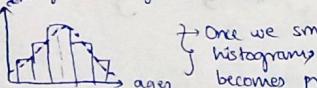
### Multiplication Rule: {Independent & Dependent Events}

→ **Independent:** 2 events are independent if they don't affect one another.  
 ex: tossing a coin {H & then Tails}  $\sim P(H) = \frac{1}{2}, P(T) = \frac{1}{2}$   
 $P(H \text{ and } T) = P(H) \cdot P(T) \therefore \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$

→ **Dependent:** 2 events are dependent if they affect each other.  
 ex: take a king from the deck & then the queen card from the deck  
 $P(K) = \frac{4}{52}, P(Q) = \frac{4}{51}$   
 $P(K \text{ and } Q) = P(K) \cdot P(Q)$  (Conditional probability)  
 $= \frac{4}{52} \times \frac{3}{51}$

## PROBABILITY DISTRIBUTION FUNCTION & TYPES OF DISTRIBUTION

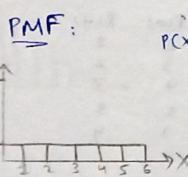
### Probability Distribution Function:

↳ describes the likelihood of different outcomes in a random event.  
 ex: probability ages = {---} → continuous rand. var.  
  
 Once we smoothen the histogram, the y-axis becomes prob. density

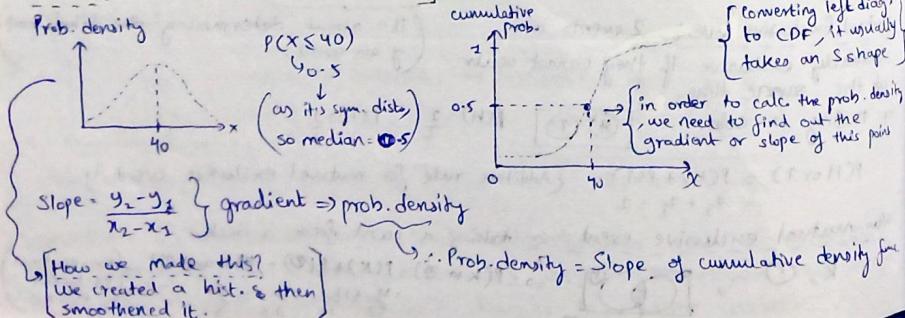
Types of probability dist. func. :  
 ↳ Probability Mass func. (PMF): used for discrete random variables  
 ↳ Cumulative distribution func. (CDF):  
 ↳ Probability Density func. (PDF): used for continuous random variables

### → PMF: [Discrete Random Variables]

ex: rolling a dice  $X = \{1, 2, 3, 4, 5, 6\}$   
 $P(1) = P(2) = \dots = P(6) = \frac{1}{6}$



### → PDF: [Continuous Random Variable]



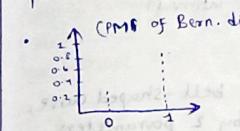
Properties:  
 → Non negativity  $f(x) \geq 0$  for all  $x$   
 (func. which converts prob. density to CDF)  
 → Total area under the PDF curve is equal to 1.  $\Rightarrow \int_{-\infty}^{\infty} f(x) dx = 1$   
 \* wrt diff. distribution,  $f(x)$  func. is going to change.

### Types of probability distribution:

- ↳ Bernoulli dist. → outcomes are binary (pmf) → discrete random variables
- ↳ Binomial dist. → pmf
- ↳ Normal/Gaussian dist. → pdf
- ↳ Poisson dist. → pmf
- ↳ Log Normal dist. → pdf
- ↳ Uniform dist. → pmf

### Bernoulli Distribution:

It's the simplest discrete prob. dist. which represents the probab. dist. of a random variable that has exactly 2 possible outcomes: success (with prob. ' $p$ ') & failure (with prob. ' $1-p$ '). It's used to model binary outcome.



Parameters:  $0 \leq p \leq 1$   
 $q = 1-p$   
 $K = \{0, 1\} \Rightarrow 2 \text{ outcomes}$   
 $P(\text{success}) \propto K+1$   
 $P(\text{failure}) \propto K+2$

$$\text{PMF} = \begin{cases} q = 1-p & \text{if } K=0 \\ p & \text{if } K=1 \end{cases}$$

Mean of bern. dist. = ' $p$ '

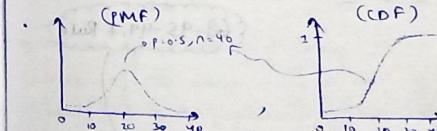
$$\text{Median of bern. dist.} = \begin{cases} 0 & \text{if } p < \frac{1}{2} \\ [0, 1] & \text{if } p = \frac{1}{2} \\ 1 & \text{if } p > \frac{1}{2} \end{cases} \quad \begin{cases} \text{median} = 0 & \text{if } q > p \\ \text{median} = 0.5 & \text{if } q = p \\ \text{median} = 1 & \text{if } q < p \end{cases}$$

$$\text{Mode} = \begin{cases} p & p > q \\ q & q > p \end{cases}$$

$$\text{Variance} (\sigma^2) = pq \quad \text{Std. dev. } \sigma = \sqrt{pq}$$

### Cumulative Density func. (CDF)

It's a discrete prob. dist. that gives the no. of successes in ' $n$ ' independent trials, each with a success prob. ' $p$ ' & failure prob. ' $q$ ' ( $q=1-p$ ). Every outcome of the experiment is binary. It's like collection of bernoulli distribution.



Notation:  $B(n, p)$

Parameters:  $n \in \{0, 1, 2, \dots\} \Rightarrow \text{no. of trials}$   
 $p \in [0, 1] \Rightarrow \text{prob. of success}$   
 $q = 1-p$

Support:  $K \in \{0, 1, 2, \dots, n\} \Rightarrow \text{no. of successes}$

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

$$\text{for } k = 0, 1, 2, \dots, n \text{ where } {}^n C_k = \frac{n!}{k!(n-k)!} \quad \text{binomial coefficient}$$

$$\text{Mean} = np$$

$$\text{Variance} = npq \quad \text{Std. dev.} = \sqrt{npq}$$

Ex: Quality Control: Scenario: inspecting 10 items in a factory where each item has a 10% chance of being defective.

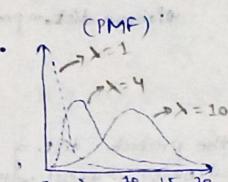
$n=10$ , prob. of success ( $p$ ) = 0.1 (defective item), no. of successes ( $k$ ) = varies from 0 to 10.

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

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

### Poisson Distribution:

It's a discrete prob. dist. that expresses the prob. of a given no. of events occurring within a fixed interval of time or space, assuming these events occur with a const. mean rate ( $\lambda$ ) & independently of the time since the last event. (ex: counting how many cars pass through a checkpoint per hour.)



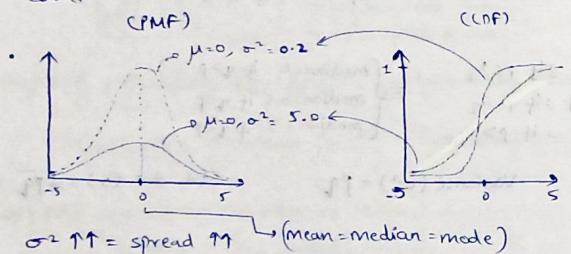
$\lambda$  = expected no. of events occurring at every time interval

$$\text{PMF} \Rightarrow P(X=x) = \frac{e^{-\lambda} (\lambda)^x}{x!}$$

$$\text{Mean} = \text{Variance} = \lambda t \quad (\text{time interval})$$

### Normal / Gaussian distribution:

It's a continuous prob. dist. characterized by its bell-shaped curve, which is symmetric around its mean. It's defined by 2 parameters: the mean (avg.),  $\mu$  & the std. dev. ( $\sigma$ ), which measures the spread of the data.



$$\text{Mean} \Rightarrow \mu = \frac{\sum x_i}{n}$$

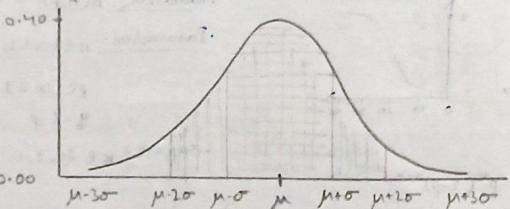
$$\text{Variance} \Rightarrow \sigma^2 = \frac{\sum (x_i - \mu)^2}{n}$$

### Empirical Rule of Normal / Gaussian Distribution:

Let size = 50, assuming the random var. follows →

$x \sim \mathcal{N}(\mu, \sigma^2)$

(Normal/Gaussian dist.)



{68-95-99.7 Rule}

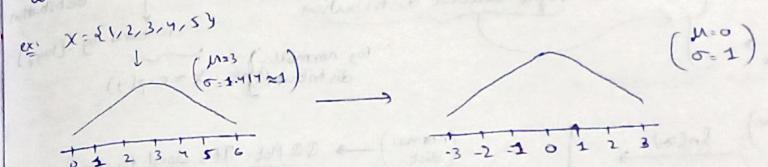
$$P(\mu-\sigma \leq X \leq \mu+\sigma) \approx 68\% \quad (\text{prob. of rand. var. to be btw } \mu-\sigma \text{ & } \mu+\sigma)$$

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

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

### Standard Normal Distribution:

Let's say that I am having a gaussian distribution with some mean  $\mu$  & std. dev.  $\sigma$ . And now I want to convert this distribution into a different distribution where I get mean = 0 & std. dev. = 1. Then this distribution is specifically called as 'standard normal distribution'. For converting, we will use the formula called as Z-Score



$Z\text{-score} = \frac{x_i - \mu}{\sigma}$

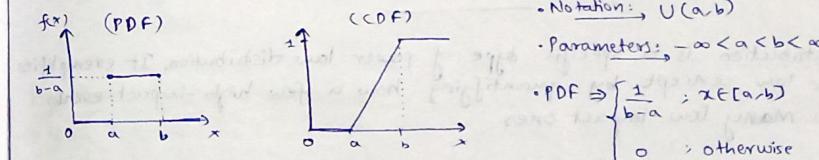
helps in bringing diff. features into 1 scale so that model can improve itself effectively.

Z-score is the result of standardization.

### Uniform Distribution:

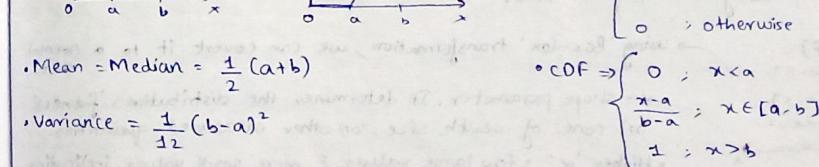
#### 1) Continuous Uniform Distribution : [Continuous Random Variable]

It's a probability distribution where all outcomes in a range are equally likely. ex: picking a random no. btw 0 and 1.

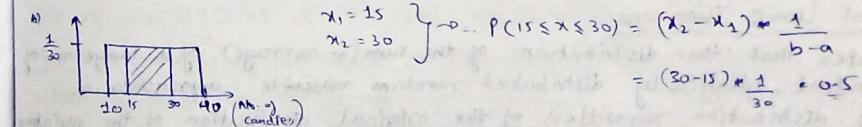


$$\text{Mean} = \text{Median} = \frac{a+b}{2}$$

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



Ex: The no. of candies sold daily at a shop is uniformly distributed with a max. of 40 candies & a min. of 10 candies. Find prob. of daily sales to fall btw 15 & 30?

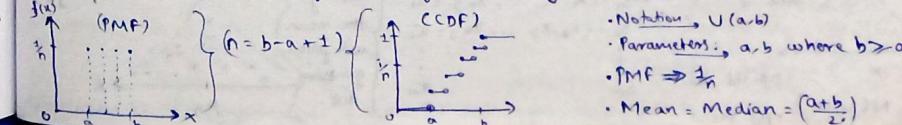


$$x_1 = 15 \quad x_2 = 30$$

$$P(15 \leq X \leq 30) = \frac{(x_2 - x_1)}{b-a} = \frac{30-15}{40-10} = \frac{1}{30} \approx 0.033$$

#### 2) Discrete Uniform Distribution:

It's a prob. distribution where each possible outcome of a discrete random variable has an equal chance of occurring. ex: rolling a dice → each face has prob. of  $\frac{1}{6}$ .



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

$$\text{Parameters} \rightarrow a, b \text{ where } b \geq a$$

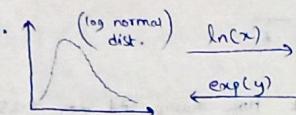
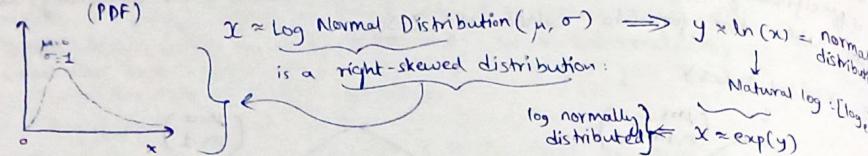
$$\text{PMF} \Rightarrow \frac{1}{n}$$

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

## → Log Normal Distribution:

It's a prob. distribution of a random variable whose logarithm is distributed. This means if you take the natural logarithm of the values, they will follow a normal distribution.

(PDF)

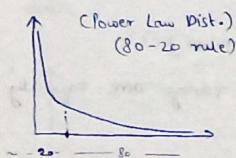


→ QQ Plot: It is used to check whether a distribution is log normal or not.

ex: wealth distribution of the world, salaries of employees in a company.

## → Power Law Distribution:

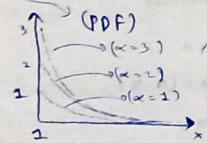
It's a type of statistical distribution where the freq. of an event decreases as the event size increases. In other words, small events are very common, & large events are rare but significant.



ex: 80% of the wealth are distributed with 20% of the total population.



Pareto distribution is a specific type of power law distribution. It exemplifies the power law concept by quantifying how a few high-impact events dominate many low-impact ones.



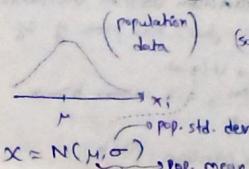
- Using 'Box-Cox' transformation, we can convert it to a normal dist.
- $\alpha$  = shape parameter. It determines the distribution's 'fatness' (or) conc. of wealth, size, (or) other measured quantities.
- Higher ' $\alpha$ ' = few large values & more small values indicating a steeper drop-off.

ex: 80% of the defects can be solved if we solve 20% of the defect.

## → Central Limit Theorem:

It states that the distribution of the sum (or average) of a large no. of independent, identically distributed random variable approaches a normal distribution, regardless of the original distribution of the variables.

### (1) When random variable follows gaussian distribution:



$n$

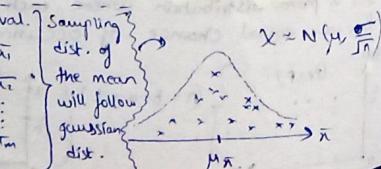
$s_1$

$s_2$

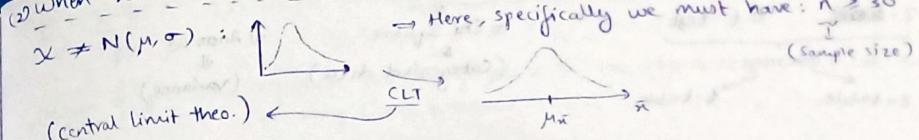
$s_m$

$s_n$

$s_{\bar{x}}$



(2) When random variable doesn't follow gaussian dist.



→ Here, specifically we must have:  $n \geq 30$   
(sample size)

Estimates: It is a specified observed numerical value used to estimate an unknown population parameter.

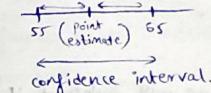
Type:

→ Point Estimate: single numerical value used to estimate an unknown population parameter.

ex: Sample mean is a point estimate of a population mean.

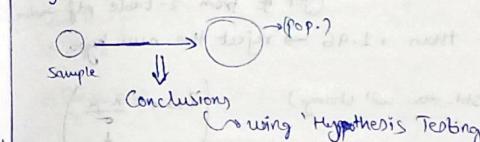


→ Interval Estimate: range of values used to estimate the unknown population parameter.



## INFERENTIAL STATS AND HYPOTHESIS TESTING

The main aim of inferential stats is to come up with conclusions (or) inferences from a sample data about an unknown population parameter.



### Hypothesis Testing Mechanism:

1) Null Hypothesis ( $H_0$ ): The assumption you are beginning with.

2) Alternate Hypothesis ( $H_1$ ): opposite of null hypothesis

3) Experiments: statistical analysis (ex: p-value, significance test).

4) Accept or reject the null hypothesis

### P-Value:

It's a no. that describes how likely you are to have found a particular set of observations if the null hypothesis were true.

Hypothesis testing: ex: Coin is fair or not {100 times}  $\rightarrow P(H) = 0.5, P(T) = 0.5$

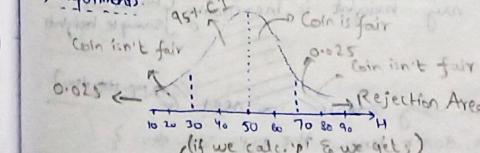
$$P(H) = 0.6, P(T) = 0.4$$

$$P(H) = 0.7, P(T) = 0.3$$

⇒ Null Hypothesis:  $H_0$  = Coin is fair

⇒ Alternate Hypothesis:  $H_1$  = Coin is not fair

3) Experiments:



if we calc.  $p$  & we get:

5) Conclusions: if  $p < \text{significance val.}$  = reject the null hypo., else fail to reject the null hypo.

4) Significance Value ( $\alpha$ ): Let  $\alpha = 0.05$

threshold for determining whether a result is statistically significant

∴ Confidence Interval:  $1 - 0.05 = 0.95$

## Hypothesis Testing & Statistical Analysis:

1) Z-test      2) t-test  
 average  
 z-table      t-table

3) Chi Square  
 (categorical data)  
 ↓  
 (Variance)

4) Anova  
 ↓  
 (Variance)

→ Z-test: We must have: pop. std. dev.  $\approx$  n(sample size)  $\geq 30$

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

(a) State null & alternate hypothesis

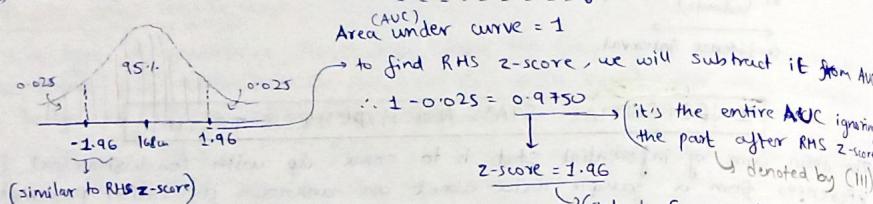
(b) At a 95% confidence level, is there enough evidence to reject the null

Ans)  $\mu = 168\text{cm}$ ,  $\sigma = 3.9$ ,  $n = 36$ ,  $\bar{x} = 169.5\text{cm}$ ,  $CI = 0.95 \rightarrow \alpha = 1 - CI = 0.05$

1) Null Hypo.:  $H_0 \Rightarrow \mu = 168\text{cm}$

2) Alternate Hypo.:  $H_1 \Rightarrow \mu \neq 168\text{cm}$

3) Based on CI, we will draw "Decision Boundary":



If ' $Z'$ ' is less than  $-1.96$  or greater than  $+1.96 \rightarrow$  reject the null hypo.

\* Z-test:  $Z_d = \frac{\bar{x} - \mu}{\left(\frac{\sigma}{\sqrt{n}}\right)}$  (for sample data, std. dev. will change)  $\therefore \left( Z_{\text{score}} = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \right)$   
 $= \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = 2.31$   $\therefore Z_{\text{score}}$  → It tells how many std. dev.'s a data point is from the mean. It's a measure of relative position within a distribution.

∴ Conclusion:  $2.31 > 1.96 \rightarrow$  reject the null hypo.  $\therefore p < 0.05$   $\therefore$  significance value.

Now constructing new graph from above:  $\therefore$  (2nd way of cross-checking p-value)  
 $Z_{\text{score}} \text{ of } 2.31 = 0.98956$  (from z-table pdf)  
 $\therefore AUC (\text{iii})$

$$\therefore AUC = 1 - 0.98956 = 0.01044 \quad \therefore AUC \text{ b/w } [-2.31, 2.31] = 1 - 2 \times 0.01044 = 0.9791$$

$$\therefore p\text{-value} = 0.01044 + 0.01044 \quad \therefore p < 0.05 \quad \therefore \text{reject the null hypo.}$$

Final Conclusion: The avg.  $\neq 168\text{cm}$ .

Now since  $Z = 2.31 > 1.96 \rightarrow$  means my height is increasing based on sample height

→ Student t-distribution:

In Z-stat when we perform any analysis using z-score, we need  $\sigma$  (pop. std. dev.)

How do we perform our analysis when we don't know the pop. std. dev.?  $\therefore$  (Student t-distribution)

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \quad (s = \text{sample standard dev.}) \quad \therefore t\text{-table} \Rightarrow t\text{-test}$$

Degree of freedom: refers to the no. of val.'s in the final calc. that are free to vary (with t-test).  
 $\therefore n-1$  ( $n = \text{sample size}$ )

→ T-test: (this prob. aka one sample t-test)

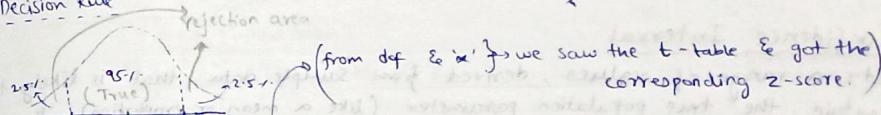
ex: In the pop., the avg. IQ = 100. A team of researchers want to test a new medication to see if it has either a (+ve) or (-ve) effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a std. dev. of 20. Did the medication affect intelligence? (Let  $\alpha = 0.05 \rightarrow CI = 0.95$ )

Ans)  $\mu = 100$ ,  $n = 30$ ,  $\bar{x} = 140$ ,  $s = 20$ ,  $CI = 95\%$ ,  $\alpha = 0.05$

1) Null Hypo.:  $H_0 \Rightarrow \mu = 100$       2) Alt. Hypo.:  $H_1 \Rightarrow \mu \neq 100$  (2 tail test)

3) Significance Level:  $\alpha = 0.05$       4) Dof  $\rightarrow n-1 = 30-1 = 29$   $\therefore$  means it can either be more or less

v) Decision Rule



∴ If t-test is less than  $-2.045$  or greater than  $2.045$ , reject the null hypo. (due to symmetry)

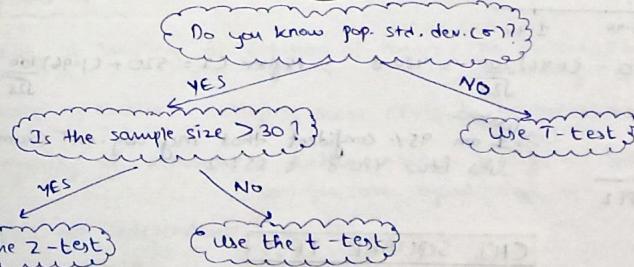
v) Calculate test statistics:  $t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{140 - 100}{\frac{20}{\sqrt{30}}} = 10.96$

Since  $t = 10.96 > 2.045 \rightarrow$  reject the null hypo.

Conclusion: medication used has affected the intelligence.

∴ Since 't' val. is  $> 2.045 \rightarrow$  Medication has increased the intelligence.

When to use Z-test & t-test:



Type I & Type II errors:

Reality: Null Hypothesis is true (or) null hypo. is false.

Decision: Null Hypo. is true (or) null hypo. is false.

↳ (this is based on our result of hypo. testing)

- Outcome 1: We reject the null hypo., when in reality it's false  $\rightarrow$  TYPE I ERROR
- Outcome 2: We reject the null hypo., when in reality it's true  $\rightarrow$  TYPE 2 ERROR
- Outcome 3: We retain the null hypo., when in reality it's false  $\rightarrow$  TYPE 2 ERROR
- Outcome 4: We retain the null hypo., when in reality it's true  $\rightarrow$  GOOD

### Bayes' Theorem:

Probability  $\rightarrow$  Independent events: ex: rolling a dice

Dependent events: ex: picking a king & then picking a queen without placing the card in the deck.

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

$$\text{(Bayes' Theorem)} \quad P(B|A) = \frac{P(B) * P(A|B)}{P(A)} \quad (\text{or}) \quad P(A|B) = \frac{P(A) * P(B|A)}{P(B)}$$

(A, B = events)  $\Rightarrow P(A|B)$  = prob. of A given B is true  $\Rightarrow P(A), P(B)$  = independent prob.  
 $P(B|A)$  = prob. of B given A is true

### Confidence Interval:

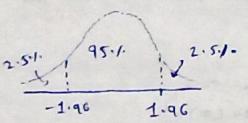
It's a range of values, derived from sample data, that is likely to contain the true population parameter (like a mean or proportion) with a certain level of confidence. For ex., a 95%-confidence interval means we're 95% confident that the true values lies within that range.

$$\cdot CI = \text{Point Estimate} \pm \text{Margin of Error}$$

$$\cdot Z_{\text{crit}} \Rightarrow \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} ; \text{ if we use t-test } \sim \text{replace } (Z_{\alpha/2}) \text{ by } (t_{\alpha/2})$$

ex: On the verbal section of CAT exam, the std. dev. is known to be 100. A sample of 30 test takers has a mean of 520. Construct 95%-CI about the mean.

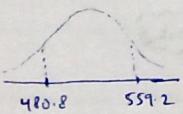
$$\alpha = 0.05$$



$$1 - \alpha/2 = 0.9750 \quad \text{(2-table)} \quad \uparrow$$

AUC wrt that = 1.96

$$\text{Lower CI} = 520 - (1.96) \frac{100}{\sqrt{30}} = 480.8 ; \text{ Higher CI} = 520 + (1.96) \frac{100}{\sqrt{30}} = 554.2$$



Conclusion: I am 95% confident, that my avg. CAT exam score lies btw 480.8 & 554.2

### CHI-SQUARE TEST

The chi-square test for goodness of fit test claims about population proportions. It's a non-parametric test that is performed on categorical [ordinal & nominal] data. ex: there's a population of male who likes diff. color bikes,

Theory	Sample
Yellow ~ $\frac{1}{2}$	28
Red ~ $\frac{1}{2}$	70

(Theory Categorical Distribution) Goodness of fit test: wrt sample info. that we have collected, does it supports the theory? This is what basically we are going to do with the help of chi-square test.

ex: In 2010, census of the city  $\rightarrow$  weights of individuals:  
 In 2020, weight of n=500 individuals were sampled  
 Using  $\alpha=0.05$ , would you conclude the population differences of weight has changed in the last 10 years?

<50kg 50-75 >75  
20% 30% 50%

<50kg	50-75	>75
140	160	200

2010  $\rightarrow$  expected, 2020  $\rightarrow$  observed

expected:

<50	50-75	>75
$0.2 \times 500$ = 100	$0.3 \times 500$ = 150	$0.5 \times 500$ = 250

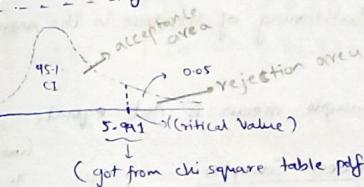
1) Null Hypo.: The data meets the expectation.

Alt. Hypo.: The data does not meet the expectation.

3) Degree of freedom  $\rightarrow$  def =  $k-1$ ; {k = no. of categories (here k=3)}

$$= 3-1 = 2$$

4) Decision Boundary: {Here in chisquare test  $\rightarrow$  we have right skewed graph}



If  $\chi^2 > 5.99 \rightarrow$  reject null hypothesis  
else, we fail to reject null hypothesis

5) Calculate Chi Square Test Statistics:  $\chi^2 = \sum \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}}$

$$\therefore \chi^2 = \frac{(140-100)^2}{100} + \frac{(160-150)^2}{150} + \frac{(200-250)^2}{250} = 26.66 \quad \left. \begin{array}{l} \therefore 26.66 > 5.99 \\ \therefore \text{reject null hypothesis} \end{array} \right.$$

Conclusion: The weights of 2020 population are different than those expected in 2010 pop.

### ANOVA

It's a statistical method used to compare the means of 2 (or) more groups.

1) factors (variable)      2) Levels      3) ex: Mode of payment (factor)  
 (Analysis of Variance)      PAYTM      3) (levels)

### Assumptions in ANOVA

- Normality of Sampling distribution of mean: The distribution of sample mean is normally distributed.
- Absence of Outliers: Outlying score needs to be removed from the dataset.
- Homogeneity of Variance: Population variance in different levels of each independent variable are equal. i.e.  $\sigma_1^2 = \sigma_2^2 = \sigma_3^2$
- Samples are independent & random.

### Types of ANOVA:

- One Way ANOVA: one factor with atleast 2 levels & these levels are independent
- Doctor wants to test a new medication to decrease headache. They split the participants in 3 conditions [10mg, 20mg, 30mg]. Doctor asks the participants to rate the headache [1-10]

Medication $\rightarrow$ Factor	10mg	20mg	30mg
long	5	7	2
3	4	6	

2) Repeated Measures ANOVA: one factor with at least 2 levels & levels are dependent.

e.g.: Running  $\rightarrow$  Factor 3

	Day 1	Day 2	Day 3
Level	8	5	4
	7	4	9
			$\rightarrow$ 2nd person

3) Factorial ANOVA: 2 or more factors (each of which with at least 2 levels, levels can be independent (or) dependent).

e.g.: Running  $\rightarrow$  Factor 3  
(dependent)

Gender  $\rightarrow$  Factor 3 D1 D2 D3

	Male	Female
Independent	8 9 7	5 4 8
Dependent	4 3 3	6 6 3

• Hypothesis Testing in ANOVA: {Partitioning of Variance in the ANOVA}

Null Hypo.:  $H_0 \Rightarrow \mu_0 = \mu_1 = \mu_2 = \dots = \mu_k$

Alt. Hypo.:  $H_1 \Rightarrow$  at least one of the sample mean is not equal.

Test Statistics: F Test =  $\frac{\text{Variance between Samples}}{\text{Variance within Samples}}$

$$\begin{matrix} \text{var. btw. samples} \\ \text{var. within samples} \\ \text{samples} \end{matrix} \quad \begin{matrix} x_1 & x_2 & x_3 \\ 1 & 6 & 5 \\ 2 & 7 & 5 \\ 3 & 3 & 6 \end{matrix}$$

e.g.: Doctors want to test a new medication which reduces headache. They split the participant into 3 cond.'s [15mg, 30mg, 45mg]. Later on, the doctor asks the patient to rate the headache btw [1-10]. Are there any differences btw the 3 conditions using  $\alpha = 0.05$ ?

① Define Null & Alt. Hypo.?:  $H_0: \mu_1 = \mu_2 = \mu_3$   
 $H_1: \text{not all } \mu \text{ are equal.}$

② Significance:  $\alpha = 0.05 \rightarrow \text{CI} = 0.95$

③ Calc. degree of freedom:  $N = 21, a = 3, n = 7$

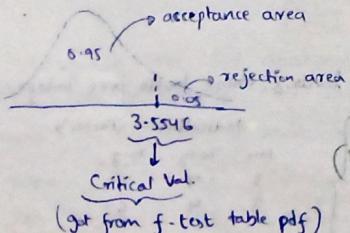
(total sample size)  $\nearrow (7 \times 3)^3$  (total no. of categories)  $\searrow$  (Sample size)

$df_{\text{between}} = a - 1 = 3 - 1 = 2$   $\nearrow$  range:  $(2, 18)$  used in F-table +  $\alpha \rightarrow$  critical value

$df_{\text{within}} = N - a = 21 - 3 = 18$   $\nearrow (2 + 18 = df_{\text{total}})$

$df_{\text{total}} = N - 1 = 20$

④ Decision Boundary: (F-test follows right skewed distribution)



Decision Rule: If F is greater than 3.5546, reject the null hypo.

⑤ Calc. F test Statistics:

We need to fill the following table (sum of squares)

	SS	df	MS	F
Between				
Within				
Total				

$$\begin{aligned} \text{SS}_{\text{between}} &= \frac{\sum (\Sigma a_i)^2}{n} - \frac{T^2}{N} \\ &= \frac{57^2 + 47^2 + 21^2}{7} - \frac{(57+47+21)^2}{21} \\ &= 98.67 \end{aligned}$$

$$\begin{aligned} 15\text{mg} &= 9+8+7+8+9+8 = 57 \\ 30\text{mg} &= 7+6+6+7+8+7+6 = 47 \\ 45\text{mg} &= 4+3+2+3+4+3+2 = 21 \end{aligned}$$

$$\begin{aligned} \text{SS}_{\text{within}} &= \sum y^2 - \frac{\sum (\Sigma a_i)^2}{n} ; \sum y^2 = \text{sum of all values squared} \\ &= 853 - \left( \frac{57^2 + 47^2 + 21^2}{7} \right) \\ &= 10.29 \end{aligned}$$

$$\begin{aligned} &= 9^2 + 8^2 + 7^2 + 8^2 + 9^2 + 8^2 + 7^2 + 6^2 + 7^2 + 8^2 + 7^2 + 6^2 + 4^2 + 3^2 + 2^2 \\ &= 853 \end{aligned}$$

	SS	df	MS = (SS/df)
Between	98.67	2	49.34
Within	10.29	18	0.54
Total	108.96	20	49.88

$$F \rightarrow \frac{\text{Var. btw. Samp.}}{\text{Var. within Samp.}} = \frac{MS_{\text{btw}}}{MS_{\text{within}}}$$

$$\therefore F\text{-test} = \frac{MS_{\text{btw}}}{MS_{\text{within}}} = \frac{49.34}{0.54} = 86.56$$

Since 86.56  $>$  3.5546  $\rightarrow$  we reject  $H_0$ . Conclusion: There are differences btw those 3 conditions.

## DIFFERENTIAL CALCULUS

$$\text{Slope} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{\Delta y}{\Delta x} = \frac{dy}{dx}$$

Interpretation of Slope:

• (+)ve slope: if slope  $> 0$ , the line rises from L  $\rightarrow$  R as it moves. Larger slope = steeper line

• (-)ve slope: if slope  $< 0$ , line falls as it moves L  $\rightarrow$  R. More (-)ve the slope = steeper line in downward dirn.

• Zero slope: if slope = 0, line is horizontal  $\rightarrow$  no vertical change as line moves from L  $\rightarrow$  R.

• Undefined slope: line is vertical  $\&$  slope is undefined.

Slope of eqn. of line:  $y = mx + c$

(slope)  $\downarrow$  (intercept)  $\rightarrow$  point where the line crosses Y-axis.

$$\text{Slope: } \frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1} = \frac{f(x_0+h) - f(x_0)}{x_0+h - x_0} \rightarrow \text{Slope of secant line} = \frac{f(x_0+h) - f(x_0)}{h}$$

$$\therefore \text{Derivative of } f(x) = f'(x) = \lim_{h \rightarrow 0} \frac{f(x_0+h) - f(x_0)}{h} = \frac{d(f(x))}{dx}$$

Power Rule: If  $f(x) = x^n$  where  $n \neq 0 \rightarrow f'(x) = n \cdot x^{n-1}$

$$\frac{d}{dx}(c) = 0 \quad \begin{matrix} \downarrow \\ (\text{constant}) \end{matrix}$$

$$\frac{d}{dx}(c \cdot f(x)) = c \cdot \frac{d}{dx}(f(x)) = c \cdot f'(x)$$

$$\frac{d}{dx}(f(x) + g(x)) = \frac{df(x)}{dx} + \frac{dg(x)}{dx}$$

$$\begin{aligned} f(x) &= \sin(x) \\ f'(x) &= \cos(x) \\ g(x) &= \cos(x) \\ g'(x) &= -\sin(x) \end{aligned}$$

$$f(x) = \cos(x)$$

$$f'(x) = -\sin(x)$$

$$f(x) = \ln(x)$$

$$\therefore f'(x) = \frac{1}{x}$$

$$f(x) = e^x$$

$$\therefore f'(x) = e^x$$

$$\frac{d}{dx}(h(x) \cdot f(x)) = h'(x)f(x) + h(x)f'(x)$$

(Product Rule)

Chain Rule:

If  $y = f(g(x))$  where  $y = f(u)$  &  $u = g(x)$ , then  $\frac{dy}{dx} = \frac{df}{du} \times \frac{du}{dx} = f'(g(x))$

$$\text{ex: } y = \sqrt{\sin(3x)}$$

Step 1: Identify the functions: Outer func. :  $f(u) = \sin u = u^n$   
 Middle func. :  $g(x) = \sin v$   
 Inner func. :  $v = h(x) = 3x$

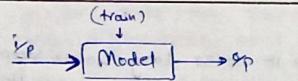
$$\text{Step 2: } f'(u) = \frac{1}{2}(u)^{-\frac{1}{2}} = \frac{1}{2\sqrt{u}} \quad g'(v) = \cos v, \quad h'(x) = 3$$

$$\text{Step 3: Apply the chain rule: } \frac{du}{dx} = \frac{df}{du} \times \frac{du}{dv} \times \frac{dv}{dx} = f'(g(h(x))) \cdot g'(h(x)) \cdot h'(x)$$

$$\therefore \frac{du}{dx} = \frac{3 \cos(3x)}{2\sqrt{\sin(3x)}}$$

- Applications:
  - Backpropagation in neural networks (training DL Models)
  - Gradient Descent Optimization → Feature Engineering
  - Regularization Techniques: Overfitting & Underfitting.

### APPLICATION OF LINEAR ALGEBRA, STATS, & DIFFERENTIAL CALCULUS



(simple linear regression)  $\rightarrow$  (since we have only 1 IP)  
 (we call it 'simple'). If multiple IPs  $\rightarrow$  Multiple LR

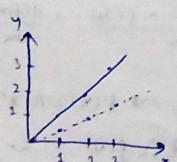
Cost function:  $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$   
 (We call this cost func.)  
 (as Mean Squared Error)

(we have to minimize this)

Our main aim is to basically come up with a best fit line wherein when I try to calculate cost func., it should be minimal.

$$\text{ex: } h_\theta(x) = \theta_0 + \theta_1 x \quad \text{let } \theta_0 = 0 \Rightarrow h_\theta(x) = \theta_1 x$$

x	y
1	1
2	2
3	3



$$\text{Let } \theta_1 = 1 \text{ (line)}$$

$$\therefore h_\theta(x) = 1; x=1$$

$$h_\theta(x) = 2; x=2$$

$$h_\theta(x) = 3; x=3$$

$$\text{Let } \theta_1 = 0.5 \text{ (dotted)}$$

$$\therefore h_\theta(x) = 0.5; x=1$$

$$h_\theta(x) = 1; x=2$$

$$h_\theta(x) = 1.5; x=3$$

Now:  $\theta_0 = 1$   
 $\therefore J(\theta_0, \theta_1) = \frac{1}{2 \times 3} [(1-1)^2 + (2-2)^2 + (3-3)^2] = 0$

$\theta_0 = 0.5$   
 $\therefore J(\theta_0, \theta_1) = \frac{1}{2 \times 3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \approx 0.58$

CPasses via (0,0)

Now:  $\theta_0 = 1$   
 $\therefore J(\theta_0, \theta_1) = \frac{1}{2 \times 3} [(1-1)^2 + (2-2)^2 + (3-3)^2] = 0$

$\theta_0 = 0.5$   
 $\therefore J(\theta_0, \theta_1) = \frac{1}{2 \times 3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \approx 0.58$

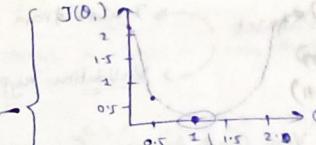
$\theta_0 = 0.5$   
 $\therefore J(\theta_0, \theta_1) = \frac{1}{2 \times 3} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \approx 0.58$

$$\theta_1 = 0$$

$$\therefore J(\theta_0) = \frac{1}{2 \times 3} [(0-1)^2 + (0-2)^2 + (0-3)^2]$$

$$\approx 2.3$$

(Graph when we plot  $J(\theta_0)$  &  $\theta_0$  for diff. values)



Error has been minimized

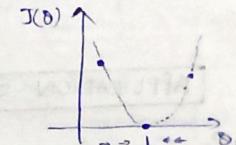
Convergence Algorithm: → optimize the changes of  $\theta_1$  value

. Repeat until Convergence: { (we are going to change  $\theta_1$  value)

$$\theta_j = \theta_j - \alpha \left[ \frac{\partial J(\theta_j)}{\partial \theta_j} \right]$$

It controls the speed at which the convergence should happen. If it's a very very small val., it will take more to converge. If it's a very big val., then it may continuously jump & may never converge. So select a smaller val.  $\approx \alpha = 0.001$

{ we are calculating slope at that point }



if we are on LHS we must go to RHS VERSO

### Multiple Linear Regression

ex: No. of rooms Size of House Price  
 (Independent features) (represents) (Op feature)

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots$$

where,  $\theta_0$  = Intercept;  $\theta_1, \theta_2, \theta_3, \dots, \theta_k$  = co-efficient

Performance Metrics in LR: (Linear Regression)

R-Squared =  $1 - \frac{\text{Sum of square}}{\text{SS total}}$

$$= 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

≈ gives us accuracy.

### Adjusted R-Squared:

ex: Size of House Price

$$R^2 = 75 \Rightarrow 0.75$$

+ No. of bedrooms Price

$$R^2 = 85 \Rightarrow 0.85$$

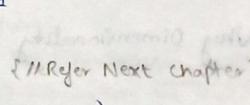
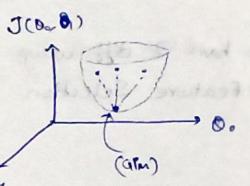
+ Gender Price

$$R^2 = 87 \Rightarrow 0.87$$

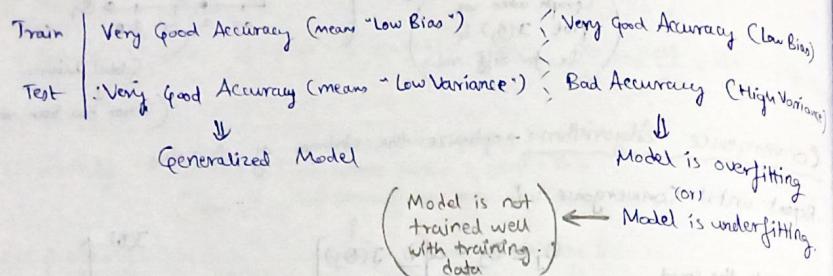
Even though we don't have a feature (Gender) that is directly correlated with the op feature, then also accuracy ↑. To prevent this we use adjusted R-square.

$$\text{Adjusted R Square} = 1 - \frac{(1-R^2)(N-1)}{N-p-1}$$

where N = no. of data points,  
 p = no. of independent features



- Training Dataset → Trains train the model
- Validation: hyperparameter tuning of model.



### APPLICATION OF LINEAR ALGEBRA IN DIMENSIONALITY REDUCTION:

#### Principal Component Analysis (PCA):

1) Curse of Dimensionality: Suppose: 500 features → (~dimensions)

I have many ML models : 3 features      6 features      15 feat.      50 feat.      200 feat.

Acc.1 < Acc.2 < Acc.3 > Acc.4 > Acc.5

We have 2 diff. ways to remove CoD:

- 1) Feature Selection
  - 2) Dimensionality Reduction (PCA)
- ↓  
Feature Extraction

• Why Dimensionality Reduction?  
 ↗ Prevent CoD.  
 ↗ Improve the performance of the model.  
 ↗ visualize & understand the data.

Feature Selection: Helps us select most imp. feat.)

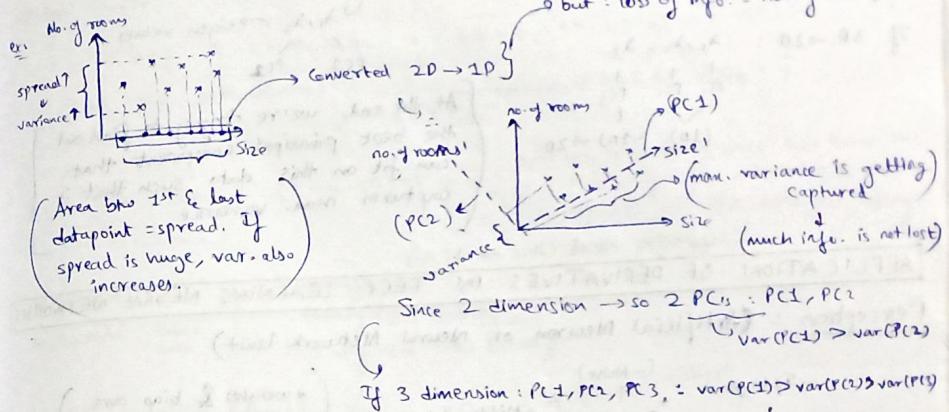
Techniques:  
 1)  $\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1}$  ↗ no → No relationship  
 ↗  $x \perp\!\!\!\perp y$

2) Pearson's Correlation =  $\frac{\text{Cov}(x, y)}{\sigma_x \cdot \sigma_y} = [-1, 1]$  ↗ The more towards +1, the more +ve correlated 'x' & 'y' are. & similar for -1 v/s.

#### Feature Extraction:

ex: Room Size | No. of rooms | Price  
 ↗ 2 feature → 1 feature  
 (Transformation to extract new feature)  
 ↓  
 House Size | Price

PCA Geometric Intuition:  
 → Goal of PCA: find out principal comp. in such a way that max. variance will be captured.



#### Maths Intuition:

$$\text{Proj}_{P_i} u = \frac{P_i \cdot u}{\|P_i\|} = P_i \cdot \frac{u}{\|u\|} \quad \begin{array}{l} (\text{Proj. of } u \text{ on } P_i) \\ (\text{Scalar Val.}) \end{array}$$

$$\therefore \text{Max. Variance} = \sum_{i=1}^n (x_i - \bar{x})^2 \quad \begin{array}{l} (\text{Goal: find the best unit vector which}) \\ \text{captures max. variance}) \end{array}$$

(cost func.)

#### Eigen Vectors & Eigen Values:

- 1) We need to find out covariance matrix b/w features.
- 2) Eigen vectors & Eigen Values will be found out from this covariance matrix.
- 3) Eigen Vector → eigen value (magnitude of eigen vector) → capture the max. var.

$$[ ] \times [v] = \lambda \cdot v$$

↓  
 (suppose we have a matrix)  
 (eigen vector) gives max. magnitude → main eigen vector → best principal component

(already mathematically proven)

#### Steps to calculate eigen vectors & eigen values:

1) Covariance of features :  $\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{N-1}$ ;  $A = \begin{bmatrix} x & y \\ y & z \end{bmatrix}$

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

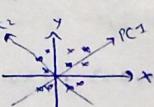
$$\left( \begin{array}{l} \text{Cov}(x, x) = \text{Var}(x) \\ \text{Cov}(y, y) = \text{Var}(y) \end{array} \right)$$

2)  $A \cdot v = \lambda \cdot v$  ↗ eigen values  $\Rightarrow P_1, P_2$

#### Revising:

Given:  $\rightarrow 2D \rightarrow 1D$

1) Standardize the data:  $\rightarrow \lambda_1 \rightarrow \text{mag. of eigen vector}$



$\rightarrow$  Covariance Matrix of  $X \& Y$  :  $A = \begin{matrix} & x & y \\ x & \text{Var}(x) & \text{Cov}(x,y) \\ y & \text{Cov}(y,x) & \text{Var}(y) \end{matrix}$

$\rightarrow$  find out Eigen Vectors & values,  $Av = \lambda v$   $\rightarrow \lambda_1, \lambda_2 \rightarrow$  eigen values  
 $\rightarrow$   $P = \begin{pmatrix} p_{11} & p_{12} \\ p_{21} & p_{22} \end{pmatrix}$

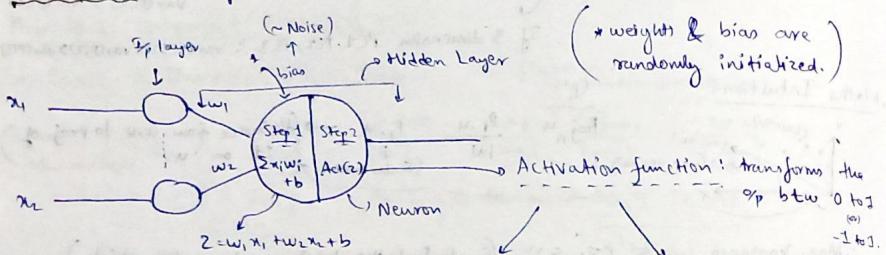
If  $3P \rightarrow 2D$ :  $\lambda_1, \lambda_2, \lambda_3$   
 $\begin{pmatrix} 1 & 1 & 1 \\ p_{11} & p_{12} & p_{21} \\ \sqrt{\lambda_1} & \sqrt{\lambda_2} & \sqrt{\lambda_3} \end{pmatrix}$   
 $(1D) \rightarrow (2D)$

At the end, we're trying to find out the best principal component that can fit on this data such that it captures max. variance.

// Refer Next Chap  $\rightarrow$  PCA

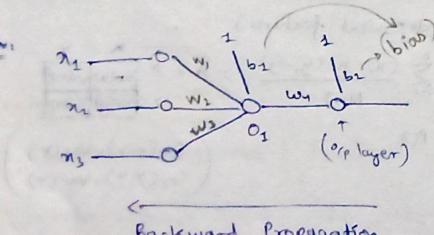
## APPLICATION OF DERIVATIVES IN DEEP LEARNING NEURAL NETWORKS

Perception: (Artificial Neuron or Neural Network Unit)



(Perception is a linear classifier)

Simple Layer Perception Model



Steps:  
 $\rightarrow$  Fwd prop.:  $Act(z)$  (In hidden layer 1)

$$Z = 95 \times 0.01 + 4 \times 0.02 + 4 \times 0.03 + 1 \times 0.001 = 1.151$$

	IQ	Study hrs.	Play hrs.	pass/fail
	95	4	4	1
	100	5	2	1
	95	2	7	0

Let:  $[0.01, 0.02, 0.03, 0.001]$   
 $\begin{matrix} w_1 \\ w_2 \\ w_3 \end{matrix}$

$\rightarrow$  Activation ( $z$ ): Sigmoid =  $\frac{1}{1+e^{-z}}$   $\rightarrow f(z) = \frac{1}{1+e^{-0.751}} = 0.759 = O_2$

Hidden Layer 2: Let  $w_4 = 0.02, b_2 = 0.03$

$$\rightarrow z = 0 + w_4 + b_2 = 0.02 \times 0.759 + 0.03 = 0.04518$$

$$\rightarrow \text{Step 2: Activation } (z) = \frac{1}{1+e^{-(0.04518)}} = 0.51129 = O_2 = \hat{y}$$

(actual exp)  
 $\rightarrow$  Loss func. / Error =  $(1 - 0.51129) \approx 0.49$   
 $\rightarrow$  we have to reduce this error, the only way to minimize is to update all the weights using back propagation.

Regression

MSE, MAE, Huber Loss

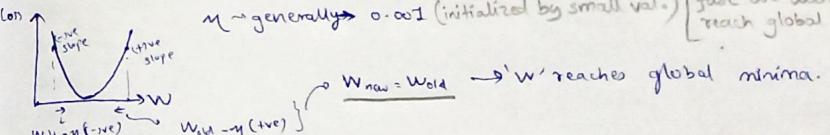
Classification

Binary Cross Entropy,  
Categorical Cross Entropy.

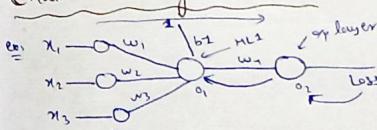
Weight Updation Formula:  $W_{\text{new}} = W_{\text{old}} - \eta \left( \frac{\partial h}{\partial W_{\text{old}}} \right)$  (slope)

Optimizers: To reduce the loss value.

(Learning rate)  $\rightarrow$  decides the step size towards global minima i.e. how fast we want to reach global minima

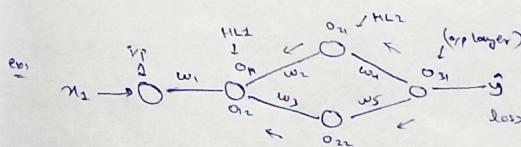


Chain Rule of Derivatives:



$$W_{\text{new}} = W_{\text{old}} - \eta \left( \frac{\partial h}{\partial W_{\text{old}}} \right)$$

$$\frac{\partial h}{\partial W_{\text{old}}} = \frac{\partial h}{\partial O_2} \sim \frac{\partial O_2}{\partial O_{11}} \sim \frac{\partial O_2}{\partial O_{12}} \sim \frac{\partial O_2}{\partial O_{13}}$$



$$W_{\text{new}} = W_{\text{old}} - \eta \left( \frac{\partial h}{\partial W_{\text{old}}} \right)$$

$$\begin{aligned} \frac{\partial h}{\partial W_{\text{old}}} &= \left[ \frac{\partial h}{\partial O_{31}} \times \frac{\partial O_{31}}{\partial O_{21}} \times \frac{\partial O_{21}}{\partial O_1} \times \frac{\partial O_1}{\partial W_{101}} \right] \\ &\quad + \left[ \frac{\partial h}{\partial O_{32}} \times \frac{\partial O_{32}}{\partial O_{22}} \times \frac{\partial O_{22}}{\partial O_1} \times \frac{\partial O_1}{\partial W_{101}} \right] \end{aligned}$$

// Refer. Next Chap  $\rightarrow$  Deep Learning