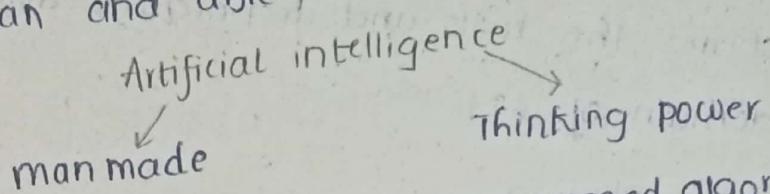


## Introduction to AI:

The AI (artificial intelligence) branch of computer science or science by which we can create intelligent machine, which can behave like a human and think like human and able to make decisions.



- \* AI can create a machine which programmed algorithm which can work with own intelligence.
  - \* AI can not need to use pre program machine to do same work  
(or)
  - \* AI is a branch of computer science which deals with helping machines finding solution to complex problems like human
- Humanbeings  
↓ → convert  
Algorithms  
↓ → implemented  
machine
- \* AI is not related to computer science it is also related to maths, biology, [psychology] psychology etc. etc.  
ex:- automated driving systems (self driving cars)  
To calculate 100 students marks  
\* human take more time but computer can perform very fast.

Need of AI :-

- \* AI can create software or devices which can solve real world problems very easily and accuracy.  
ex:- Health issues, supermarkets, traffic issues.
- \* AI can create your personal virtual assistance such as

Gemini, Google assistance, Siri etc.

- \* AI can create [your] robots it works just like a human.

ex: Sofia, Alexa

Applications of artificial intelligence (AI):

Every branch of science, engineering and technology shares the tools and techniques available in the domain of AI.

#### 1. Game playing :

Game playing is one of the leading domains where AI has been applied with great success.

#### 2. Expert Systems :

- \* An expert system is a software that manipulates encoded knowledge to problems in a specialization domain that normally requires human experts.
- \* An expert system is an AI program in which the system knowledge is obtained from an expert source such that intelligence advice or intelligence decision in solving problems.

#### 3. Natural language processing :

NLP is a technique that builds ability in machines to read and understand the languages that human speaks.

#### 4. Image understanding :

- \* many of AI programs are engineered to solve problem without humans.
- \* 2D array contains grey levels can be used to receive digital images that recognized by video camera.

#### 5. Robotics :

- \* AI is applied in robotics in order to see, hear and react to other sensory stimulation.

Ex: Erica and Sophia.

#### 6. Finance:

- \* The banks we use AI to perform different Operations
- \* Organize Operations (credit, debit)
- \* Invest in stocks
- \* Financial institutions have used AI to detect changes or claims etc.

#### 7. Music:

composition of song, performance sounds and research in music

#### 8. Transportation:

fuzzy logic controllers have been developed for automatic gearbox in automobiles.

ex:- Audi TT

VW TORG

VW corrol

} automatic gearbox

#### 9. Hospitals:

A medical clinic use AI system to organise bed schedule make a rotation, heart sound analysis identify tumors

#### 10. Computer vision:

- \* It is able to extract information from its vision computer vision plays its role here to recognise the object as an image and identify the taste.
- \* The image data only can be in the form of pictures, videos, multidimensional data from a medical scanner or multiple cameras.

#### -Advantages of AI:

High accuracy with less error

High speed

High reliability

Useful for risk data

Digital assistant

Useful as a public utility



## Disadvantages of AI:

- \* High cost
- \* No feelings, emotions
- \* No original creations
- \* can't think out of box.

## → machine Learning :-

- \* machine learning is the branch of artificial intelligence concerned with the design and development of algorithms.
- \* It allows the computer to behave in a way based on the empirical data, such as from electronic sensors or database.
- \* machine learning is a field of studying that gives computers a capability to learn without being explicitly programmed.

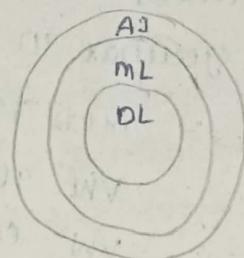
e.g. Online shopping.

(or)

Machine learning is a subfield of artificial intelligence that provides system that ability to automatically and improve from experience without being explicitly programmed.

- Features of machine learning :-
- \* machine learning is used to detect various patterns in a given dataset.
- \* It can learn from past data and improve performance perfectly.
- \* It is a data driven technology
- \* ML is similar to Datamining and it can work large amount of data.
- Types of ML :-

There are three types of machine learning :-



1. supervised learning

2. unsupervised learning

3. Reinforcement learning

1. supervised Learning :-

- \* In supervised learning we teach the machine using label data to predict outcome for new unseen data.
- \* An external super vision is there.
- \* The labelled data means some input data is already tagged with current output.
- \* If any wrong data is given then send feedback to machine.
- \* There are two types of supervised learning algorithm.

1. Regression

2. classification.

a. Regression : It is used for prediction of continuous variables.

ex: price, Gender, age and salary

b. classification : It is used for prediction of Output variables.

ex: true or false, positive or negative, high or low.

2) Unsupervised learning :-

Unsupervised learning is a self learning algorithm with no supervision.

\* The training is provided to the machine with the set of data that not been classified and the algorithms needs to act on that data without any supervision.

\* The input data is unlabelled data.

\* No supervision and no feedback.

\* There are two types of unsupervised learning algorithms:

1. clustering

2. Association

a. Clustering : Grouping of objects into clusters.

b. Association : To find out relationship between two objects.

### 3. Reinforcement Learning :-

- \* It is a feedback based learning approach. Here an agent learn the environment and perform the action to get result of actions.
- \* The reinforcement learning is an example of semi supervised learning because it works based on supervised data and unsupervised learning.
- \* There are two types of reinforcement learning:-
  1. positive reinforcement learning
  2. Negative reinforcement learning

### Applications of machine learning :-

- \* machine learning is a buzzword for today technology, and it is growing very rapidly day by day.
- \* We are using machine learning in our daily life even without knowing it such as Google maps, Google assistant, Alexa etc.
- \* Below are some most trending real-word applications of machine learning:
  1. Image recognition:
    - \* Image recognition is the one of the most common application of machine learning
    - \* It is used to identify objects, persons, places, digital images etc.
  2. Speech recognition:
    - \* While using Google, we get an option of "search by voice", it comes under speech recognition, and it's a popular application of machine learning
    - \* Examples : Google assistant, Siri, Cortana and Alexa
  3. Traffic prediction:
    - \* If we want to visit a new place, we take help of Google maps, which shows us the correct path with the

shortest route and predicts the traffic conditions

\* It predicts the traffic conditions such as whether traffic is cleared, slow-moving, or heavily congested with the help of two ways:

→ Real time location of vehicle from google map app and sensors.

→ Average time has taken on past days at the same time.

#### 4. Product recommendations:

• machine learning is widely used by various e-commerce and entertainment companies such as Amazon, Netflix etc for product recommendation to the user.

#### 5. Self driving cars:

\* One of the most exciting applications of machine learning is self-driving cars.

\* machine learning plays a significant role in self-driving cars.

\* Tesla; for the most popular car manufacturing company is working on self-driving car.

#### 6. Virtual personal assistant:

\* we have various virtual personal assistants such as Google assistant, Alexa, Cortana, Siri.

#### 7. Medical diagnosis:

\* In medical science, machine learning is used for disease diagnosis.

\* with this, medical technology is growing very fast and able to build 3D models that can predict the exact position of lesions in the brain.

#### • Advantages of machine learning:-

\* easily identifies trends and patterns.

\* continuous improvements.

\* No human interaction needed.

\* Handling multi-dimensional data.

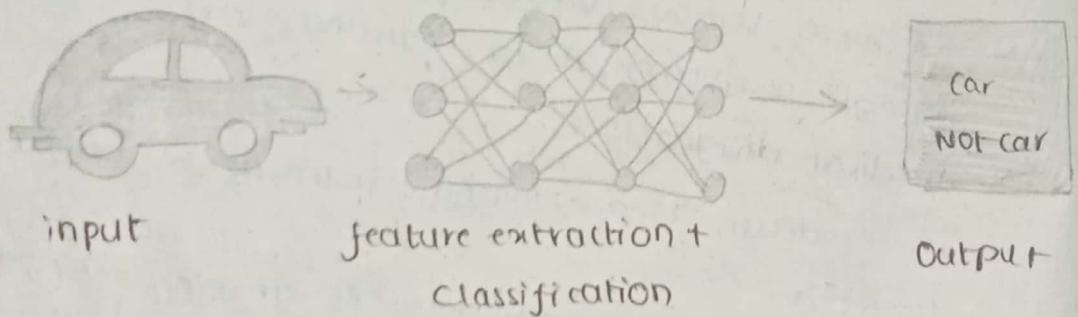
\* wide range of applicability

## Disadvantages of machine learning :-

- \* changing nature of jobs.
- \* Highly expensive
- \* High error chances
- \* results interpretations
- \* Time & resources
- \* Data acquisition.

## ⇒ Deep Learning :-

- \* Deep learning is a synonym for deep[learning] structure hierarchical learning.
- \* It is a subfield of machine learning where supervised, unsupervised & semi supervised learning methods are adopted to learn from data representation.
- \* It contains set of algorithms which have been inspired from the structure and function of human brain.
- \* Deep learning is used to feature extraction



- \* Deep learning will do imitate the human brain.
- \* Deep learning is a subfield of ML for learning feature hierarchies that are actually on artificial neural networks.
- \* Deep learning is implemented by the help of deep networks, which is nothing but neural networks with multiple hidden layers.
- \* Neural Networks : Neural network is also known as artificial neural networks (ANN) simulated

neural network.

- \* most popular techniques used for implementing deep architectures are:

1. Artificial neural networks (or) multilayer perception (or) feed forward network
2. convolutional neural network or deep forward N/W.
3. recurrent Neural network.

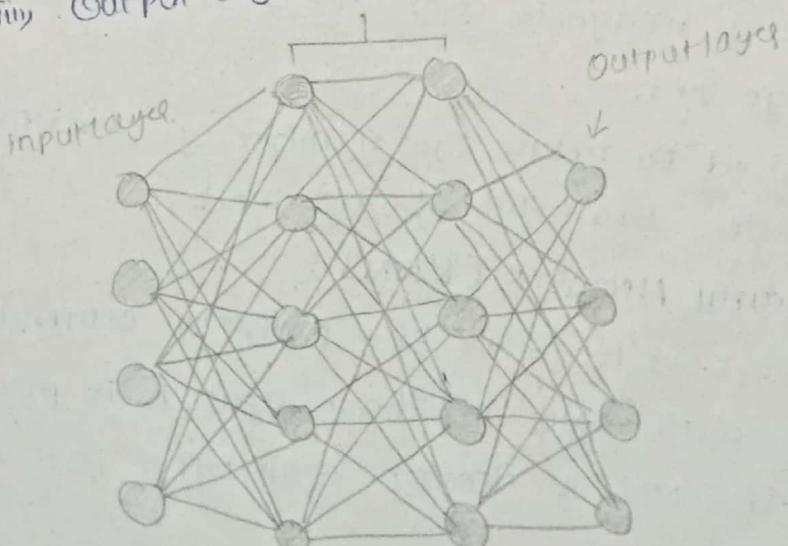
### i) Artificial neural networks :

- \* Artificial neural network is inspired by the working of human brain.
- \* The human brain has neurons interconnected to one another. ANN also have neurons that are interconnected to one another in various layers of the networks.
- \* These neurons are known as nodes or units.
- \* Neural network are a set of algorithm that tries to recognize the patterns, relationships and information from the data.
- \* Artificial neural network is also called as multilayer perception or feed forward network.

### Components of ANN :

A simple neural network consists of three (3) components

- i) input layer
- ii) hidden layer (middle layer)
- iii) output layer



**Input Layer:** input nodes receives inputs/information from the outside world.

**Hidden Layer:**

- \* Hidden layer is a set of neurons where all the computations are performed on the input data.
- \* Any number of hidden layers are used in neural network. But the simple network consists of only one hidden layer.

**Output Layer:**

- \* The output layer is the output/conclusion derived from all the computations performed.
- \* There can be single or multiple nodes in the o/p layer.
- \* If we give binary classification problem the output node is 1.
- \* If we are give input as multi class classification, the output nodes are more than one.

**2) Convolutional neural Networks:**

- \* convolutional neural network (CNN) is similar to a multilayer perceptron network.
- \* The difference is network learns structures and purpose they are mostly used.
- \* CNN can be applied in domain of computer vision problems / image processing etc.
- \* If the given image is color then use three (3) as size of gray image use 1.
- \* The CNN used to reduce the number of parameters & speed up the training of the model.

**3) Recurrent Neural Network (RNN):**

- \* The RNN is used to solve the problem occurs in CNN problems with CNN. we have input & it produce output. In CNN it doesn't maintain internal

- \* Every input is independent from other input
- \* CNN cannot solve problems like sentences, stock prices and time series.
- \* In RNN, each neuron or unit of RNN uses its internal memory to maintain information about the previous O/P. When it required to predict the txt of word of sentence the previous word are required.
- \* The RNN is used to remember the previous output.
- \* The most and main important features of RNN is hidden layers.
- \* The hidden layers have a memory which remembers all the information what has been calculated.
- \* It can also reduce the complexity of parameters.

### Applications of Deep Learning :-

- \* Automatic colouring of black and white images.
- \* automatically adding sounds to silent movies.
- \* automatically adding text to photographs.
- \* automatic machine translation.
- \* object classification & detection in photographs.
- \* automatic text generations.
- \* recommendation engine.
- \* chatbots and speech recognition.
- \* Image recognitions, IoT, Computer Vision.

### → main challenges of machine learning :-

- \* There are a lot of challenges that machine learning professionals face to inculcate ML skills and create an application from scratch.
- 1. poor quality of Data:  
Data plays a significant role in the machine learning process.
- \* One of the significant issues that machine learning professionals face is the absence of good quality data. Unclean and noisy data can make the whole process extremely exhausting.
- \* we don't want our algorithm to make inaccurate or faulty predictions, hence the quality of data is

essential to enhance the output.

## 2. Underfitting Of training data :-

- \* The process occurs when data is unable to establish an accurate relationship between Input and Output variables.
- \* To Overcome this issue:
  - i. maximizing the training data.
  - ii, enhance the complexity of the model.
  - iii, Add more features to the data.
  - iv, reduce regular parameters.
  - v. Increasing the training time of model.

## 3. Overfitting of training data:-

- \* An overfitting occurs when the low bias and high variance. That means the model training is very well but testing is very low.
- \* Training is properly done and but testing is not properly done.

### How to avoid the Overfitting :-

- a) cross validation
  - b) Training with more data
  - c) removing features
  - d) early stopping of training
- \* Both overfitting and underfitting cause the degraded performance of the machine learning model.

## 4) Lack of training Data:

- \* The most important task you need to do in the machine learning process is to train the data to achieve an accurate output.
- \* less amount training data will produce inaccurate results.

## 5. slow implementation:

- \* This is one of the common issues faced by machine learning professionals.

- \* The machine learning models are highly efficient in providing accurate results.
  - \* But it takes more time to provide accurate results.
- Difference between artificial intelligence vs machine learning vs Deep learning.

artificial intelligence	machine learning	Deep learning
AI stands for artificial intelligence, and is basically the study/ process which enables machines to mimic human behaviour through particular algorithm.	ML stands for machine learning, and is the study that uses statistical methods enabling machines to improve with experience	DL stands for deep learning is the study that makes use of neural networks (similar to neurons present in human brain) to imitate functionality just like a human brain
AI is the broader family consisting of ML and DL as its components	ML is the subset of artificial intelligence	DL is the subset of machine learning
AI is a computer algorithm which exhibits intelligence through decision making	ML is an AI algorithm which allows system to learn from data	DL is a ML algorithm that uses deep (more than one layer) neural networks to analyze data and provide output accordingly.
Search trees and much complex math involved in AI	if you have a clear idea about the logic involved in behind and you can visualize the complex functionalities like k-means, support vector machines, etc. then it defines the ML aspect	if you are clear about math involved in it but don't have idea about the features. So you break the complex functionalities into linear/lower dimension features by adding more layers, then it defines DL aspect

The aim is basically increase chances of success and not accuracy	The aim is to increase accuracy not caring much about the success ratio	It attains the highest rank in terms of accuracy when it's trained with large amount of data
Three broad categories / types of AI are: artificial narrow intelligence (ANI), artificial General intelligence (AGI) and artificial SUPER intelligence (ASI)	Three broad categories / types of ML are: Supervised learning Unsupervised learning reinforcement learning	DL can be considered as neural networks with large numbers of parameters layers lying in one of the four fundamental network architectures: Unsupervised Pretrained networks, convolutional neural networks, recurrent neural networks & recursive NN.
The efficiency of AI is basically the efficiency provided by ML and DL respectively	less efficient than DL as it can't work for longer dimensions or higher amount of data.	more powerful than ML as it can easily work for larger sets of data.
Examples of AI: Google's AI-powered predictions, ride sharing Apps like Uber and Lyft, Commercial flights, Use of AI autopilot etc.	Examples of ML: Virtual personal assistant: Siri, Alexa, Google etc, Email spam and malware filtering	examples of DL: sentiment based news aggregation, image analysis and caption generation etc.
AI systems can be rule based, knowledge based or data driven	In reinforcement learning, the algorithm learns by trial and error, receiving	DL n/w consists of multiple layers of interconnected neurons that process

feedback in the form of rewards or punishments

data in hierarchical manner, allowing them to learn increasingly complex representation of data.

Q difference between supervised and unsupervised learning:

parameters	Supervised Learning	Unsupervised Learning
Definition	Learning from labeled data to predict Outputs	Learning from unlabeled data to find patterns.
Data types	requires labeled data (input - output pairs)	uses unlabeled data (only input features)
Objective	predict outcomes for new unseen data	discover hidden pattern or structure in the data.
Examples	Email spam detection, house price prediction	customer segmentation, anomaly detection
Techniques Used	regression, classification algorithms (e.g. decision trees, SVM)	clustering, dimensionality reduction (e.g. K-means, PCA)
Output	Predicts specific outputs (labels for values)	identifies groups, clusters or patterns.
Dependency on labels	highly dependent on labeled data	No labels required for training
Complexity	(General) Generally simpler as labels guide the learning process	more complex as it involves finding structure on its own
Real-world Use cases	Fraud detection, medical diagnosis	Social network analysis, segmentation

## ⇒ Sampling distribution of an Estimator [Population]:

### • Population :

population is the totality of statistical data forming subject of investigation.

### • Size :

\* The number of observations in the populations is defined to the size of population. It may be finite or infinite.

\* size of the population denoted by  $N$

\* size of the sample denoted by ' $n$ '

### • Sample :

\* most of the times study of entire population may not be possible to carry out and hence a part alone is selected from the given population.

### Sampling :

The process of selecting the sample is called sampling Sampling distribution!

\* A sample distribution is a probability of a statistics obtained from a larger number of samples drawn from a specific population.

\* Sample statistic Only estimate population parameters, mean, standard derivation.

\* The sample mean will be different to the population mean

\* A researcher will never know the exact amount of sampling. But using sampling distribution they can estimate the sampling error.

1. Number of samples with replacement =  $N^n$

2. Number of samples without replacement =  $NC_n$

\* Samples are classified in two ways :

i. Large sample :  $n > 30$

ii. Small sample :  $n \leq 30$

### Parameters :

population related mean(μ), standard deviation(σ) variance( $\sigma^2$ ), median etc are called parameters.

statistics: sample related mean( $\bar{x}$ ), standard deviation ( $s$ ), variance ( $s^2$ ), median etc are called statistics.

sample mean:

If  $x_1, x_2, \dots, x_n$  represent a random sample of size 'n'

then the sample mean is defined by  $\bar{x}$ .

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

sample variance:

If  $x_1, x_2, \dots, x_n$  represents a random sample of size 'n'

then the sample variance is defined by  $s^2$ .

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

sample error:

- \* The difference between sample mean statistic and the population parameter is called sample error.

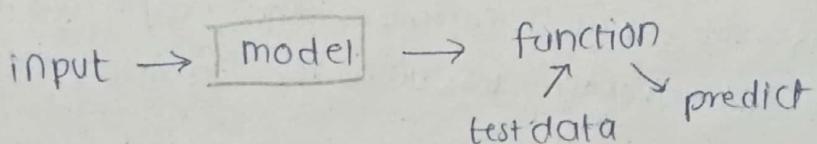
Estimator:

- \* A tool in ML that uses data to make predictions or classify new instances.

⇒ Training and testing loss:

tip amount prediction

Bill amount	Tip amount
500	45
550	50
600	55
500	45
560	50
650	65
600	60
550	50
500	48
600	50



linear regression:

i/p → model linear → function  
e.g. offline

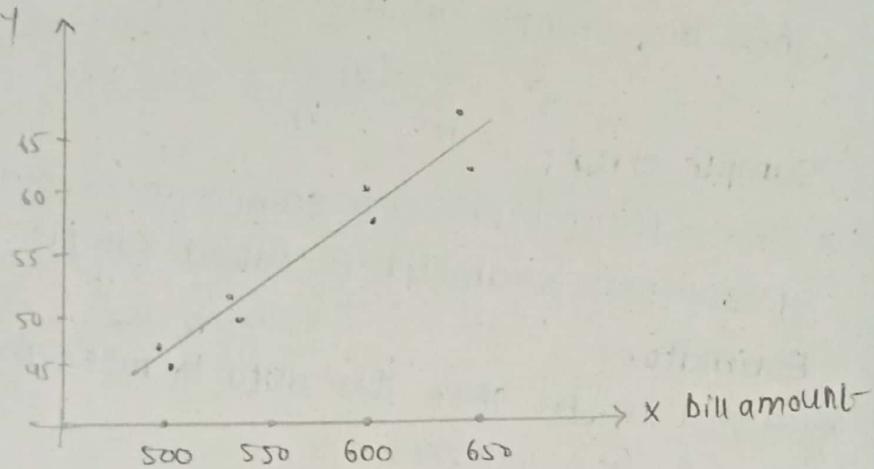
$$\hat{y} = 0.1x + 2$$

$$\text{loss} = \sum_{i=1}^n \frac{(y - \hat{y})^2}{n}$$

for testing :  $T_{\text{test}} = \frac{(7)^2 + (4)^2 + (12)^2}{3}$

 $T_{\text{train}} = \frac{(7)^2 + (7)^2 + \dots}{7}$

Graph: tipamount



\* Training loss and test loss are used to measure how well a model performs.

Training loss: A measure of how well a model performs on the training data. It is used to guide the model's learning process.

Test loss: A measure of how well a model performs on new unseen data. It's an independent measure of the model's generalisation performance.

⇒ Empirical risk management (ERM) minimization:  
 Empirical risk minimization is a fundamental principle in statistical learning theory used to find a predictive model by minimizing the error or loss on a given training dataset. The objective is to choose a hypothesis from a hypothesis set that minimizes the empirical risk, which is the average loss on the training data.

- key concepts:

- Hypothesis Set: A set of possible functions (models) from

which we aim to choose the best one.

Loss function: measures the error between the predicted and actual values.

Common example include mean squared error for regression and cross-entropy for classification.

Empirical risk: The average loss over the training samples if the training set consists of  $N$  samples and loss function is  $L$ , the empirical risk is given by:

$$R_{\text{emp}}(h) = (1/N) \sum L(h(x_i), y_i)$$

where  $h(x_i)$  is the prediction of the model 'h' for the input  $x_i$  and  $y_i$  is the actual output

Example:

Let's consider a simple example with linear regression.

1. Data: suppose we have a training set with three data points:

$$\{(1, 2), (2, 2.5), (3, 3.5)\}$$

2. Hypothesis set: We consider linear function of the form

$$h(x) = w_1 x + b$$

3. Loss function: mean squared error (MSE):

$$L(h(x_i), y_i) = (h(x_i) - y_i)^2$$

4. Empirical risk calculation:

suppose our model is  $h(x) = 1.2x + 0.5$

calculate predictions:

$$n(1) = 1.2(1) + 0.5 = 1.7$$

$$n(2) = 1.2(2) + 0.5 = 2.9$$

$$n(3) = 1.2(3) + 0.5 = 4$$

calculate empirical risk:

$$R_{\text{emp}}(h) = \frac{1}{3} [(1.7 - 2)^2 + (2.9 - 2.5)^2 + (4.1 - 3.5)^2]$$

$$R_{\text{emp}}(h) = \frac{1}{3} (0.09 + 0.16 + 0.36)$$

$$= \frac{1}{3} (0.61) = 0.203$$

$$\therefore R_{\text{emp}}(h) = 0.203$$

Why ERM?

practical : provides a concrete criterion to select the best model based on observed data.

Theoretical foundation: forms the basis for many machine learning algorithms, ensuring they generalize well to unseen data.

Flexibility : can be applied to various types of loss functions and hypothesis tests.

- \* ERM is a core concept in machine learning and helps in building models that perform well on training data, hoping they will also perform well on new, unseen data.

Advantages of Empirical risk minimization:

- \* Foundation of supervised learning: ERM is a core principle in supervised learning, providing a systematic approach to model training and evaluation.
- \* Flexibility: ERM can be applied to various types of models and loss functions, making it versatile across different problem domains.
- \* Optimization framework: provides a clear objective for optimization minimization the empirical risk (arg min loss on training data).
- \* Generalization: when combined with techniques like cross-validation and regularization, ERM help models generalize well on unseen data.
- \* Adaptability: can be applied to different learning algorithms from linear regression to deep neural networks, by choosing appropriate loss functions and hypothesis sets.

Disadvantages Empirical risk minimization:

1. Overfitting: focusing solely on minimizing empirical risk may lead to overfitting, where the model performs well on training data set but poorly on new data.

- 2. Dependence on training data quality : ERM relies heavily on the quality and representativeness of the training data. If the data is biased or noisy, the model's performance will suffer.
- 3. Computational complexity : For large datasets or complex models, minimizing empirical risk can be computational-ally intensive, requiring significant resources.
- 4. Choice of loss function : The effectiveness of ERM depends on the choice of loss function, which can be challenging to select and may not always capture the true cost of prediction errors.
- 5. Local minima : In non-convex optimization problems ERM can get stuck in local minima, leading to suboptimal models.

#### ⇒ Estimating Risk Statistics :

- Risk management plays a crucial role in ML of building effective and reliable.
- At a high level the statistical risk measure the quality of the learning algorithm.
- In ML risk occurs when training data is more and testing data is less otherwise training data is less and testing data is more.

#### Techniques of risk statistics :

- Empirical risk minimization (ERM)
- [Regulation] regularization risk minimization (RRM)
- Cross validation
- Loss function etc

#### Empirical risk minimization [ERM]

- As we don't know parameter 'p' we can't compute the true risk but we can compute the empirical risk based on  $x_i, y_i$  where  $i = 1, 2, 3, 4, \dots, n$
- It is a statistical learning algorithm used to find the optimal solution out of a set of possible solutions based on sample data.



$$L_{\text{emp}}(h) = \frac{1}{N} \sum_{i=1}^N L(y_i, h(x_i))$$

- \* It is flexible and easy to implement.
- \* Risk occurs when there is a Overfitting and Underfitting in data models.

- Overfitting : It occurs when there is a more training data.
- Underfitting : It occurs when there is a more testing data than training data.

- a. [Regulation risk] Regularization risk minimization:
  - \* It is a machine learning technique that uses regularization to reduce overfitting in data models.
  - \* By using set of methods (maths, statistics, computer etc.) that trade a small decrease in training accuracy for better generalizability.
  - \* There are three primary regularization methods they are :

L<sub>1</sub> regularization

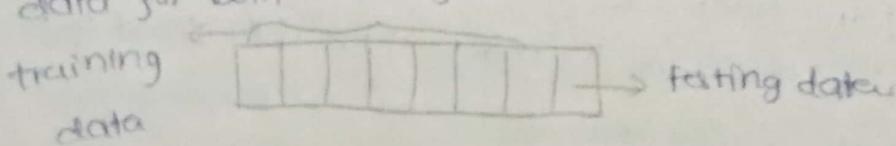
L<sub>2</sub> regularization

Elastic net

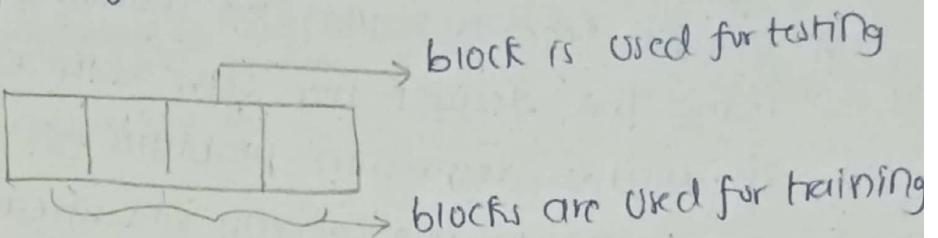
- L<sub>1</sub> regularization: the process of adding sum of absolute values to the loss function.
  - \* It is also called "Lasso".
- L<sub>2</sub> regularization: The process of adding sum of square values to the loss function.
  - \* It is also called "ridge".
- Elastic net: The combination of L<sub>1</sub> regularization (lasso) and L<sub>2</sub> regularization (ridge) is called elastic net.

### b) cross validations:

In machine language, it is not suitable to use same data for both training and testing data.



Cross validation is the mechanism to identify how much data is used for training and how much data is used for testing to get an accurate results.



Formula:

$$CV_k = \frac{1}{n} \sum_{i=1}^n \text{loss}(g_T - e(x_i), (x_i, y_i))$$

- \* consider a dataset  $T$  of size  $n$  is divided into  $k$  folds
- \* loss function of training set  $\ell(g_T)$

#### 4. Loss function:

- \* loss function quantifies the difference between the predicted values and actual values.
- \* it is also known as error function (or) cost function

Categories of loss functions:

- classification model
- regression model

Regression model: It is a powerful tool for understanding & predicting relationships between the variables.

Classification model:

Classification loss quantifies the error between prediction class labels and true class labels.

## Unit 2

⇒ Decision tree Induction (or) Decision tree:

- \* Decision tree is supervised learning technique and by using this decision tree solve both classification problems as well as regression problems.

- \* Decision tree will represents classification model and regression model in the form of tree structure.
- \* Decision tree contains 3 types of nodes:

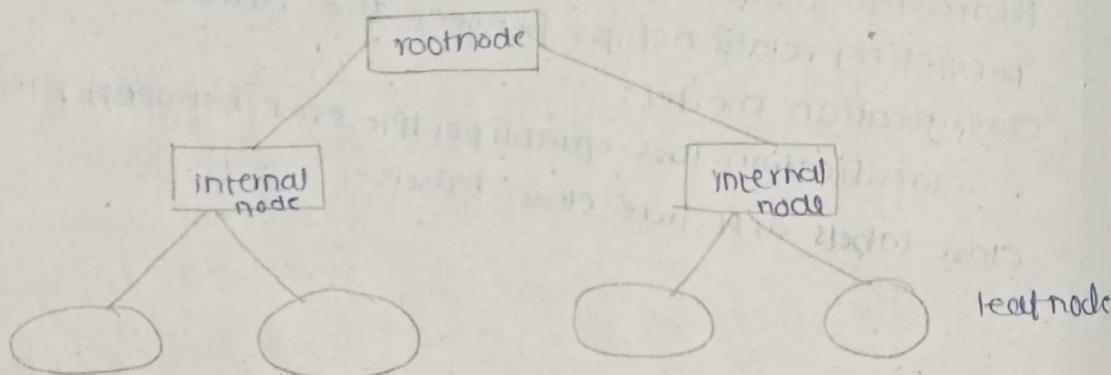
1. root node
2. internal node
3. leaf node.

1. Root node: It has no incoming edge and zero or more outgoing edges.

2. Internal node: It has exactly one incoming edge and two or more outgoing edges.

3. Leaf node: It has exactly one incoming edge and no outgoing edges.

\* each leaf node represents a class.



\* Basically a decision tree are two types:

1. regression tree
2. classification tree

1. regression tree:

\* For continues quantitative target variables  
ex: predictive, revenue, predicting rainfalls, predicting marks etc.

2. classification tree:

\* For discrete categorical target variables  
ex: predicting positive or negative, predicting yes or no

- \* There are two types of algorithm in classification tree
  1. ID tree (iterative discodemizers)
  2. CART (classification regression tree)
- \* In ID tree algorithm by using entropy we built classification tree.
- \* In CART algorithm by using GINI index and impurity measure
- Entropy : Entropy is measure of randomness or disorder -ness (or)
- \* Entropy measures the impurity of a collection of examples or points. It depends from the distribution of the random variables ( $P$ ).

Formula:

$$\text{Entropy}(S) = -P_+ \log_2 \frac{(P_+)}{P_-} - P_- \log_2 \frac{(P_-)}{P_+}$$

Here  $S$  = collection of training examples / records  
 $P_+$  = it is a proportion of positive examples in ' $S$ '.  
 $P_-$  = it is a proportion of negative examples in ' $S$ '.

- Information gain:
- \* Given entropy as a measure of the impurity in a collection of training examples, the information gain, is simply the expected reduction in entropy caused by partitioning the examples according to an attribute.
- \* The information gain,  $\text{Gain}(S, A)$  of an attribute ( $A$ ) relative to a collection of examples ( $S$ ) is defined as

$$\text{Gain}(S, A) = \text{entropy}(S) - \sum_{\text{values } v(A)} \frac{|S_v|}{|S|} \text{entropy}(S_v)$$

Here,  $S_v \rightarrow$  subset of  $S$

$v(A) \rightarrow$  set of all possible values for the attribute  $A$ .

Decision tree algorithm - ID3 solved examples.

Instance	Classification	a1	a2
1	+	T	T
2	+	T	T
3	-	T	F
4	+	F	F
5	-	F	T
6	-	F	T

Attribute : a<sub>1</sub>

We have 6 examples or points or records. There are class labels + (+), -.

$$S[3+, 3-]$$

$$\text{Entropy}(S) = -P_+ \log_2(P_+) - P_- \log_2(P_-)$$

$$= 1.0$$

$$\therefore \text{entropy}(S) = 1.0$$

1st attribute : True values {+, -}

$$S[2+, 1-]$$

$$\therefore \text{entropy}(S_T) = -P_+ \log_2(P_+) - P_- \log_2(P_-)$$

$$= -\frac{2}{3} \log_2(\frac{2}{3}) - \left(\frac{1}{3}\right) \log_2\left(\frac{1}{3}\right) = 0.9183$$

2nd attribute : False, values {+, -, -}

$$\therefore S[1+, 2-]$$

$$\therefore \text{entropy}(S_F) = -P_+ \log_2(P_+) - P_- \log_2(P_-)$$

$$= -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = 0.9183$$

$$\therefore \text{Information Gain}(S, a_1) = \text{entropy}(S) - \sum \frac{|S_r|}{|S|} * \text{entropy}(S_r)$$

$$= 1 - \frac{3}{6} \text{entropy}(S_T) - \frac{3}{6} \text{entropy}(S_F)$$

$$= 1 - \frac{3}{6} (0.9183) - \frac{3}{6} (0.9183)$$

$$= 1 - 0.45915 - 0.45915 = 1 - 0.9183$$

$$= 0.0817$$

Attribute : a<sub>2</sub>

1st attribute : True values = {T, T, F, F}

$$\therefore S[3+, 3-] = 1$$

$$S_T = [2+, 2-]$$

$$\therefore \text{entropy}(S_T) = 1$$

$$\text{2nd attribute: } S_F = [1+, 1-]$$

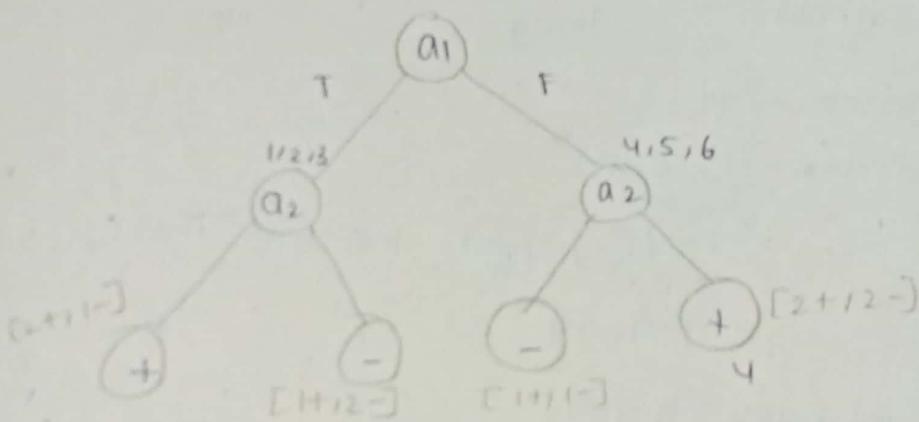
$$\text{entropy}(S_F) = 1$$

$$\text{information gain}(S, a_2) = \text{entropy}(S) - \frac{\sum (S_v)}{|S|} \cdot \text{entropy}(v)$$

$$= 1 - \frac{4}{6} \times 1 - \frac{2}{6}(1)$$

$$= 1 - 0.6666 - 0.3333$$

$$= 0.0001$$



→ GINI Index:

- The GINI index is used in CART algorithm
- Decision tree algorithms are famous decision for classification and regression tree.
- GINI index is also called as GINI coefficient or GINI impurity.
- GINI coefficient is a significant impurity measure utilized in decision tree algorithm.
- The GINI index is a proportion of impurity or inequality or in statistical settings
- They will investigate the idea of GINI index exclusively it's a numerical formula and it's applications in machine learning.

$$\text{Formula: GINI index} = 1 - \sum_{i=1}^k (P_i)^2$$

here  $P_i = i^{\text{th}}$  probability

- Identify the attribute that will act as the root node of decision tree to predict g+t play for following data base with index. Indicate an intermediate step &

outlook	wind	playgolf
Rain	strong	No
sunny	weak	yes
Overcast	weak	yes
Rain	weak	yes
sunny	strong	yes
rain	strong	No
Overcast	strong	No

$$\text{Gini Index} = [4+13-]$$

$$\begin{aligned}\text{Gini Index} &= 1 - \sum_{i=1}^n (p_i)^2 \\ &= 1 - \left(\frac{4}{7}\right)^2 - \left(\frac{3}{7}\right)^2 = 1 - 0.3265 - 0.1836 \\ &= 0.4899\end{aligned}$$

$$\therefore \text{Gini Index} = 0.4899$$

Attribute 1: Outlook.

values = {rain, sunny, Overcast}

$$\text{Srain} = [1+32-]$$

$$\text{Gini}(\text{rain}) = 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 = 0.4444$$

$$\therefore \text{Gini}(\text{sunny}) = [2+0-]$$

$$\therefore \text{Gini}(\text{sunny}) = 0$$

$$\rightarrow \text{Gini}(\text{Overcast}) = [1+1-] = 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 = 0.5$$

$$\therefore \text{Gini index (outlook)} = \frac{3}{7} (\text{Gini}(\text{rain})) + \frac{2}{7} \text{Gini}(\text{sunny}) + \frac{2}{7} (\text{Gini}(\text{Overcast}))$$

$$= \frac{3}{7} (0.4444) + \frac{2}{7} (0) + \frac{2}{7} (0.5)$$

$$= 0.3333$$

Attribute 2: wind

values = {strong, weak}

$$\therefore \text{Gini index}(\text{strong}) = [1+13-]$$

$$= 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 1 - 0.0625 - 0.75$$

$$= 0.375$$

$$\therefore \text{Gini}(\text{wind}) = \frac{4}{7} \times \text{gini(strong)} + \frac{3}{7} (\text{gini(weak)})$$

$$= \frac{4}{7} (0.375) + \frac{3}{7} (0)$$

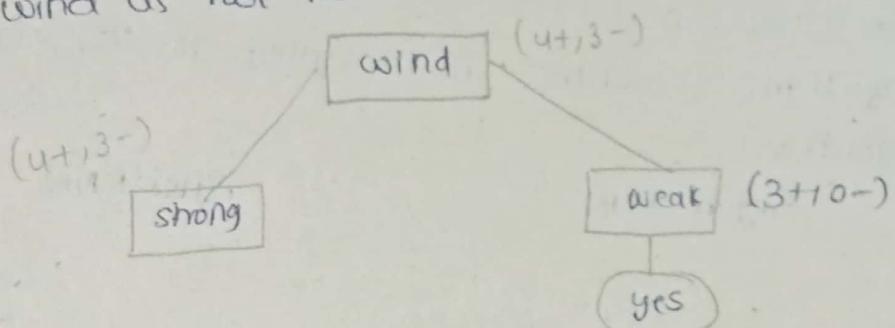
$$= 0.2142$$

from two attributes;

$$\text{Gini}(\text{outlook}) = 0.3333$$

$\text{Gini}(\text{wind}) = 0.2142$

∴ here  $\text{Gini}(\text{wind})$  is smaller. hence we consider the wind as root node.



outlook	wind	play golf
rain	strong	no
sunny	strong	yes
Rain	strong	no
Overcast	strong	no

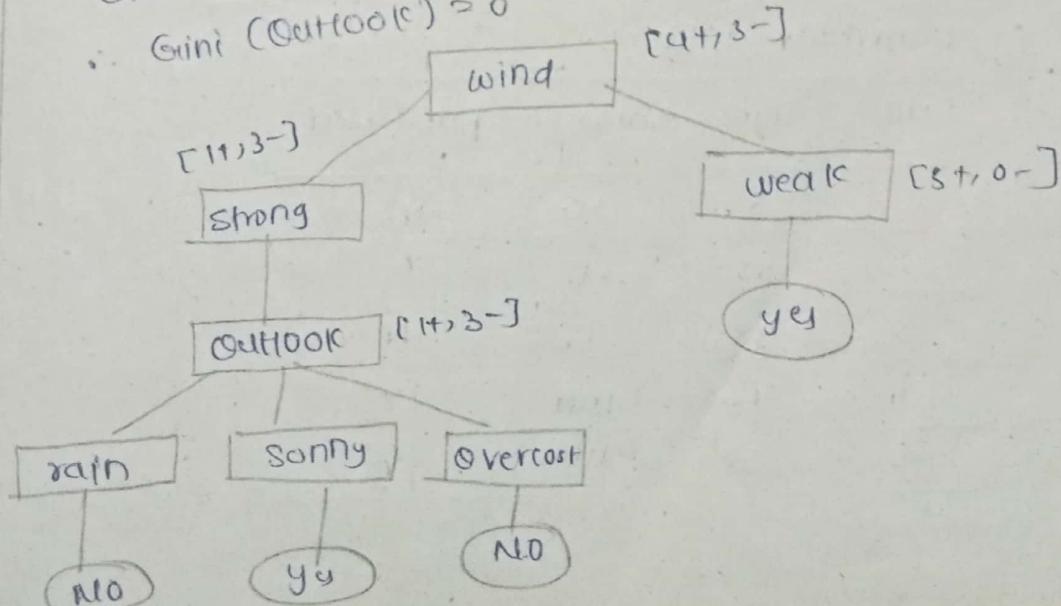
Attribute 1 : outlook , values = {rain, sunny, Overcast}

$$\text{Gini}(\text{rain}) = [0+1, 2-] = 0$$

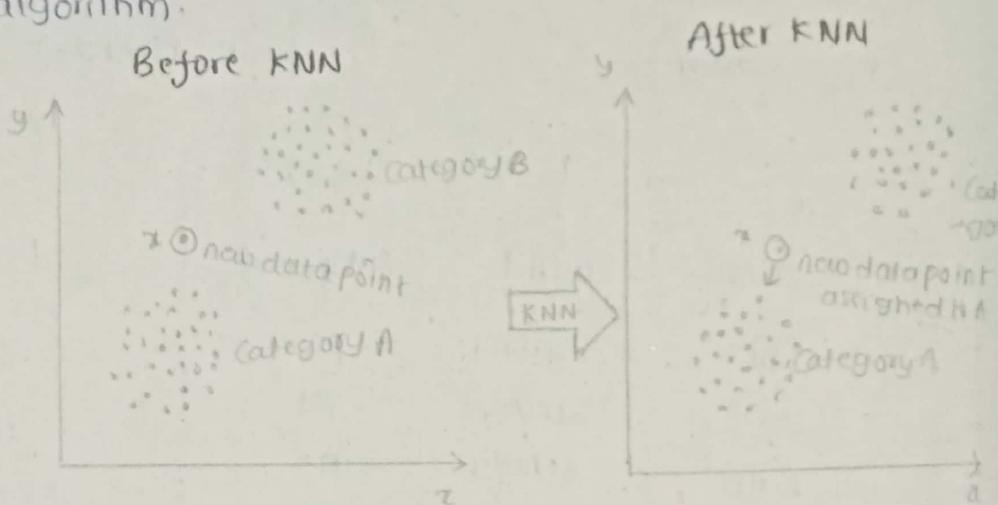
$$\text{Gini}(\text{sunny}) = [1+1, 0-] = 0$$

$$\text{Gini}(\text{Overcast}) = [0+1, 1-] = 0$$

$$\therefore \text{Gini}(\text{outlook}) = 0$$



- ⇒ K-Nearest Neighbours [Page 17]
- \* K-Nearest Neighbour is one of the simplest machine learning algorithm based on supervised learning technique.
  - \* It can be used in many applications.
  - \* It stores all of the training data and classify the new data "x" by finding the training data ( $x_i, y_i$ ) similar to the new data.
  - \* The classifies the data into a category that is much similar to the new data.
  - \* KNN algorithm is also called as non-parametric algorithm. Sometimes it is called as lazy-learner algorithm.



- \* There are two points category A and category B and we have a new datapoint. So, this data point will lies in which of these category.
- \* To solve this type of problem we can use KNN algorithm
- \* With the help of KNN, we can easily identify the class of particular dataset.

SNO	age	salary	purchased
1	25	20	N
2	63	120	Y
3	33	75	N
4	42	100	Y
5	50	100	?

sol: A first calculate the distance between test instance and training instance using euclidian distance,  $k=3$

SNO	age	salary	purchased	Distance
1	25	20	N	122.57
2	63	120	Y	23.85
3	33	75	N	67.186
4	42	100	Y	40.79
5	50	140	? (N) Y	

1.  $(x_1, y_1) = (25, 20)$  and  $(x_2, y_2) = (50, 140)$

By using euclidian distance =  $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

$$= \sqrt{(50-25)^2 + (140-20)^2} = \sqrt{625 + 14400} = \frac{45.4422}{122.5765}$$

2.  $(x_1, y_1) = (63, 120)$  and  $(x_2, y_2) = (50, 140)$

Euclidian distance =  $\sqrt{(50-63)^2 + (140-120)^2}$

$$= \sqrt{169 + 400} = 23.8537$$

3)  $(x_1, y_1) = (33, 75)$  and  $(x_2, y_2) = (50, 140)$

Euclidian distance =  $\sqrt{(50-33)^2 + (140-75)^2} = \sqrt{289 + 4225} = 67.1863$

4)  $(x_1, y_1) = (42, 100)$  and  $(x_2, y_2) = (50, 140)$

Euclidian distance =  $\sqrt{(50-42)^2 + (140-100)^2}$

$$= \sqrt{64 + 1600} = 40.7921$$

⇒ without giving 'k' values:-

Height	Weight ( $w$ )	Class	Distance
167	51	Underweight	6.70
182	62	Normal	13
176	69	Normal	13.41
173	64	Normal	7.61
172	65	Normal	8.24
174	52	Underweight	4.12
169	58	Normal	1.41
173	57	Normal	3
170	55	Normal	2
170	55	?	0

$$1) (x_1, y_1) = (167, 51) \text{ and } (x_2, y_2) = (170, 57)$$

$$\therefore \text{Euclidian distance} = \sqrt{(170-167)^2 + (57-51)^2} = \sqrt{9+36} \\ = \sqrt{45} = 6.708$$

$$2) (x_1, y_1) = (182, 62) \text{ and } (x_2, y_2) = (170, 57)$$

$$\therefore \text{Euclidian distance} = \sqrt{(170-182)^2 + (57-62)^2} = 13$$

$$3) (x_1, y_1) = (176, 69) \text{ and } (x_2, y_2) = (170, 57)$$

$$\therefore \text{Euclidian distance} = \sqrt{(170-176)^2 + (57-69)^2} = 13.41$$

$$4) (x_1, y_1) = (173, 64) \text{ and } (x_2, y_2) = (170, 57)$$

$$\text{Euclidian distance} = \sqrt{(170-173)^2 + (57-64)^2} = 7.61$$

$$5) \text{ Euclidian distance} = \sqrt{(170-172)^2 + (57-65)^2} = 8.24$$

$$6) (x_1, y_1) = (174, 56)$$

$$\therefore \text{Euclidian distance} = \sqrt{(170-174)^2 + (57-56)^2} = 4.123$$

$$7) (x_1, y_1) = (169, 58)$$

$$\therefore \text{Euclidian distance} = \sqrt{(170-169)^2 + (57-58)^2} = 1.41$$

$$8) (x_1, y_1) = (173, 57)$$

$$\text{Euclidian distance} = \sqrt{(170-173)^2 + (57-57)^2} = 3$$

$$9) (x_1, y_1) = (170, 55)$$

$$\text{Euclidian distance} = \sqrt{(170-170)^2 + (57-55)^2} = 2$$

$$10) (x_1, y_1) = (170, 57)$$

$$\text{Euclidian distance} = \sqrt{(170-170)^2 + (57-57)^2} = 0$$

⇒ Naive Bayes Algorithm:

→ Naive Bayes algorithm is a supervised learning algorithm, which is based on Bayes theorem and used for solving classification problems.

→ It is mainly used in text classification that includes a high dimensional training data set.

→ Naive Bayes classifier is one of the simple and most effective classification algorithm.

\* It is a probabilistic classifier, which means it predicts on the basis of the probability of an object.

\* Examples: Spam filters, Sentimental analysis and text classification, article categorization etc.

• Naïve: It assumes that the occurrence of a certain feature is independent of the occurrence of other features.

$$P(A|B) = P(B/A) \cdot P(A)$$

• Bayes: It depends on the principle of Bayes theorem.

• Bayes theorem (or) Bayes rule (or) Bayes law:

$$\text{Equation: } P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Here  $P(A|B)$  = posterior probability

$P(B|A)$  = likelihood probability

$P(A)$  = prior probability

$P(B)$  =

• Advantages:

\* It can be used for binary as well as multiclass classification

\* It is the most popular choice for text classification problems.

\* It is one of the fast and easy machine learning algorithm

\* It performs well in multiclass predictions as compared to the other algorithms.

• Disadvantages:

\* It assumes that all features are independent (or) unrelated ex: Estimate conditional probability of each attribute {color, legs, height, smile} for species classes:

{m, h} Using the data given in table.

Using these probability estimate the probability value for new instance

{color: green, legs: 2, height: tall, smile: No}



SNO	Color	legs	height	smile	Species
1	white	3	short	yes	M
2	Green	2	tall	No	M
3	Green	3	short	yes	M
4	white	3	short	yes	M
5	Green	2	short	No	H
6	white	2	tall	No	H
7	white	2	tall	No	H
8	white	2	short	yes	H

Sol:

First we try to [solution] Calculate conditional probability as well as prior [probab] probability to up with probability of new instances.

new instance: {color: green; leg: 2, height: tall;  
smile: No}

$$P(M) = \frac{4}{8} = \frac{1}{2} = 0.5$$

$$P(H) = \frac{4}{8} = \frac{1}{2} = 0.5$$

We need to calculate conditional probability of each of these attributes

- First attribute: color (white, Green)

color	M	H
white	$\frac{2}{4}$	$\frac{3}{4}$
Green	$\frac{2}{4}$	$\frac{1}{4}$

$$\rightarrow \text{color(white, M)} = \frac{2}{4} \quad \text{color(Green, M)} = \frac{2}{4}$$

$$\text{color(white, H)} = \frac{3}{4} \quad \text{color(Green, H)} = \frac{1}{4}$$

- 2nd attribute: legs (3/2)

legs	M	H
3	$\frac{3}{4}$	$\frac{0}{4}$
2	$\frac{1}{4}$	$\frac{4}{4}$

- 3rd attribute: height (short, tan)

Height	M	H
short	3/4	2/4
tall	1/4	2/4

4th attribute: smile (Yes, No)

smile	M	H
yes	3/4	1/4
no	1/4	3/4

We need to calculate prior probability:  
new instances: {color: green, leg = 2, height: tall, smile: No}

$$\therefore P(M/\text{newinstances}) = P(M) \cdot P(\text{color} = \text{green}) \cdot P(\text{leg} = 2) \\ \cdot P(\text{height} = \text{tall}) \cdot P(\text{smile} = \text{No})$$

$$= 0.5 \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} \\ = 0.0039$$

$$\therefore P(H/\text{newinstance}) = P(H) \cdot P(\text{color} = \text{green}) \cdot P(\text{leg} = 2) \cdot P(\text{height} \\ = \text{tall}) \cdot P(\text{smile} = \text{No}) \\ = 0.5 \cdot \frac{1}{4} \cdot \frac{3}{4} \cdot \frac{2}{4} \cdot \frac{3}{4} \\ = 0.046875$$

$$\therefore P(H/\text{newinstance}) = 0.046875$$

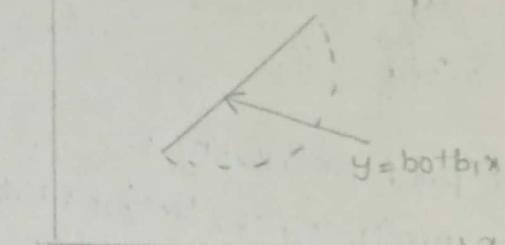
- ⇒ polynomial Regression:
- # polynomial regression is a type of regression analysis  
and used on statistics and machine learning
- # the relationship b/w independent variable 'x' and  
dependent variable 'y' is modelled by the relation
- $y = a_0 + a_1 x_1 + a_2 x_2^2 + a_3 x_3^3 - \dots$
- # if the relationship between x and y are linear then linear  
regression
- # if the relationship between x and y are not linear then  
non linear regression can't be used as it will result in  
large errors.
- # The problem of non linear regression can be solved by two

## Methods:

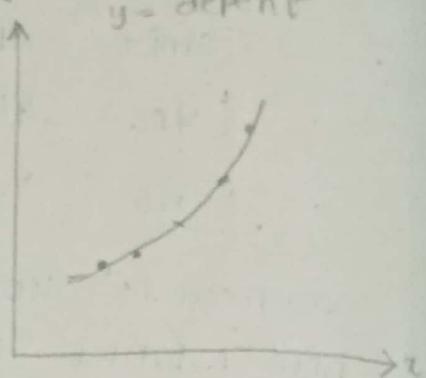
• Transformation of non-linear to linear data so that the linear regression can handle the data.

or Using polynomial regression.

simple linear model



x-independent  
y-dependent



### • Transformation:

- This converts non-linear data to linear data that can be handled using linear regression method.
- Let us consider an exponential function  $y = a e^{bx}$
- + To avoid that large errors.
- \* The transformation can be done by applying function on both sides to get:

$$\ln(y) = \ln(a) + \ln(e^{bx})$$

$$\ln(y) = \ln(a) + bx \rightarrow \ln(e) \quad (\because \log_e = 1)$$

$$\ln(y) = \ln(a) + bx \rightarrow \text{linear function.}$$

### • Polynomial Regression:

- It can handle non-linear relationship among variables by using  $n^{\text{th}}$  degree of polynomial.

- For example, the second-degree polynomial is given as  $y = a_0 + a_1 x + a_2 x^2$  and third-degree polynomial is called cubic transformation given as

$$y = a_0 + a_1 x + a_2 x^2 + a_3 x^3$$

- Generally polynomial of maximum degree  $N$  are used as higher order polynomials take some different (or) strange shapes.

- Consider the polynomial of 2nd degree.
- Then polynomial equation (or) given by

$$y = a_0 + a_1 x + a_2 x^2$$

- The coefficient of  $a_0, a_1$  and  $a_2$  are calculating using the formula.

$$\therefore a = X^{-1}B$$

where

$$X = \begin{bmatrix} n & \sum x_i & \sum x_i^2 \\ \sum x_i & \sum x_i^2 & \sum x_i^3 \\ \sum x_i^2 & \sum x_i^3 & \sum x_i^4 \end{bmatrix} \quad B = \begin{bmatrix} \sum y_i \\ \sum (y_i - \bar{y}) \\ \sum (x_i^2 - \bar{x}^2) \end{bmatrix}$$

Eg : independent(x) dependent(y)

1	1
2	4
3	9
4	15

$x_i$	$y_i$	$x_i y_i$	$x_i^2$	$x_i^2 y_i$	$x_i^3$	$x_i^4$
1	1	1	1	1	1	1
2	4	8	4	16	8	16
3	9	27	9	81	27	81
4	$\frac{15}{29}$	$\frac{60}{96}$	$\frac{16}{30}$	$\frac{240}{338}$	$\frac{64}{100}$	$\frac{256}{354}$
$\bar{x}$	$\frac{29}{4}$	$\frac{96}{4}$	$\frac{30}{4}$	$\frac{338}{4}$	$\frac{100}{4}$	$\frac{354}{4}$

$$\begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} 4 & 10 & 30 \\ 10 & 30 & 100 \\ 30 & 100 & 354 \end{bmatrix}^{-1} \quad X = \begin{bmatrix} 29 \\ 96 \\ 338 \end{bmatrix}$$

$$= \begin{bmatrix} -0.75 \\ 0.95 \\ 0.75 \end{bmatrix}$$

Final polynomial regression

## Linear Regression :-

- 1) linear regression is the most one of the most easiest and popular machine learning.
- 2) it is a statistical method that is used for predictive analysis.
- 3) Linear regression makes prediction for continuous and numeric variables such as sales, temperature, time etc.
- 4) LR algorithm shows a linear relationship between dependent variables and one or more independent variables. Hence it is called a linear regression.
- 5) LR is a mathematical representation:

$$y = a_0 + a_1 x + \epsilon$$

here  $a_0$  = intercepts of the line

$x, y$  = variables / datapoints

$a_1$  = coefficient of line of regression

$\epsilon$  = random error.

- 6) There are two types of LR:

1. simple LR | single linear regression  
more independent / one dependent variable

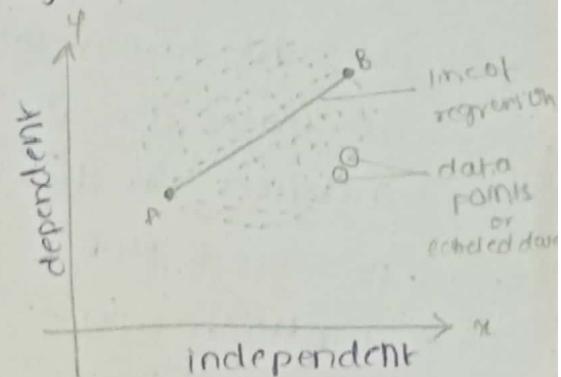
2. multiple linear regression,

1. simple linear regression: if a single independent variable is used to predict the value of a numeric dependent variable, then such a LR algorithm is called SLR (simple linear regression).

2. multiple linear regression: if more than one independent variable is used to predict the value of a numerical dependent variable, then such a LR algorithm is called mLR.

linear regression line:

A linear line showing the relationship between the independent and dependent variables is called a regression line.



there are two types of regression lines:

1. Positive linear relationship

2. Negative linear relationship

1. Positive linear relationship: If the dependent variable increases on the y-axis and independent variable increases on x-axis, then such a relationship is called PLR.

2. Negative linear]

→ Equation :  $y = a_0 + a_1 x$

2. Negative linear relationship: If the dependent variable decreases on the y-axis and independent variable increases on the x-axis, then such a relationship is called NLR.

→ Equation :  $y = - (a_0) + (a_1 x) = - a_0 + a_1 x$

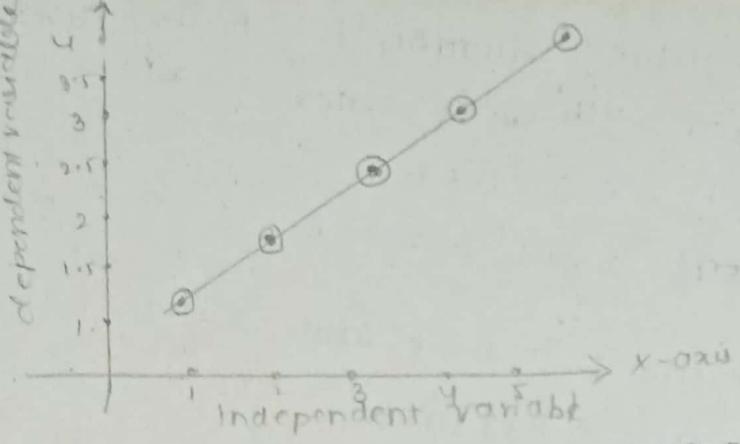
example: let us consider the 5 weeks sales data is given below.

$x_i$ (week)	$y_i$ (sales in thousands)
1	1.2
2	1.8
3	2.6
4	3.2
5	3.8

Apply linear regression technique to predict the seventh & 12th week sales.

Sol: In linear regression algorithm,  
we plot the independent variable and dependent variable

independent variable is represented as x-axis  
dependent variable is represented as y-axis



Linear regression equation;  $y = a_0 + a_1 x + \epsilon$

$$a_0 = \bar{y} - a_1 \bar{x}$$

$$a_1 = \frac{\bar{xy} - (\bar{x})(\bar{y})}{\bar{x}^2 - \bar{x}^2}$$

Here there are 5 items : 1, 2, 3, 4, 5

$x_i$	$y_i$	$x^2$	$xy$
1	1.2	1	1.2
2	1.8	4	3.6
3	2.6	9	7.8
4	3.2	16	12.8
5	4.0	25	20.0
$\bar{x}$	$\bar{y}$	$\bar{x}^2$	$\bar{xy}$
3	2.5	12.5	44.4

$$\bar{x} = \frac{\sum x_i}{n} = \frac{15}{5} = 3$$

$$\bar{y} = \frac{12.5}{5} = 2.5$$

$$[\bar{x}^2 = \frac{55}{5} = 11]$$

$$\bar{xy} = \frac{\sum xy}{n} = \frac{44.4}{5} = 8.88$$

$$a_0 = \bar{y} - a_1 \bar{x}$$

$$= 2.5 - a_1(3)$$

$$a_1 = \frac{\bar{xy} - (\bar{x})(\bar{y})}{\bar{x}^2 - \bar{x}^2} = \frac{8.88 - (3)(2.5)}{11 - 12.5}$$

$$= \frac{1.38}{-44} = -0.031$$

$$\therefore a_0 = 2.5 - (-0.031)(3) = 2.593$$

a linear regression equation is  $y = a_0 + a_1x + \epsilon$   
the predicted 7th week sales  $x=7$

$$y = 2.593 - 0.031(7)$$
$$= 2.593 - 0.217 = 2.376$$

to predict the 12th week sales  $x=12$

$$y = 2.593 - 0.031(12)$$
$$= 2.593 - 0.372 = 2.221$$

∴ 7th week sales is a 2.376 and 12th week sales  
is a 2.221

## 2. Multi Linear Regression:

- \* A single linear regression model involves one independent variable and one dependent variable.
- \* In multiple linear regression model involves one or more independent variables and one dependent variable.
- \* This is the extension of the linear regression

equation:  $y = a_0 + a_1x_1 + a_2x_2 + \epsilon$

Here  $x_1, x_2$  = independent variables

$y$  = dependent variable

$a_0$  = inception of line

$a_1, a_2$  = coefficient of line of regression

$\epsilon$  = random error.

example: Let us consider weekly sales along with sales for products  $x_1, x_2$  are provided.

$x_1$ (product 1 sales)	$x_2$ (product 2 sales)	weekly sales ( $y$ )
1	4	1
2	5	6
3	8	8
4	2	12

How to find the multiple regression for the values.

$x_1, x_2$  are two independent variables,  $y$  is a dependent variable. The value of  $y$  depends on  $x_1$  and  $x_2$  in the particular case.

We use matrix approach for finding multiple regression. The matrix for  $x$  and  $y$  are given as follows:

$$x = \begin{bmatrix} 1 & 1 & 4 \\ 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 2 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ 6 \\ 8 \\ 12 \end{bmatrix}$$

The coefficient of the multiple regression equation is given

$$\therefore a = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix}$$

The regression coefficient for multiple regression is calculated

$$\therefore \hat{a} = ((x^T x)^{-1} x^T) y.$$

$$\text{Here } x^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{bmatrix}$$

$$x^T \cdot x = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{bmatrix} \begin{bmatrix} 1 & 1 & 4 \\ 1 & 2 & 5 \\ 1 & 3 & 8 \\ 1 & 4 & 2 \end{bmatrix}$$

$$= \begin{bmatrix} 1(1)+1(1)+1(1)+1(1) & 1(1)+1(2)+1(3)+1(4) & 1(4)+1(5)+1(8)+1(2) \\ 1(1)+2(1)+3(1)+4(1) & 1(1)+2(2)+3(3)+4(4) & 1(4)+2(5)+3(8)+4(2) \\ 4(1)+5(1)+8(1)+2(1) & 4(1)+5(2)+8(3)+2(4) & 4(4)+5(5)+8(8)+2(1) \end{bmatrix}$$

$$x^T \cdot x = \begin{bmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{bmatrix}$$

$$\therefore (x^T \cdot x)^{-1} = \frac{\text{Adjoint}}{\det A} = \begin{bmatrix} 0 & -1 & 1 \\ -1 & 0 & -1 \\ 1 & -1 & 0 \end{bmatrix}$$

$$\rightarrow \det A = 4(30 \times 109 - 46 \times 46) - 10(10 \times 109 - 46 \times 19) + 19(10 \times 46 - 30 \times 19)$$

$$= 4616 - 2160 - 2090$$



$$\therefore \det A = 366$$

$\therefore \text{adj } A :$

$$\det A = \begin{bmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{bmatrix}$$

$$A_1 = +4 \begin{bmatrix} 30 & 46 \\ 46 & 109 \end{bmatrix} = 4 [30 \times 109 - 46 \times 46] \\ = 4(1154) = [4616]$$

$$B_1 = -10 \begin{bmatrix} 10 & 46 \\ 19 & 109 \end{bmatrix} = -10 [10 \times 109 - 46 \times 19] = -216 \\ = -216$$

$$A_3 c_1 = 19 \begin{bmatrix} 10 & 30 \\ 19 & 46 \end{bmatrix} = 19 [10 \times 46 - 30 \times 19] = +110 \\ = -110$$

$$A_2 = -10 \begin{bmatrix} 10 & 19 \\ 46 & 109 \end{bmatrix} = -10 [10 \times 109 - 19 \times 46] = +216$$

$$B_2 = 30 \begin{bmatrix} 4 & 19 \\ 19 & 109 \end{bmatrix} = 30 [4 \times 109 - 19 \times 19] = 75$$

$$C_2 = -46 \begin{bmatrix} 4 & 10 \\ 19 & 46 \end{bmatrix} = -46 [4 \times 46 - 10 \times 19] = -(-6) = 6$$

$$A_3 = + \begin{bmatrix} 10 & 19 \\ 30 & 46 \end{bmatrix} = + [10 \times 46 - 19 \times 30] = +110$$

$$B_3 = - \begin{bmatrix} 4 & 19 \\ 10 & 46 \end{bmatrix} = - [4 \times 46 - 19 \times 10] = -(-6) = 6$$

$$C_3 = + \begin{bmatrix} 4 & 10 \\ 10 & 30 \end{bmatrix} = + [4 \times 30 - 10 \times 10] = 20$$

$$\text{adj } A = \begin{bmatrix} 1154 & +216 & +110 \\ -216 & 75 & +6 \\ -110 & +6 & 20 \end{bmatrix} \times \frac{1}{366}$$

$$\frac{\text{adj } A}{\det A} = \begin{bmatrix} 1154 & -216 & -110 \\ -216 & 75 & 6 \\ -110 & 6 & 20 \end{bmatrix} \begin{bmatrix} y \\ 366 \end{bmatrix}$$

$$\frac{\text{adj } A}{\det A} = \begin{bmatrix} 3.1530 & -0.590 & -0.3005 \\ -0.590 & 0.204 & 0.0163 \\ -0.300 & 0.01 & 0.05 \end{bmatrix}$$

$$\therefore \frac{\text{adj } A}{\det A} = \begin{bmatrix} 3.1530 & -0.590 & -0.3005 \\ -0.590 & 0.204 & 0.0163 \\ -0.300 & 0.01 & 0.05 \end{bmatrix}$$

$$(x^T \cdot x)^{-1} = \begin{bmatrix} 3.15 & -0.59 & -0.30 \\ -0.59 & 0.20 & 0.01 \\ -0.30 & 0.01 & 0.05 \end{bmatrix}$$

$$\therefore [x^T \cdot x]^{-1} x^T = \begin{bmatrix} 3.15 & -0.59 & -0.30 \\ -0.59 & 0.20 & 0.01 \\ -0.30 & 0.01 & 0.05 \end{bmatrix} \begin{bmatrix} 1 & 2 & 1 & 1 \\ 1 & 2 & 3 & 4 \\ 4 & 5 & 8 & 2 \end{bmatrix}$$

$$= [3.15(1) + (-0.59)(2) + (-0.30)(4)]$$

$$= \begin{bmatrix} 1.36 & 0.47 & -1.02 & 0.19 \\ -0.35 & -0.14 & 0.09 & 0.23 \\ -0.09 & -0.03 & 0.13 & -0.16 \end{bmatrix}$$

$$\therefore [(x^T \cdot x)^{-1} x^T] y = \begin{bmatrix} 1.36 & 0.47 & -1.02 & 0.19 \\ -0.35 & -0.14 & 0.09 & 0.23 \\ -0.09 & -0.03 & 0.13 & -0.16 \end{bmatrix} \begin{bmatrix} 1 \\ 6 \\ 8 \\ 12 \end{bmatrix}$$

$$= \begin{bmatrix} 1.36 + 2.82 - 8.16 + 2.28 \\ -0.35 + (-0.84) + 0.72 + 2.76 \\ -0.09 + (-0.18) + 1.04 - 1.92 \end{bmatrix} = \begin{bmatrix} -1.7 \\ 2.29 \\ -1.15 \end{bmatrix}$$

$$\therefore y = a_0 + a_1 x_1 + a_2 x_2$$

$$= (-1.7) + 2.29(1) + (-1.15)(4)$$

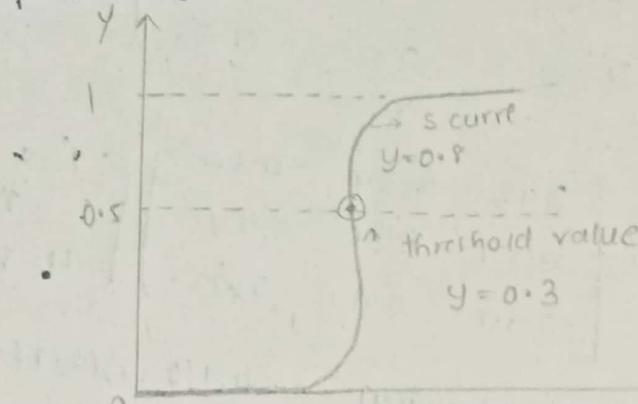
$$= -1.7 + 2.29 - 4.6 = -4.01$$

$$\boxed{y = -4.01}$$



## ⇒ Logistic Regression:

- Logistic regression is a supervised machine learning algorithm used for classification tasks (true or false, Yes or No, + or -)
- where the goal is to predict the probability than an instance belongs to given class or not.
- logistic regression is a statistical algorithm which analyzes the relationship between two data factors.
- logistic regression is used for binary classification where we use sigmoid function (or) logistic function.
- It takes input as independent variables and produces a probability value between '0' and '1'.



- We have two classes class "0" and class "1".
- If the value of the logistic function for an input is greater than 0.5.
- Then it belongs to class "1" Otherwise it belongs to class "0".
- The sigmoid function is a mathematical function used to map the predicted values to probabilities
- It maps any real value into another value within a range of "0" and "1".
- In logistic regression, instead of fitting a regression line, we fit an "S" shaped logistic function, which predicts two maximum value (0 or 1).

Sigmoid function formula:

$$S(x) = \frac{1}{1 + e^{-x}}$$

Simple linear and multiple linear equations

$$Z = a_0 + a_1 x_1 + a_2 x_2 + \dots$$

logistic regression equation:

$$h(x) = \frac{1}{1+e^{-x}}$$

main aim: logistic regression is used when the dependent variable is categorical value (or) discrete value.  
ex: to predict whether an email is spam "1".

#### Types of logistic regression:

There are three(3) types of logistic regression:

1. Binomial

2. multinomial

3. Ordinal

1. Binomial: There can be only two(2) possible types of the dependent variables, such as "0" or "1", pass or fail.

2. multinomial:

There can be three(3) or more possible unordered types of the dependent variable, such as "cat", "dog" or "sheep".

3. Ordinal:

There can be three(3) or more possible ordered type of dependent variables, such as "low", "medium", or "high".

Q Difference between linear regression and logistic regression.

feature	linear regression	logistic regression
purpose	used for predicting continuous values	Used for classification (categorical output)
output type	produces continuous numerical values	produces probabilities mapped to classes (0/1)
mathematical function	uses a linear equation $y = mx + c$	uses a sigmoid function $g(x) = \frac{1}{1+e^{-x}}$

Type of problem	regression problems (eg: predicting sales, temperatures)	classification, On problems (eg : Spam detection, fraud detection.)
Algorithm type	Based On least squares method	Based On maximum likelihood estimation.
Linearity	assumes a linear relationship between input and output	Does not assume a linear relationship; applies a non-linear transformation (sigmoid function)
Example Use cases	predicting house prices, stock prices etc	Email spam detection, disease diagnosis (yes/no)

