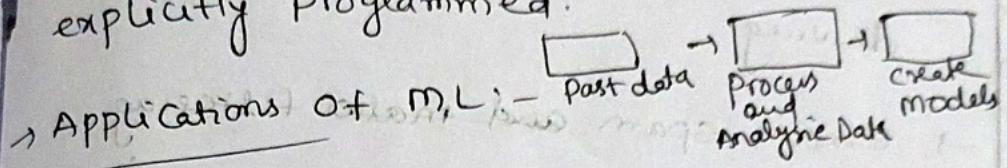


## Machine Learning

### 1) machine learning:-

→ machine learning is an application of AI that enables systems to learn and improve from experience without being explicitly programmed.



### 1) speech Recognition:-

→ while using Google, we get an option of search by voice, it comes under speech recognition.

→ Google Assistant, Siri, Alexa are using speech recognition technology to follow the voice instructions.

### 2) 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.

→ whenever we search for some product on Amazon, then we start getting an advertisement for the same product while internet surfing on the browser.

## → self driving cars :-

- machine learning plays a significant role in self driving cars (Tesla),
- It is using unsupervised learning method to train the car models to detect people and objects while driving.

## → Email spam and email filtering

- whenever we receive a new email, it is filtered automatically as important normal and spam.
- we always receive an important mail in our inbox with the important symbol and spam emails in our spam box.

## → Medical Diagnosis :-

- In medical science, machine learning is used to diseases diagnoses.
- with this medical technology can predict the exact position of disease in the brain.

- Automatic language Translation:-
  - translates the text from one language to another language.
- Types of machine learning:
  - 1) supervised learning
  - 2) unsupervised learning
  - 3) Reinforcement Learning
- 1) supervised machine learning:
  - The supervised learning technique, we train the machines using the labelled data set, and based on the training, the machine predicts the output.
  - we train the machine with input and corresponding output, and then we ask the machine to predict the output using the test dataset.
  - suppose we have an input dataset of cats and dog images. so first, we will provide the training to the machine to understand the images, such as the shape and size of the tail of cat and dog, shape of eyes, height (dogs are taller than cats are smaller) etc.

- After completion of training, we input the picture of cat and ask the machine to identify and predict the output.
- The main goal of the supervised learning techniques is to map the input variable ( $x$ ) with the output variable ( $y$ ).  
Ex:- Application: Fraud detection, Spam filtering
- Categories of supervised machine learning
- Supervised machine learning can be classified into two types of problems.
  - 1) Classification
  - 2) Regression
- 1) Classification algorithms are used to solve the classification problem in which the output variable is categorical. such as "Yes" or "No", male or female etc.
- 2) Regression:  
~~Regression finds correlation between dependent (target) and independent (predicted) variables.~~  
→ Regression algorithm

Regression algorithms are used to solve regression problems in which there is a linear (straight line) relationship between input and output variables.

→ These are used to predict continuous output variables, such as market trends, weather prediction, etc.

→ unsupervised machine learning:-

→ In unsupervised machine learning the machine is trained using the unlabeled dataset.

→ The main aim of the unsupervised learning algorithm is to group or categories the unsorted dataset according to the similarities, patterns and differences.

→ Machines are instructed to find the hidden patterns from the input dataset.

→ Suppose there is a basket of fruit images and we input it into the machine learning model.

→ The images are totally unknown to the model and the task of the machine is to find the patterns and categories of the objects.

→ So, now the machine will discover its patterns and such as colour, shape and predict the output.

When it is tested with test dataset.

e.g. Network analysis:- identity Plagiarism  
Recommendation system:- e-commerce websites

→ Categories of unsupervised Machine Learning

1) Clustering

2) Association.

→ Clustering:-

The clustering technique is used when we want to find the inherent groups from the data.

→ It is a way to group the objects into a cluster such that the objects with the most similarities remain in one group, and have fewer or no similarities with the objects of other groups.

→ An example of the clustering algorithm is grouping the customers by their purchasing behaviour.

### → Association:-

- Association rule learning is a technique that checks for the dependency of one data item on another data item and maps accordingly.
- It tries to find some interesting relations or associations among the variables of dataset.
- It is based on different rules to discover the interesting relations between variables.
- For example if a customer buys bread, he most likely can also buy butter, eggs or milk. So these products are stored within a shelf or mostly nearby.

### → Reinforcement Learning:-

- It works on a feedback-based process, in which an AI agent (A software component) automatically explores its surroundings by hitting & trail, taking action, learning from experiences and improving its performance.
- Agent gets rewarded for each good action and gets punished for each bad action.

Hence the goal of reinforcement learning agent is to maximize the rewards.

- In reinforcement learning there is no labelled data like supervised learning and agents learn from their experience only.
- The reinforcement learning process is similar to a human being, for example, a child learns various things by experiences in his day to day life.
- Applications of RL are Video Games, Robotics (intelligent Robots)
- Basic concepts in machine learning

→ Features:- attributes that describes an instances in the data set  
Ex:-  
Customer number of purchases  
age  
product brought

→ Feature Selection:- (optimisation)  
selecting best features  
process of selecting the optimal features to include in the training phase

→ DATA SET :- data used to train the system and detect patterns to design algorithms

→ Learning (or) Training :- Patterns of a dataset are detected

→ Tunning :- Optimising the parameters of an algorithm to find the best combination for your specific data set. (Hyper parameter optimisation)

→ models :- models to make predictions or draw conclusions from work

→ validation :- It is the process of evaluating a trained model on test data set.

→ machine learning process :-

1) Gathering data

2) Preparing the data

3) choosing a model

4) Training

5) Evaluation

6) Hyperparameter tuning

7) Prediction.

## → Testing Machine Learning Algorithms

→ Goal of ML testing:-

1) Quality assurance:-

It is required to make sure that the software system works according to the requirement.

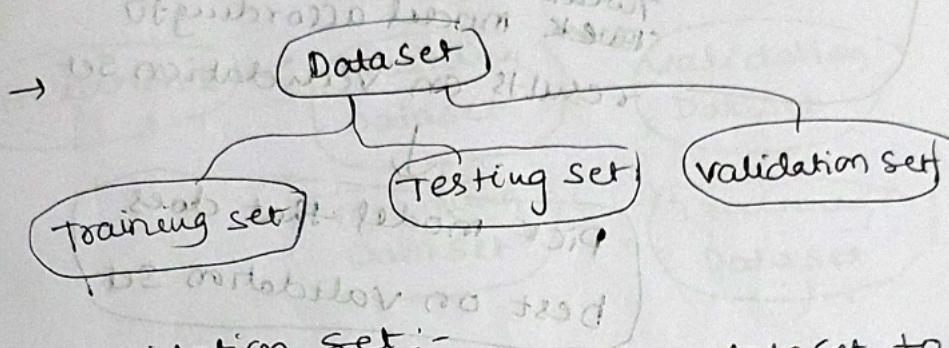
2) The purpose of ML testing is first of all to ensure that the learned logic will remain consistent, no matter how many times we call the program.

→ Model evaluation in ML testing:-

1) Unit tests:- The program is broken down in blocks, and each unit is tested separately.

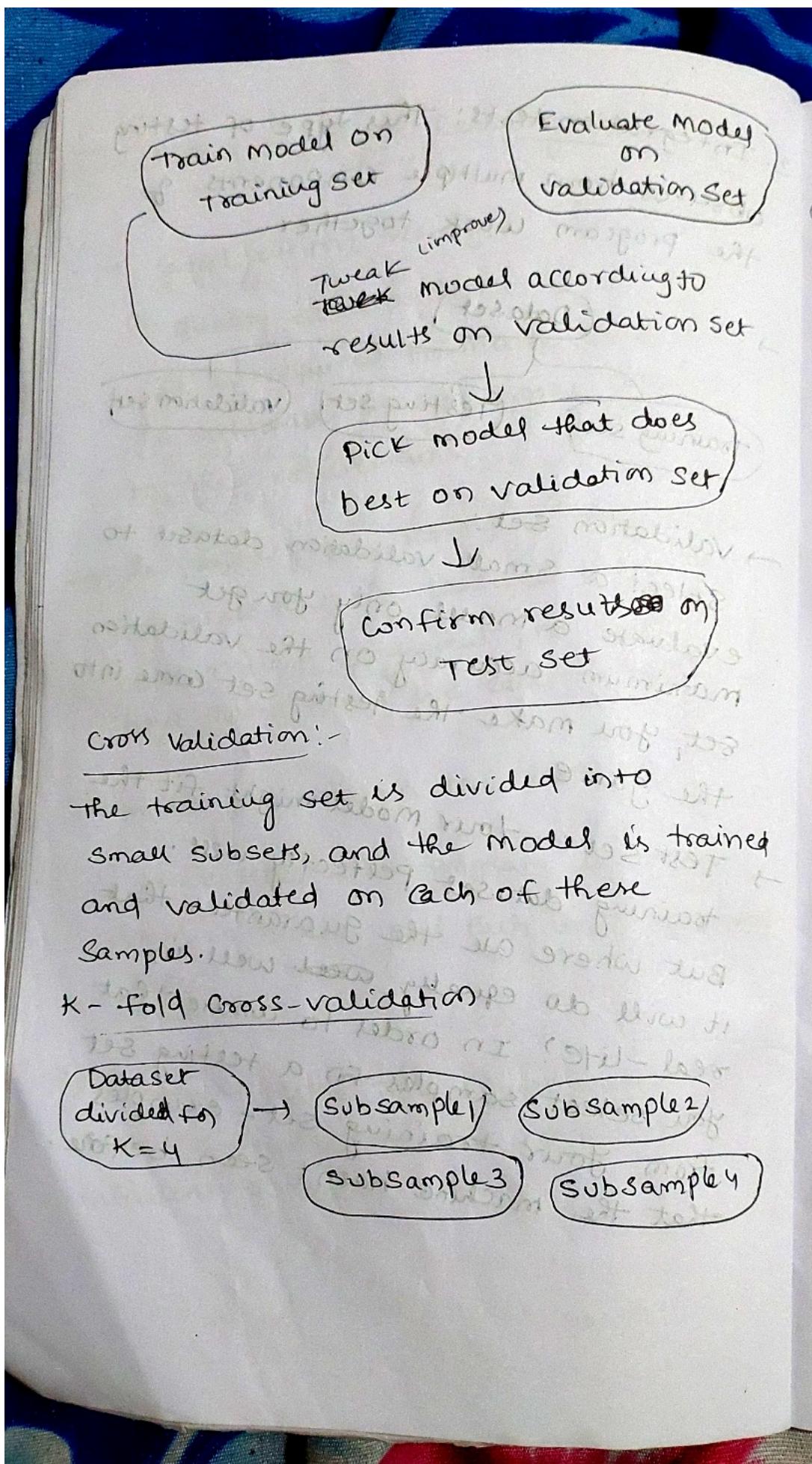
2) Regression tests: They cover already tested software to see if it doesn't suddenly break.

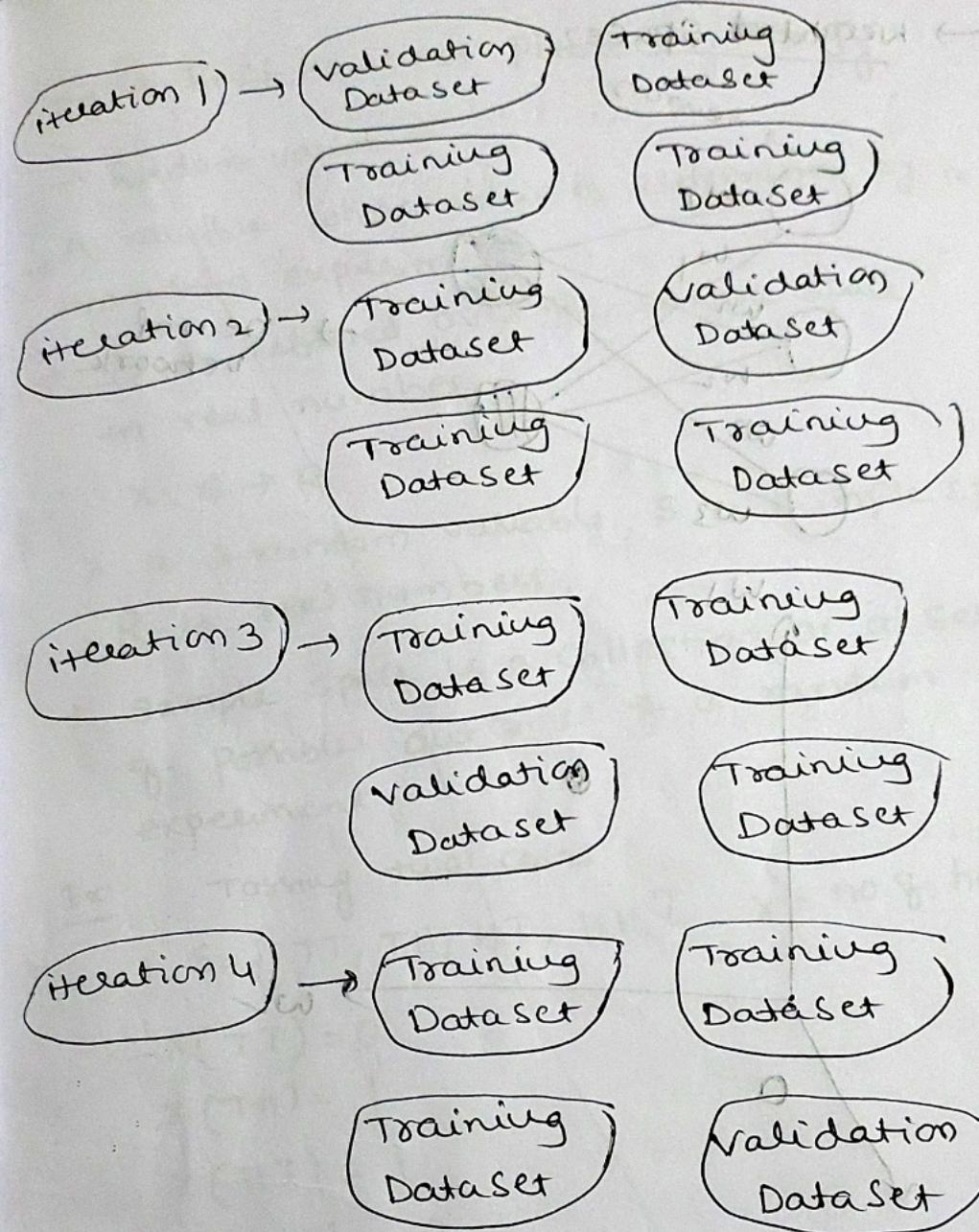
3. Integration tests:- This type of testing observes how multiple components of the program work together.



→ Validation set:- Select a small validation dataset to evaluate a model. Only you get maximum accuracy on the validation set, you make the testing set come into the game.

→ Test set:- Your model might fit the training dataset perfectly well. But where are the guarantees that it will do equally ~~well~~ well in real-life? In order to assure that you select samples for a testing set from your training set - examples that the machine hasn't seen before.

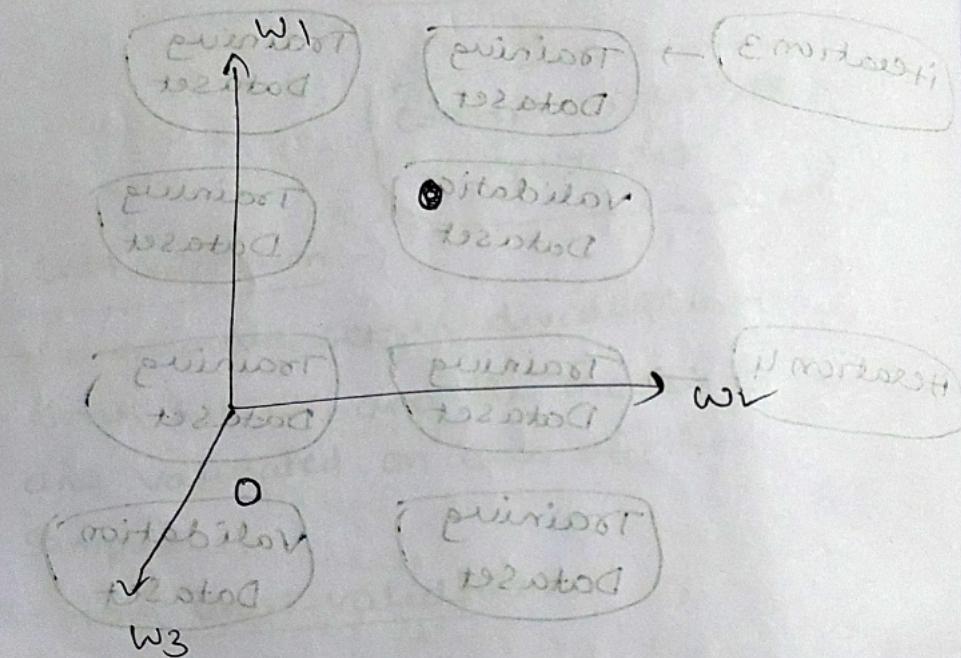
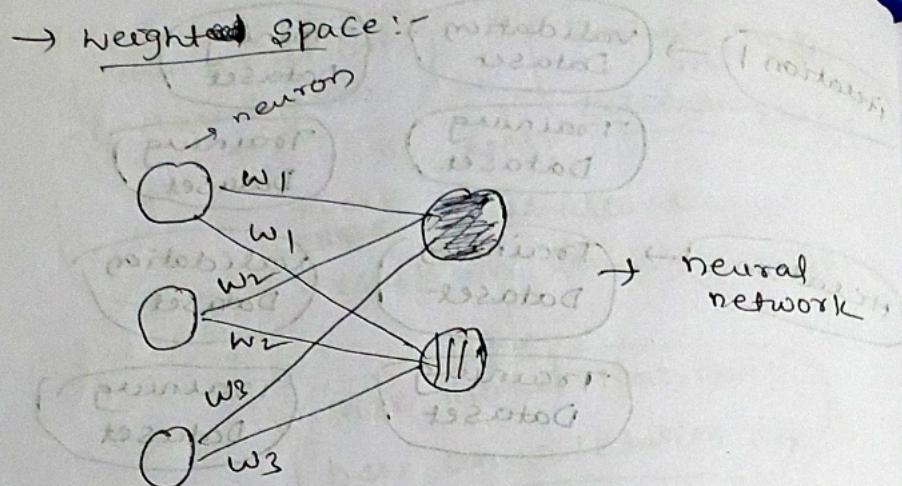




for each test step we will start on 71

of the new elements will go into one of

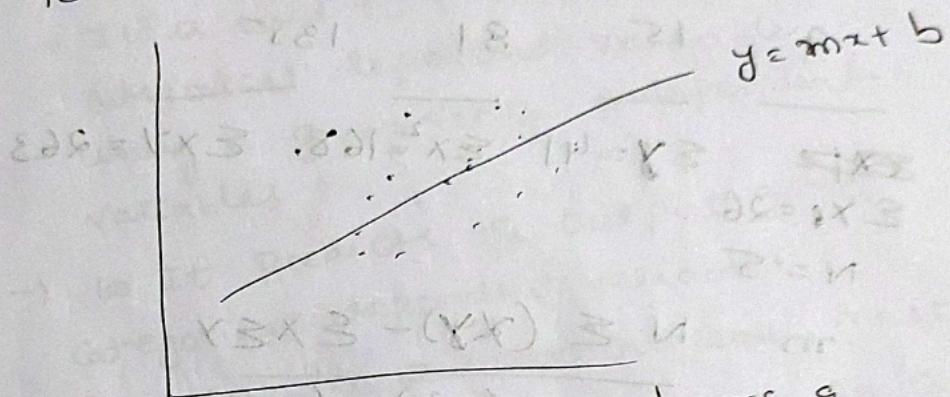
the two groups so we will do this



If we treat the weights that get fed into one of the neurons as a set of coordinates it is known as weight space.

→ Lineal regression:-

It is a statistical method used to find the line of best fit of the form of an equation such as  $y = mx + b$  to the given data.



Ex- Sam found how many hours of Sunshine vs how many ice creams were sold at the shop from Monday to Friday:

X hours of Sunshine	Y Ice creams sold
2	4
3	5
5	7
7	10
9	15

using least square method, find the regression line of the above data.

$$2x3 - 12, 3x3 - 12, 5x5 - 12, 7x7 - 12, 9x9 - 12$$

$$\begin{array}{cccc} x & y & x^2 & xy \\ \hline 2 & 4 & 4 & 8 \\ 3 & 5 & 9 & 15 \\ 5 & 7 & 25 & 35 \\ 7 & 10 & 49 & 70 \\ 9 & 15 & 81 & 135 \end{array}$$

$$\sum x = 26 \quad \sum y = 41 \quad \sum x^2 = 168 \quad \sum xy = 263$$

$$N = 5$$

$$m = \frac{N \sum (xy) - \sum x \sum y}{N \sum (x^2) - (\sum x)^2}$$

$$= \frac{5 \times 263 - 26 \times 41}{5 \times 168 - (26)^2}$$

$$= \frac{1315 - 1066}{840 - 676} = 1.5183$$

$$b = \frac{\sum y - m \sum x}{N}$$

$$= \frac{41 - 1.5183 \times 26}{5}$$

$$= 0.3049$$

$$y = mx + b$$

$$y = 1.518x + 0.305$$

suppose weather forecast says  
we expect 8 hours of sun tomorrow

$$1.518 \times 8 + 0.305 = 12.45$$

if interval of work labor predicted

then what formulae will predict more

→ Logistic regression:

→ It is used to predict the

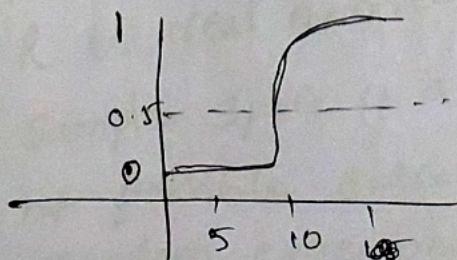
categorical dependent variable  
using a given set of independent  
variables

→ It Predicts the output of a  
categorical dependent variable.

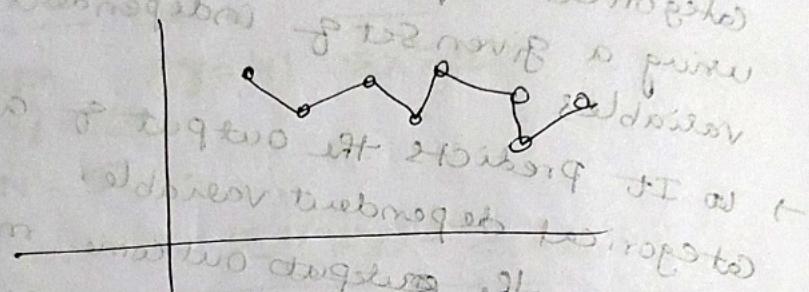
therefore the output outcome must  
be a categorical value. It can  
be either YES or NO, 0 & 1, true or false,  
etc. but instead of giving the exact  
value as 0 and 1. It gives the  
probabilistic values which lies  
between 0 and 1.

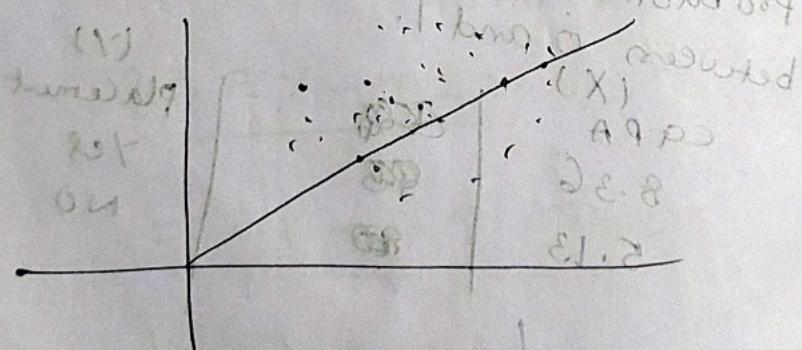
(X)	(Y)
CAPA	YES
8.36	0.9
5.13	0.2

$$y = \frac{1}{1+e^{-x}}$$



## → Overfitting and Underfitting

- Overfitting occurs when our machine learning model tries to cover all the more than the required data points present in the given dataset.
- ML model works well on training data but not on test data.
- 
- Underfitting occurs when our ML model is not able to capture the underlying trend of the data. (unable to capture the data points present in the plot)
- ML model not able to capture the relationship in training data.



→ curse of dimensionality:

2	4	6	10	12	100	1000
$m_1$	$m_2$	$m_3$	$m_4$	$m_5$	$P$	$P$
					$M_6$	$M_7$

↑ threshold

• Red ball → clarity

orange sphere

eatable

play

Red



white ball



X  
as it is eatable  
it is not ball



(not eatable  
so it ball)

→ As the number of features increases  
your model get confused.

→ A Brief Review of Probability Theory

1) Random variable:-

→ A variable whose value is determined by a random experiment.

→ It is defined over the sample space

to real numbers.

$X: S \rightarrow R$

$X$  is a random variable,  $S$  is Sample Space

$R$  is real numbers.

→ Sample space is a collection or a set of possible outcomes of a random experiment.

Ex:- Tossing two coins

$$S = \{ TT, TH, HT, HH \}$$

$X = \text{no of heads}$

$$X(TT) = 0$$

$$X(TH) = 1$$

$$X(HT) = 1$$

$$X(HH) = 2$$

All the elements of sample space are assigned to real numbers.

→ Discrete Random variables

It takes only finite number of distinct values : Ex: 0, 1, 2, 3, 4

→ Continuous random variables

Infinite and uncountable set of values

Ex:- Interest rates of loans in a country.

→ conditional probability :-

If A and B are two events in a

sample space 'S', then Conditional

probability of A given B is

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \text{ when } P(B) \neq 0$$

↓  
event B is already occurred

Bayes rule & theorem:

It determines the conditional probability of an event A given that event B has

already occurred, or occurred

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Posterior

evidence

likelihood

prior

marginal

$P(A)$  = how likely A happens (Prior knowledge)  
(the probability of a hypothesis is true  
before any evidence is present).

$P(B)$  = how likely B happens

The probability of observing the  
evidence (marginal)

$P(A|B)$  = how likely A happens given that  
B has happened (Posterior)

$P(B|A)$  = how likely B happens given that  
A has happened (likelihood)

→ Independence and Conditional Independence

→ Two random variables  $x$  and  $y$  are  
said to be statistically independent  
if and only if

$$P(x,y) = P(x)P(y)$$

Ex:- Independent - X: Throw of a dice

Y: Toss of a coin

NOT Independent - X: Height Y: weight

In general as height increases weight  
increases

→ Conditional Independence :-

Two random variables X and Y are said to be independent given Z if and only if

$$P(X|Y|Z) = P(X|Z)P(Y|Z)$$

Ex - X: Height Y: vocabulary Z: Age

→ Not independent unless I gave some condition. (If I tell you this person is just 2 feet tall, it automatically means it is most probably that this person must be child and have low vocabulary.

mean and variance :-

→ Mean is the average of given set of numbers

→ Variance gives the variation from the expected value.

→ Variance also measures amount of fluctuation of the variable.

$$\text{variance} = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

$x_i$  = value of the one observation

$\bar{x}$  = mean value of all observations

$n$  = number of observations.

→ Confusion Matrix:- It is a table used to define the performance of a classification algorithm.

		1	0	
		True Positive (TP)	False Negative (FN)	Actual
Actual	1	True Positive (TP)	False Negative (FN)	Prediction
	0	False Positive (FP)	True Negative (TN)	

True positive:- The model has predicted yes, and actual value was also yes.

True Negative: model has predicted no and actual value was also no.

False positive:- The model has predicted yes but the actual value was no.

False negative: The model has predicted no but actual value was yes.

→ Accuracy:- The accuracy is used to find the portion of correct classified values

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad \begin{matrix} (\text{no of correct prediction}) \\ \text{Total no of predictions} \end{matrix}$$

Precision = what proportion of Predicted positives is truly positive (Suppose if I have to classify shoes added Nike)

$\text{Precision} = \frac{TP}{TP + FP}$  Out of all dog predicting how many you got it right.

Recall:- what proportion of actual positives is correctly classified.

$\text{Recall} = \frac{TP}{TP + FN}$  R. Recall is out of all dog truth how many you got it right.

$F_1 \text{ Score} = \frac{2PR}{P+R}$  (Harmonic mean of precision and Recall)  
↓ average of precision and Recall.

→ Consider the given dataset Apply Naive-Bayes algorithm and predict that if a fruit has the following properties, then which type of fruit it is

Fruit	Yellow	Sweet	Long	Total
Orange	350	450	0	650
Banana	400	300	350	400
Others	50	100	50	150
TOTAL	800	850	400	1200

$$\frac{N_1 + N_2 + N_3 + N_4}{N_1 + N_2 + N_3 + N_4 + N_5} = \text{accuracy}$$

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

$$P(\text{yellow}|\text{orange}) = \frac{P(\text{orange}|\text{yellow}) \cdot P(\text{yellow})}{P(\text{orange})}$$

probability of yellow  
given that fruit is orange

$$= \frac{\frac{350}{800} \times \frac{800}{1200}}{\frac{650}{1200}}$$

$$= 0.53$$

$$P(\text{sweet}|\text{orange}) = \frac{P(\text{orange}|\text{sweet}) \cdot P(\text{sweet})}{P(\text{orange})}$$

$$= \frac{\frac{450}{850} \times \frac{850}{1200}}{\frac{650}{1200}} = 0.69$$

$$P(\text{long}|\text{orange}) = \frac{P(\text{orange}|\text{long}) P(\text{long})}{P(\text{orange})}$$

$$= \frac{\frac{0}{400} \times \frac{400}{1200}}{\frac{650}{1200}} = 0$$

$$P(\text{fruit}|\text{orange}) = 0.53 \times 0.69 \times 0 = 0$$

Note:

Probability of event to happen

$$P(E) = \frac{\text{No of favourable outcomes}}{\text{Total No of outcomes}}$$

$$P(\text{yellow} | \text{Banana}) = \frac{P(\text{orange} | \text{Banana}) P(\text{yellow})}{P(\text{Banana})}$$

$$= \frac{\frac{400}{800} \times \frac{800}{1200}}{\frac{400}{1200}} = 1$$

$$P(\text{sweet} | \text{Banana}) = \frac{P(\text{Banana} | \text{sweet}) P(\text{sweet})}{P(\text{Banana})}$$

$$= \frac{\frac{300}{850} \times \frac{850}{1200}}{\frac{400}{1200}} = 0.75$$

$$P(\text{long} | \text{Banana}) = P(\text{Banana} | \text{long}) P(\text{long})$$

$$= \frac{\frac{350}{450} \times \frac{400}{1200}}{\frac{400}{1200}} = 0.875$$

$$P(\text{others} | \text{yellow}) P(\text{Fruit} | \text{Banana}) = 1 \times 0.75 \times 0.875 = 0.65$$

$$P(\text{yellow} | \text{others}) = P(\text{others} | \text{yellow}) P(\text{yellow})$$

$$P(\text{others})$$

$$= \frac{\frac{50}{800} \times \frac{800}{1200}}{\frac{150}{1200}} = 0.33$$

$$P(\text{sweet} | \text{others}) = \frac{P(\text{others} | \text{sweet}) P(\text{sweet})}{P(\text{others})}$$

$$= \frac{\frac{100}{850} \times \frac{850}{1200}}{\frac{150}{1200}} = 0.66$$

$$P(\text{long} | \text{others}) = \frac{P(\text{others} | \text{long}) P(\text{long})}{P(\text{others})}$$

$$= \frac{\frac{50}{400} \times \frac{400}{1200}}{\frac{150}{1200}} = 0.33$$

$$P(\text{Fruit} | \text{others}) = 0.33 \times 0.66 \times 0.33 = 0.072$$

→ Fruit = Banana.

→ Bias:- while making predictions, a difference occurs between Expected values made by the model and actual values and this difference is known as bias.

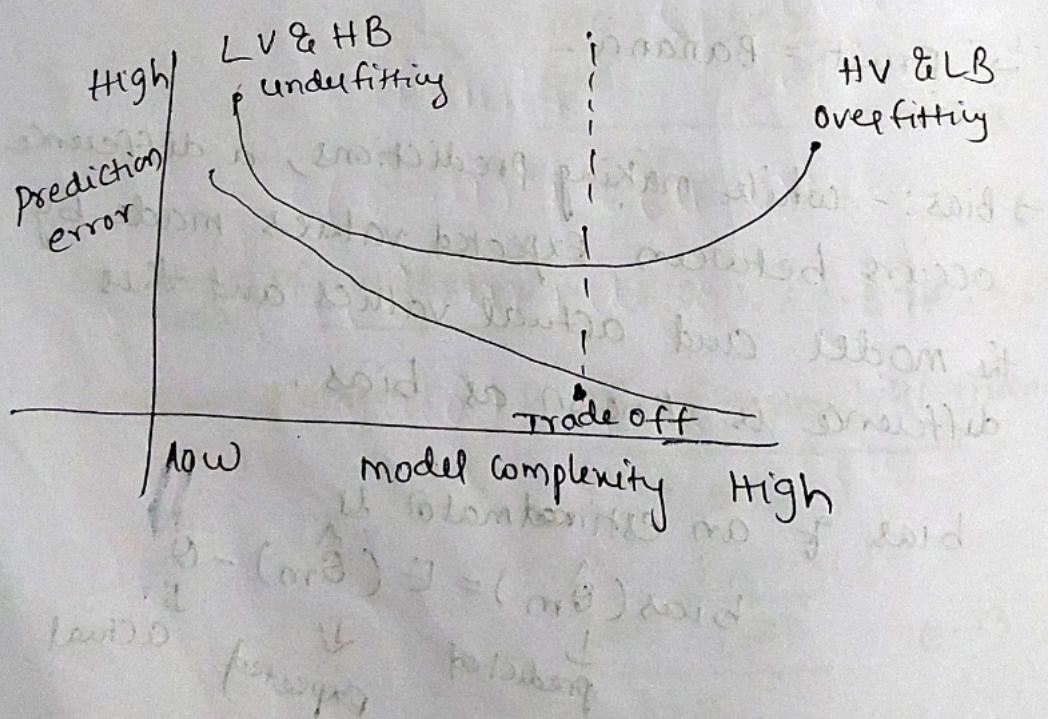
bias of an estimator is

$$\text{bias}(\hat{\theta}_m) = E(\hat{\theta}_m) - \theta$$

↓ predicted      ↓ Expected      ↓ actual

## → the Bias-variance Tradeoff:-

- Bias and variance measure two different source of errors
- while building the machine learning model it is really important to take care of bias and variance in order to avoid overfitting and underfitting in the model
- Bias-variance trade-off is about finding the ~~center~~ sweet spot to make a balance between bias and variance errors.
- the goal of any supervised machine learning algorithm is to achieve low bias and low variance



→ Linear models for Regression:-  
(Linear Basis Function models)

→ the basic linear model for regression is  
a model that involves a linear combination  
of the input variables.

$$y(w, x) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_D x_D$$

$$\text{where } \mathbf{x} = (x_1, \dots, x_D)^T$$

- the key attribute of this function is, that  
it is a linear function of the parameters  
 $w_0, w_1, \dots, w_D$ .
- It is also a linear function of the  
input variable  $x$ .

→ the Bias-variance Decomposition:-

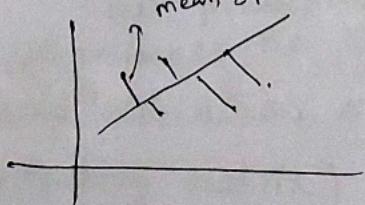
The bias-variance decomposition is a useful  
theoretical tool for understanding a learning  
algorithm's performance characteristics.

$$\text{Mean squared error} \text{ MSE}(\hat{\theta}) = \text{bias}(\hat{\theta})^2 + \text{var}(\hat{\theta})$$

$$= E[(\hat{\theta} - \theta)^2] = \cancel{\text{bias}^2(\hat{\theta}, \theta)} + \cancel{\text{var}(\hat{\theta})}$$

$\hat{\theta}$  is random variable

$\hat{\theta}$  is predicted,  $\theta$  is actual



mean squared error

mean squared error  
measures how close  
a regression line is  
to a set of data  
points.

$$E[(\hat{\theta}_s - \theta)^2] = E[\hat{\theta}_s^2] + \theta^2 - 2E[\hat{\theta}_s]$$

~~( $E(\theta)$  is constant which is ' $\theta$ ' )~~

$$\text{Bias}^2(\hat{\theta}_s, \theta) = (E[\hat{\theta}_s] - \theta)^2$$

$$= E^2[\hat{\theta}_s] + \theta^2 - 2E[\hat{\theta}_s]$$

$$\text{var}(\hat{\theta}_s) = E[\hat{\theta}_s^2] - E^2[\hat{\theta}_s]$$

$$\text{mse}(\hat{\theta}) = \text{bias}(\hat{\theta})^2 + \text{var}(\hat{\theta})$$

$$\Rightarrow E(\hat{\theta} - \theta)^2 = E((\hat{\theta} - \mu) + (\mu - \theta))^2$$

$$= E((\hat{\theta} - \mu)^2 + 2(\hat{\theta} - \mu)(\mu - \theta) + (\mu - \theta)^2)$$

$\therefore \text{let } \mu = \cancel{\theta}$   
 $E(\hat{\theta} - \mu) = 0$

$$= E(\hat{\theta} - \mu)^2 + E(\mu - \theta)^2$$

$$= \frac{E(\hat{\theta} - \mu)^2}{\text{var}(\hat{\theta})} + \frac{E(\mu - \theta)^2}{\text{bias}^2}$$

Expected value  
 of Constant  
 Constant  
 $E(\mu - \theta)^2$   
 $= (\mu - \theta)^2$

Passing number is 2.1  
 Passing number is 8

Want minimum  
 What if the n is

180000

## → Bayesian Linear Regression

- In the Bayesian view point we formulate linear regression using probability distribution rather than point estimates.
- The response,  $y$ , is not estimated as a single value, but is assumed to be drawn from a probability distribution.
- The model for Bayesian Linear Regression is the response sampled from a normal distribution
$$y \sim N(\beta^T x, \sigma^2 I)$$
- The output  $y$  is generated from a normal (Gaussian) Distribution characterized by a mean and variance.
- The mean for linear regression is the transpose of the weight matrix multiplied by the predictor matrix
$$\text{test predictions of events lead to a particular outcome}$$
- The variance is the square of the standard deviation or (multiplied by the identity matrix)

→ the aim of Bayesian linear regression is not to find the single best value of the model parameters, but rather to determine the posterior distribution for the model parameters.

→ the posterior probability of the model parameters is conditional upon the training inputs and outputs.

$$P(\beta|y, x) = \frac{P(y|x|\beta) * P(\beta)}{P(y, x)}$$

→ Here  $P(\beta|y, x)$  is the posterior probability distribution of the model parameters given the inputs and outputs.

→  $P(\beta|y, x)$  = Posterior of model parameters given specific best data distribution given situation in data.

$P(y, x|\beta)$  = Likelihood of the response features given the model and predictors

Features.

$P(\beta)$  = Prior probability of the Model parameters

$P(y, x)$  = Normalizing constant independent of Model.

(constant of proportionality)

## → Naive Bayes

- 1) For the given dataset Apply Naive-Bayes algorithm and predict the outcome for a car = {Red, Domestic, SUV}

Color	Type	Origin	stolen
Red	sports	Domestic	Yes
Red	sports	Domestic	No
Red	sports	Domestic	Yes
Yellow	sports	Domestic	No
Yellow	sports	imported	Yes
Yellow	SUV	imported	No
Yellow	SUV	imported	Yes
Yellow	SUV	Domestic	No
Red	SUV	imported	No
Red	sports	imported	Yes

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \rightarrow \text{prior}$$

↓  
Posterior      ↓  
                    Marginal

$x = \{\text{Red}, \text{Domestic}, \text{SUV}\} \rightarrow \text{Is it stolen}$

$$P(\text{Red} | \text{Yes}) = \frac{P(\text{Yes} | \text{Red}) P(\text{Red})}{P(\text{Yes})}$$

$$= \frac{\frac{3}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{3}{5}$$

$$P(\text{Domestic} | \text{yes}) = \frac{P(\text{yes} | \text{Domestic}) \cdot P(\text{Domestic})}{P(\text{yes})}$$

$$= \frac{\frac{2}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{2}{5}$$

$$P(\text{SUV} | \text{yes}) = \frac{P(\text{yes} | \text{SUV}) \cdot P(\text{SUV})}{P(\text{yes})}$$

$$= \frac{\frac{2}{5} \times \frac{4}{10}}{\frac{2}{5}} = \frac{4}{10}$$

$$P(\text{Red} | \text{no}) = \frac{P(\text{no} | \text{Red}) \cdot P(\text{Red})}{P(\text{no})}$$

$$= \frac{\frac{2}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{2}{5}$$

$$P(\text{Domestic} | \text{no}) = \frac{P(\text{no} | \text{Domestic}) \cdot P(\text{Domestic})}{P(\text{no})}$$

$$= \frac{\frac{3}{5} \times \frac{5}{10}}{\frac{5}{10}} = \frac{3}{5}$$

$$P(\text{SUV} | \text{no}) = \frac{P(\text{no} | \text{SUV}) \cdot P(\text{SUV})}{P(\text{no})}$$

$$= \frac{\frac{3}{4} \times \frac{4}{10}}{\frac{5}{10}} = \frac{3}{5}$$

$$P(\text{Red}|\text{Yes}) = \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} = \frac{6}{125}$$

$$P(\text{Yes}|X) = P(\text{Yes}) \cdot P(\text{Red}|\text{Yes}) \cdot P(\text{Survived})$$

$$P(\text{Yes}|\text{Yes}) = P(\text{Yes}) \cdot P(\text{Survived}) \\ P(\text{Domestic}|\text{Yes}) \cdot P(\text{Survived})$$

$$= \frac{8}{10} \times \frac{3}{5} \times \frac{2}{5} \times \frac{1}{5} = \frac{3}{125}$$

$$= 0.024$$

$$P(\text{Red}|\text{No}) = P(\text{No}) \cdot P(\text{Red}|\text{No}) \cdot P(\text{Domestic}|\text{No})$$

$$P(\text{No}|X) = P(\text{No}) \cdot P(\text{Survived}|\text{No})$$

$$= \frac{2}{10} \times \frac{2}{5} \times \frac{3}{5} \times \frac{3}{5}$$

$$= \frac{9}{125} = 0.072$$

### → Decision Tree

→ Decision Tree is a supervised learning

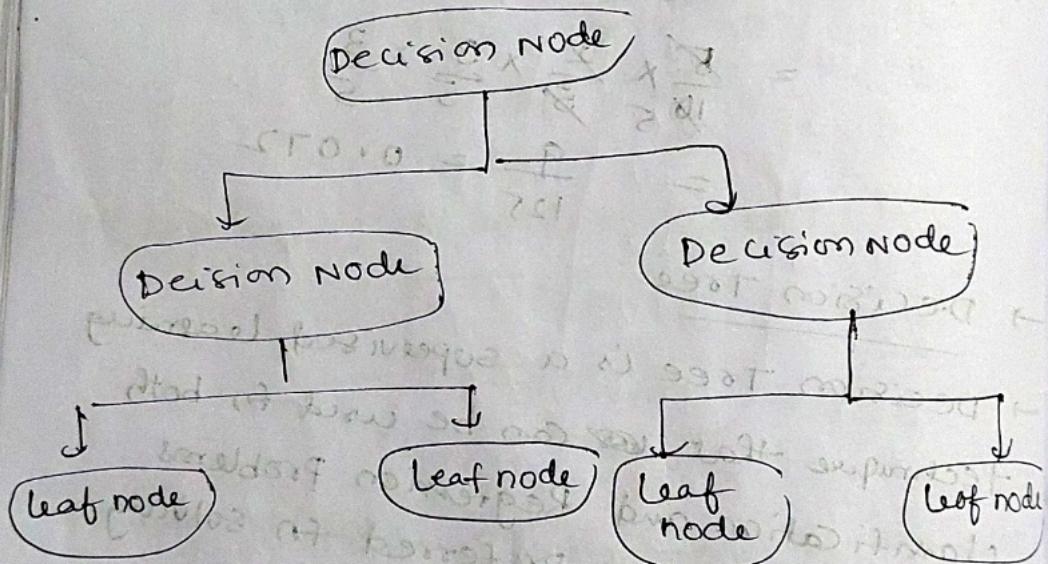
technique that can be used for both

classification and regression problems.

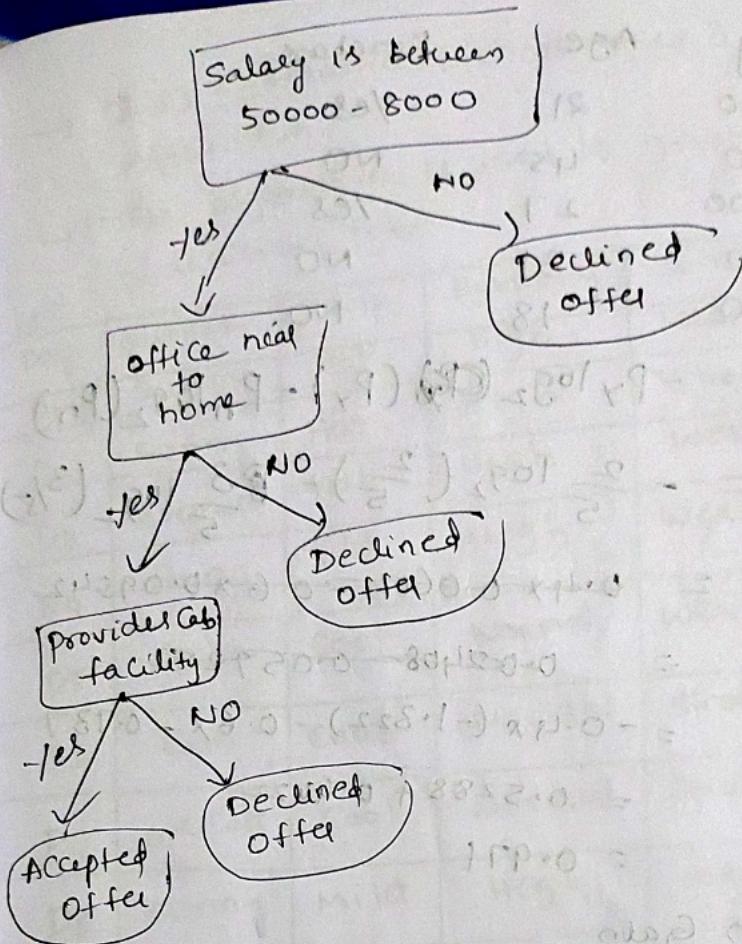
but mostly it is preferred for solving classification problems.

→ It is a tree structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.

- In a decision tree, there are two types which are the Decision Node, and leaf node. Decision nodes are used to make any decision and have multiple branches, whereas leaf nodes are the output of those decisions and do not contain any further branches.
- It is a graphical representation for getting all the possible solutions to a problem/decision based on given conditions.



Ex:- Suppose there is a candidate who gets a job offer and wants to decide whether he should accept the offer or not.



→ Entropy :-

Entropy is nothing but the measure of disorder or measure of impurity.

$$E(S) = - \sum_{i=1}^C P_i \log_2 P_i$$

$P_i$  is frequentist probability of an

element "i" in our data.

Salary	Age	Purchase
20000	21	Yes
10000	45	No
60000	27	Yes
15000	31	No
12000	18	No

$$\begin{aligned}
 E(S) &= -P_Y \log_2(P_Y) - P_N \log_2(P_N) \\
 &= -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) \\
 &= 0.4 \times 0.602 = 0.6 \times 0.09542 \\
 &\approx 0.2408 - 0.059282 \\
 &= -0.4 \times (-1.322) - 0.6 \times 0.737 \\
 &= 0.5288 + 0.4622 \\
 &= 0.971
 \end{aligned}$$

### Information Gain

→ Information Gain is a metric used to

train Decision Trees for a particular problem

→ Constructing a decision tree is all about finding attribute that returns the highest information gain

$$\begin{aligned}
 \text{Gain}(S, A) &\rightarrow \text{Attribute } A \text{ splits set } S \text{ into } n \text{ subsets} \\
 &= \text{Entropy}(S) - \sum_{\text{values}(A)} \frac{|S_v|}{|S|} \text{Entropy}
 \end{aligned}$$

→ ID3 (Iterative Dichotomiser 3)  
Algorithm

→ This algorithm is used to generate  
a decision tree.

Ex-	Day	outlook	Temp	Humidity	wind	playTennis :
D <sub>1</sub>	sunny	Hot	High	weak	NO	
D <sub>2</sub>	sunny	Hot	High	Strong	NO	
D <sub>3</sub>	overcast	Hot	High	weak	Yes	
D <sub>4</sub>	Rain	Mild	High	weak	yes	
D <sub>5</sub>	Rain	Cool	Normal	weak	yes	
D <sub>6</sub>	Rain	Cool	Normal	Strong	NO	
D <sub>7</sub>	overcast	Cool	Normal	Strong	Yes	
D <sub>8</sub>	sunny	Mild	High	weak	NO	
D <sub>9</sub>	sunny	Cool	Normal	weak	Yes	
D <sub>10</sub>	Rain	Mild	Normal	weak	Yes	
D <sub>11</sub>	sunny	Mild	Normal	Strong	Yes	
D <sub>12</sub>	overcast	Mild	High	Strong	Yes	
D <sub>13</sub>	overcast	Hot	High	Strong	Yes	
D <sub>14</sub>	Rain	Mild	High	Strong	NO	

→ Attribute: outlook = {sunny, overcast, rain}

values(outlook) = sunny, overcast, rain

$$S = [9+, 5-]$$

Whole dataset

$$\text{Entropy}(S) = -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) + \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$

$$= 0.94$$

$$S_{\text{sunny}} = [2+, 3-]$$

$$\text{Entropy}(S_{\text{sunny}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.97$$

$$S_{\text{overcast}} = [4+, 0-]$$

$$\text{Entropy}(S_{\text{overcast}}) = 0$$

$$= -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = 0$$

$$S_{\text{rain}} = [3+, 2-]$$

$$\text{Entropy}(S_{\text{rain}})$$

$$= -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.971$$

$$\text{Gain}(S, \text{outlook}) = \text{Entropy}(S) - \sum_{\text{outlook}} \frac{|\text{S}_v|}{|S|} \text{Entropy}(\text{S}_v)$$

(e.g. if S = {sunny, overcast, rain})

$\{\text{S}_v = \text{sunny, overcast, rain}\}$

$$\begin{aligned} &= \text{Entropy}(S) - \frac{5}{14} \text{Entropy}(\text{S}_{\text{sunny}}) - \frac{4}{14} \text{Entropy}(\text{S}_{\text{overcast}}) \\ &\quad - \frac{5}{14} \text{Entropy}(\text{S}_{\text{rain}}) \\ &= 0.94 - \frac{5}{14} \times 0.971 - \frac{4}{14} \times 0.971 \\ &= 0.2464 \end{aligned}$$

→ Attribute: Temp  
values(Temp) = HOT, mild, cool

$$\begin{aligned} S_{\text{HOT}} &= [2+, 2-] \\ \text{Entropy}(S_{\text{HOT}}) &= -\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} = 1.0 \\ S_{\text{mild}} &= [4+, 2-] \\ \text{Entropy}(S_{\text{mild}}) &= -\frac{4}{6} \log_2 \frac{4}{6} + \frac{2}{6} \log_2 \frac{2}{6} = 0.9183 \end{aligned}$$

$$S_{\text{cool}} = [3+, 1-]$$

$$\begin{aligned} S_{\text{cool}} &= [3+, 1-] \\ \text{Entropy}(S_{\text{cool}}) &= -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \\ &= 0.8113 \end{aligned}$$

- Gain( $S$ , Temp)

$$= \text{Entropy}(S) - \sum_{v \in \{\text{Hot, Mild, Cool}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= \text{Entropy}(S) - \frac{4}{14} \text{Entropy}(S_{\text{Hot}}) - \frac{6}{14} \text{Entropy}(S_{\text{Mild}}) - \frac{4}{14} \text{Entropy}(S_{\text{Cool}})$$

$$= 0.94 - \frac{4}{14} \times 1.0 - \frac{6}{14} \times 0.9183 - \frac{4}{14} \times 0.8113$$

$$= 0.0289$$

→ Attribute (Humidity) = High, Normal

~~so~~ = 0.5

$$S_{\text{High}} = [3+, 4-]$$

$$\text{Entropy}(S_{\text{High}}) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7}$$

$$S_{\text{Normal}} = [6+, 1-] = 0.9852$$

$$\text{Entropy}(S_{\text{Normal}}) = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7}$$

$$= 0.5916$$

Gain( $S$ , Humidity)

$$= \text{Entropy}(S) - \sum_{v \in \{\text{High, Normal}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Gain(S, Humidity)  $\rightarrow$   $H = -\sum p_i \log_2 p_i$

$$= \text{Entropy}(S) - \frac{1}{14} \text{Entropy}(S_{\text{High}}) -$$

$$\frac{1}{14} \text{Entropy}(S_{\text{Normal}})$$

Gain(S, Humidity)

$$= 0.94 - \frac{1}{14} \times 0.9852 - \frac{1}{14} \times 0.5916$$

$$= 0.1516$$

→ values(wind) = Strong, Weak

$$S_{\text{Strong}} = [3+, 3-]$$

$$\text{Entropy}[S_{\text{Strong}}] = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= 1$$

$$S_{\text{Weak}} = [6+, 2-]$$

$$\text{Entropy}[S_{\text{Weak}}] = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$= 0.8113$$

Gain(S, wind)

$$= \text{Entropy}(S) - \sum_{v \in \{\text{Strong, Weak}\}} \frac{1}{14} \text{Entropy}(S_v)$$

Gain(S, wind)

$$= \text{Entropy}(S) - \frac{6}{14} \text{Entropy}(S_{\text{Strong}}) -$$

$$\frac{8}{14} \text{Entropy}(S_{\text{Weak}})$$

$$\text{Gain}(S, \text{wind}) = 0.94 - \frac{6}{14} \times 1.0 = \frac{8}{14} \times 0.8113$$

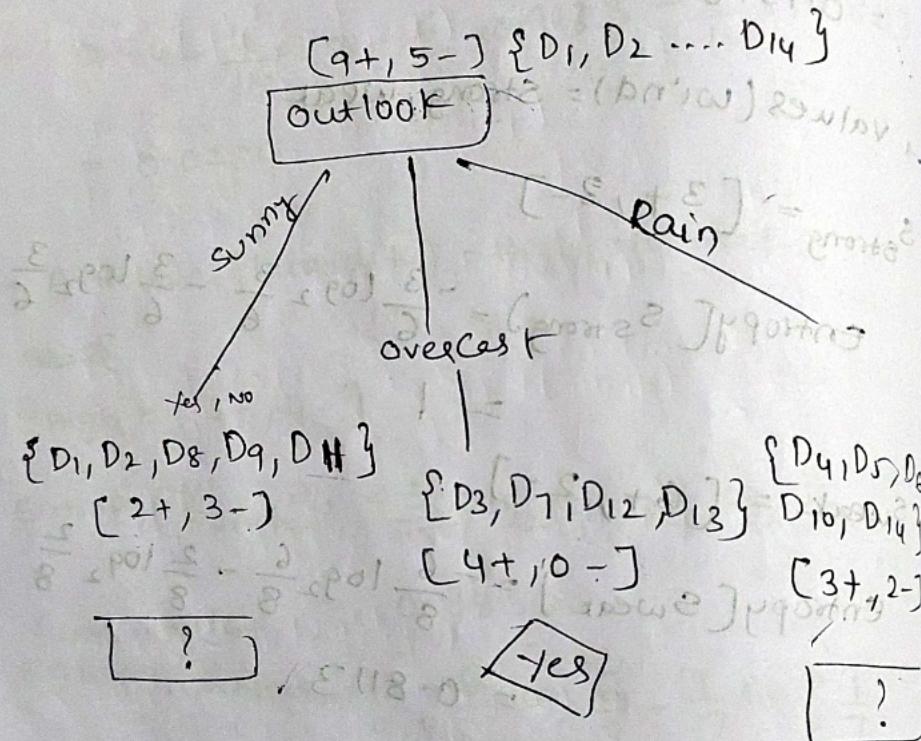
$$= 0.0478$$

$\rightarrow \text{Gain}(S, \text{outlook}) = 0.2468 \rightarrow$  (highest root node)

$$\text{Gain}(S, \text{Temp}) = 0.0289$$

$$\text{Gain}(S, \text{Humidity}) = 0.1516$$

$$\text{Gain}(S, \text{Wind}) = 0.0478$$



Day	Temp	Humidity	Wind	Play
D <sub>1</sub>	Hot	High	Weak	Tennis
D <sub>2</sub>	Hot	High	Strong	No
D <sub>8</sub>	Mild	High	Weak	No
D <sub>9</sub>	Cool	Normal	Weak	Yes
D <sub>11</sub>	Mild	Normal	Strong	Yes

→ Attributes Temp

values(Temp) = Hot, mild, Cool

$$S_{\text{Sunny}} = [2+, 3-]$$

$$\text{Entropy}(S_{\text{Sunny}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.97$$

$$S_{\text{Hot}} = [0+, 2-]$$

$$\text{Entropy}(S_{\text{Hot}}) = -\frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2}$$

$$= 0$$

$$S_{\text{mild}} = [1+, 1-]$$

$$\text{Entropy}(S_{\text{mild}}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 1$$

$$S_{\text{cool}} [1+, 0-]$$

$$\text{Entropy}(S_{\text{cool}}) = -\frac{1}{1} \log_2 1 - \frac{0}{1} \log_0$$

$$= 0$$

$$\text{Gain}(S_{\text{Sunny}}, \text{Temp}) = \frac{1}{|S|} \text{Entropy}(S_{\text{v}})$$

$$= \text{Entropy}(S) - \sum_{v \in \{\text{Hot}, \text{mild}, \text{Cool}\}} \frac{1}{|S|}$$

$$= \text{Entropy}(S) - \frac{2}{5} \text{Entropy}(S_{\text{Hot}}) - \frac{2}{5} \text{Entropy}(S_{\text{mild}})$$

$$- \frac{1}{5} \text{Entropy}(S_{\text{cool}})$$

$$= 0.97 - \frac{2}{5} \times 0 - \frac{2}{5} \times 1 - \frac{1}{5} \times 0 = 0.570$$

Attribute : Humidity

values (Humidity) = High, Normal

$$S_{\text{High}} = [0+, 3-]$$

$$\text{Entropy} = - \frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3}$$

$$S_{\text{Normal}} = [2+, 0-]$$

$$\text{Entropy}(S_{\text{Normal}}) = 0$$

Gain(  $S_{\text{sunny}}$ , Humidity )

$$= \text{Entropy}(S) - \sum_{v \in \{\text{High}, \text{Normal}\}} \frac{1}{S_v} \text{Entropy}(S_v)$$

$$= \text{Entropy}(S) - \frac{3}{5} \text{Entropy}(S_{\text{High}}) - \frac{2}{5} \text{Entropy}(S_{\text{Normal}})$$

$$= 0.97 - \frac{3}{5} \times 0 - \frac{2}{5} \times 0 = 0.97$$

→ Attribute : Wind

values (wind) = Strong, Weak

$$S_{\text{Strong}} = [1+, 1-]$$

$$\text{Entropy}(S_{\text{Strong}}) = 1$$

$$S_{\text{Weak}} = [1+, 2-]$$

$$\text{Entropy}(S_{\text{Weak}}) = - \frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.9183$$

$$\text{Gain}(S \text{ sunny, wind}) = \text{Entropy}(S) - \sum_{\text{ve (strong, weak)}} \frac{1}{18} \text{Entropy}(S_v)$$

$$= \text{Entropy}(S) - \frac{2}{5} \text{Entropy}(S \text{ strong})$$

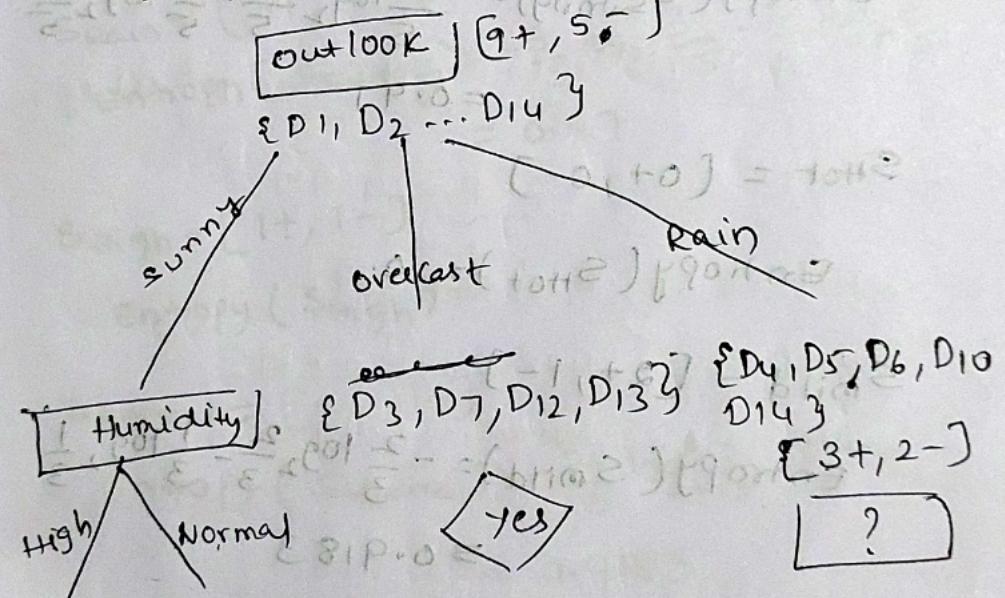
$$= \frac{3}{5} \text{Entropy}(S \text{ weak})$$

$$= 0.97 \times \frac{2}{5} \times 1.00 - \frac{3}{5} \times 0.918 = 0.0192$$

$$\rightarrow \text{Gain}(S \text{ sunny, Temp}) = 0.570$$

$$\text{Gain}(S \text{ sunny, Humidity}) = 0.97 \checkmark \text{ highest}$$

$$\text{Gain}(S \text{ sunny, wind}) = 0.0192$$



$$\{D_1, D_2, D_8\} \quad \{D_9, D_{11}\}$$

yes:  $[-1, +1] = 1.00$

no:  $1 - (100/2) = 0.50$

Attribute: Temp

values: ~~Temp~~ = HOT,

sunny = [2+, 3-]

$$\text{Entropy}(S \text{ sunny}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.97$$

Day	Temp	Humidity	Wind	Play Tennis
D4	mild	High	weak	Yes
D5	Cool	Normal	weak	Yes
D6	Cool	Normal	Strong	No
D10	mild	Normal	weak	Yes
D14	Mild	High	Strong	No

Attribute: Temp

values(Temp) = Hot, mild, Cool

$$S_{\text{Rain}} = [3+, 2-]$$

$$\text{Entropy}(S_{\text{sunny}}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.97$$

$$S_{\text{Hot}} = [0+, 0-]$$

$$\text{Entropy}(S_{\text{Hot}}) = 0.0$$

$$S_{\text{mild}} = [2+, 1-]$$

$$\text{Entropy}(S_{\text{mild}}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.9183$$

$$S_{\text{Cool}} = [1+, 1-]$$

$$\text{Entropy}(S_{\text{Cool}}) = 1.$$

$$\text{Gain}(\text{Rain}, \text{Temp}) = \text{Entropy}(S) - \sum_{v \in \{\text{Hot}, \text{Mild}, \text{Cool}\}} \frac{1}{|S|} \text{Entropy}(S_v)$$

$$= \text{Entropy}(S) - \frac{0}{5} \text{Entropy}(S_{\text{Hot}}) - \frac{3}{5} \text{Entropy}(S_{\text{Mild}})$$

$$\Rightarrow -\frac{2}{5} \text{Entropy}(S_{\text{Normal}})$$

$$\delta = 0.97 - \frac{0}{5} \times 0 - \frac{3}{5} \times 0.918 = \frac{2}{5} \times 0.18 = 0.072$$

$\rightarrow$  Attribute: Humidity

$$\text{Rain} = [3+, 2-]$$

$$\text{Entropy}(S_{\text{Sunny}}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

$$= 0.97$$

$$S_{\text{High}} = [1+, 1-]$$

$$\text{Entropy}(S_{\text{High}}) = 1 \text{mild, 1 normal}$$

$$S_{\text{Normal}} = [2+, 1-]$$

$$\text{Entropy}(S_{\text{Normal}}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3}$$

$$= 0.9183$$

$$\text{Gain}(S_{\text{Rain}}, \text{Humidity})$$

$$= \text{Entropy}(S) - \sum_{v \in \{\text{High}, \text{Normal}\}} \frac{1}{|S|} \text{Entropy}(S_v)$$

$$\therefore \text{Entropy}(S) - \frac{2}{5} \text{Entropy}(S_{\text{High}})$$

$$- \frac{3}{5} \text{Entropy}(S_{\text{Normal}})$$

$$= 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.918 = 0.0192$$

→ Attribute : Wind  
values (wind) = strong, weak

(v2)  $S_{rain} = [3+, 2-]$

$$S_{rain} = \frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5}$$

$$\text{Entropy}(S_{\text{sunny}}) = \frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5} = 0.97$$

Strong = [0+, 2-]

$$\text{Entropy}(S_{\text{strong}}) = 0$$

Weak = [3+, 0-]

$$\text{Entropy}(S_{\text{weak}}) = 0$$

$$\text{Gain}(S_{\text{rain}}, \text{Wind}) = \text{Entropy}(S) - \sum_{v \in \{\text{Strong}, \text{Weak}\}} \frac{1}{2} \text{Entropy}(S_v)$$

$$= \text{Entropy}(S) - \frac{2}{5} \text{Entropy}(S_{\text{strong}}) - \frac{3}{5} \text{Entropy}(S_{\text{weak}})$$

$$= 0.97 - \frac{2}{5} \times 0 - \frac{3}{5} \times 0 = 0.97$$

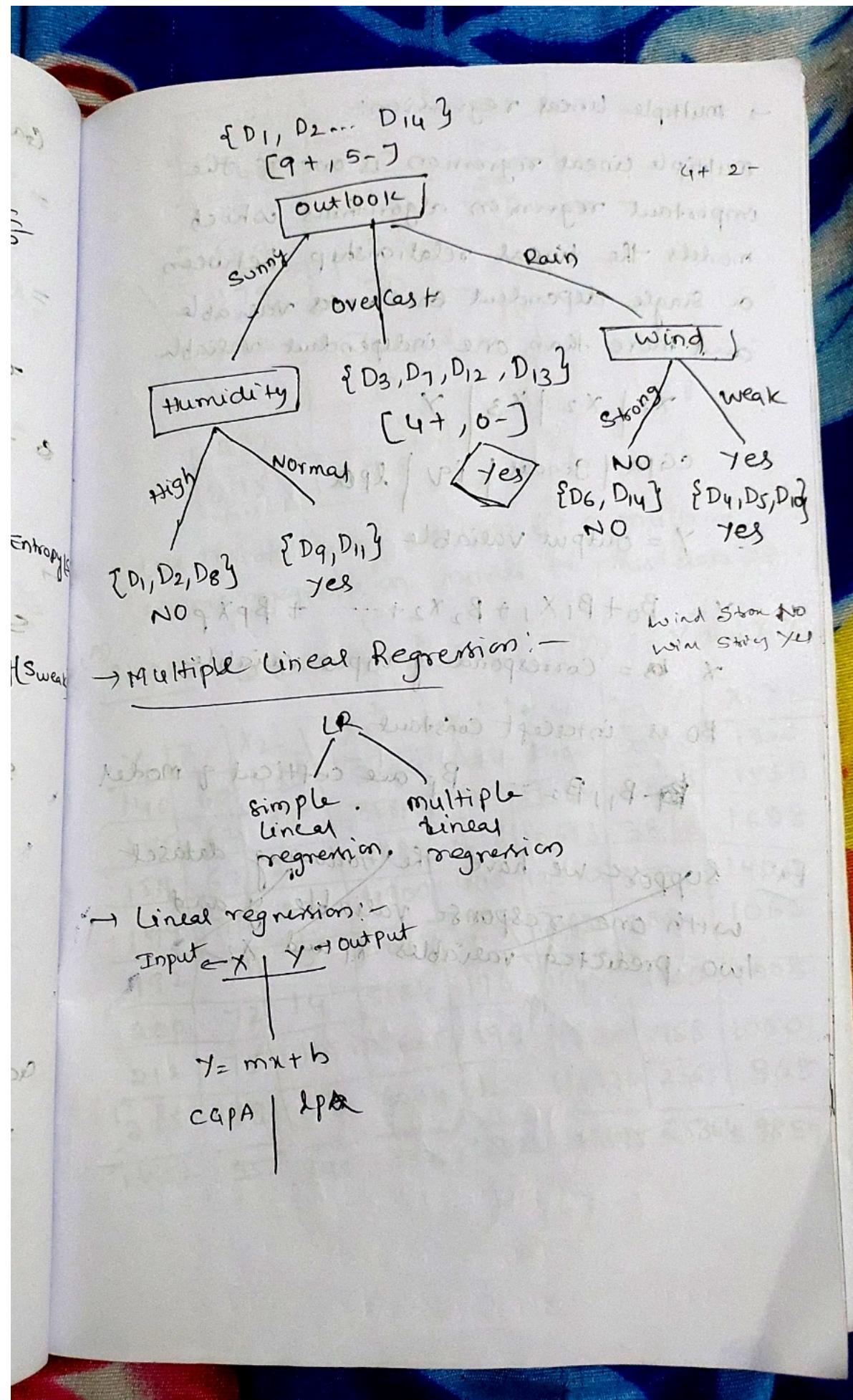
$$\rightarrow \text{Gain}(S_{\text{rain}}, \text{Temp}) = 0.0192$$

$$\text{Gain}(S_{\text{rain}}, \text{Humidity}) = 0.0192$$

$$\text{Gain}(S_{\text{rain}}, \text{Wind}) = 0.97$$

EPD =

(v2)  $\frac{1}{2} P_{\text{strong}}$  (v2)  $\frac{1}{2} P_{\text{weak}}$



→ Multiple linear regression:-

Multiple linear regression is one of the important regression algorithms which models the linear relationship between a single dependent continuous variable and more than one independent variable.

Now  $x_1 | x_2 | x_3 | Y$

Ex: CGPA | gender | iq | lpa

Y = Output variable

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

$\times \beta$  = Corresponding input variable

$\beta_0$  is intercept constant

$\beta_1, \beta_2, \dots, \beta_p$  are coefficient of Model

Ex:- Suppose we have the following dataset with one response variable y and two predicted variables  $x_1$  and  $x_2$

$$x_1 + x_2 = 1$$

$$x_1 | x_2$$

~~Step by step calculation of Regression sums~~

$$\sum x_1^2 - \frac{(\sum x_1)^2}{n}$$

$\hat{y}_P =$  predict the value of  $y$  given  $x_1$  and  $x_2$

Subject	$y$	$x_1$	$x_2$
1	-3.7	3	8
2	3.5	4	5
3	2.5	5	7
4	11.5	6	1
5	5.7	2	2

6

$$B_0 = \bar{y} - B_1 \bar{x}_1 - B_2 \bar{x}_2$$

intercept mean of  $y$  mean of  $x_1$  mean of  $x_2$

variance of  $x_1$  and  $x_2$

$$B_1 = \frac{(\sum x_1^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$B_2 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$\sum x_1^2 = \sum x_1^2 - \frac{(\sum x_1)^2}{N}$$

Subject	$\Sigma Y$	$x_1$	$x_2$	$x_1^2$	$x_2^2$	$x_1 x_2$	$x_1 Y$	$x_2 Y$
1	-3.7	3	8	9	64	24	-11.1	-29.6
2	3.5	4	5	16	25	20	14	17.5
3	2.5	5	7	25	49	35	12.5	17.5
4	11.5	6	3	36	9	18	69	34.5
5	5.7	2	1	4	1	2	11.4	5.7
	19.5	20	24	90	148	99	95.8	45.6

$$\sum x_1^2 = \sum x_1^2 - \frac{(\sum x_1)^2}{N} \rightarrow \text{regression sum}$$

variance w.r.t to 22

$$= 90 - \frac{20 \times 20}{5} = 10$$

$$\sum x_2^2 = \sum x_2^2 - \frac{(\sum x_2)^2}{N}$$

$$= 148 - \frac{24 \times 24}{5} = 32.8$$

$$\sum x_1 Y = \sum x_1 Y - \frac{(\sum x_1)(\sum Y)}{N}$$

$$= 95.8 - \frac{(20)(19.5)}{5} = 17.8$$

$$\sum x_2 Y = \sum x_2 Y - \frac{(\sum x_2)(\sum Y)}{N}$$

$$= 45.6 - \frac{(24)(19.5)}{5} = -48$$

$$\Sigma x_1 x_2 = \Sigma x_1 x_2 \frac{(\Sigma x_1)(\Sigma x_2)}{N}$$

$$= 90 - \frac{(32 \cdot 8)(3)}{5}$$

$$= 90 - \frac{(20 \cdot 24)}{5} = 3$$

$$\beta_1 = \frac{(\Sigma x_1^2)(\Sigma x_2 y) - (\Sigma x_1 x_2)(\Sigma x_2 y)}{(\Sigma x_1^2)(\Sigma x_2^2) - (\Sigma x_1 x_2)^2}$$

$$= \frac{10 \cdot 17.8 - 3(-48)}{10 \cdot 32.8 - 3^2}$$

$$= \frac{(32 \cdot 8)(17 \cdot 8) - (3)(-48)}{(10)(32 \cdot 8) - 3^2}$$

$$= 2.28$$

$$\beta_2 = \frac{(\Sigma x_1^2)(\Sigma x_2 y) - (\Sigma x_1 x_2)(\Sigma x_1 y)}{(\Sigma x_1^2)(\Sigma x_2^2) - (\Sigma x_1 x_2)^2}$$

$$= \frac{10(-48) - (3)(17 \cdot 8)}{(10)(32 \cdot 8) - (3)^2}$$

$$= -1.67$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}_1 - \beta_2 \bar{x}_2$$

$$= \frac{19.5}{5} - \frac{2.28 \times 20}{5} - \frac{-1.67 \times 24}{5}$$

$$= 2.796$$

final regression equation is

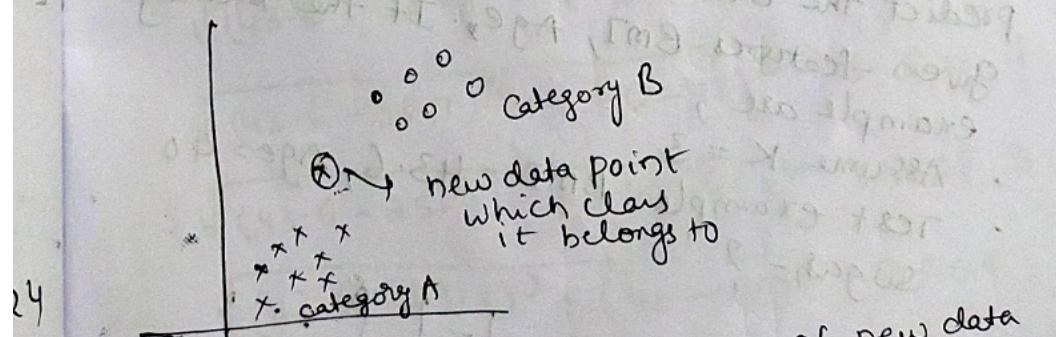
$$y = 2.796 + 2.28x_1 - 1.67x_2$$

now given  $x_1 = 3$  and  $x_2 = 2$   $y = ?$

$$y = 2.796 + 2.28 \times 3 - 1.67(2)$$

$$= 6.296$$

- KNN algorithm:-  
K Nearest Neighbour (K value must be odd)
- It is a supervised learning algorithm that classifies a new data point into the target class, depending on the features of its neighbouring data points.
- K-NN is a non-parametric algorithm which means it does not make any assumption on underlying data.
- It is called a lazy learning algorithm.



- we will calculate distance of new data point from every data point and assume particular  $K = 3$  (nearest three).

- nearest three data points - then we will choose the ones which are of maximum type.
- The  $K$  values will be decided based on trial and error so that it reduces the bias and variance.

Ex:-

	BMI	Age	Sugar	Rank
1	33.6	50	1	
2	26.6	30	0	
3	23.4	40	0	
4	43.1	67	0	
5	35.3	23	1	
6	35.9	67	1	
7	86.7	45	1	
8	25.7	46	0	3
9	23.3	29	0	
10	31	56	1	

- Apply  $K$  nearest neighbor classifier to predict the diabetic patient with the given features BMI, Age. If the training example are,
- Assume  $K = 3$
- Test example  $\text{BMI} = 43.6$ ,  $\text{Age} = 40$   
 $\text{Sugar} = ?$

$$86.7 - 1 - 1 \rightarrow \\ 37.6 - 2 - 1 \rightarrow \\ 23.3 - 3 - 0 \rightarrow$$

- ~~Distance~~  
 → First calculate the distance between the test instance and training instances  
 → Test example

BMI = 43.6, Age = 40, sugar = ?

→ Euclidean  
~~Distance~~ =  $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$

↓                    ↓                    ↓                    ↓  
 Test      Training      Test      Training  
 data      data      data      data  
 variable   variable   variable   variable

Distance	Rank
$\sqrt{(43.6 - 33.6)^2 + (40 - 50)^2}$	14.14 2
$\sqrt{(43.6 - 26.6)^2 + (40 - 30)^2}$	19.72
$\sqrt{(43.6 - 23.4)^2 + (40 - 40)^2}$	20.20
$\sqrt{(43.6 - 43.1)^2 + (40 - 67)^2}$	27.00
$\sqrt{(43.6 - 35.3)^2 + (40 - 23)^2}$	18.92
$\sqrt{(43.6 - 35.9)^2 + (40 - 67)^2}$	28.08
$\sqrt{(43.6 - 36.7)^2 + (40 - 45)^2}$	8.52 1
$\sqrt{(43.6 - 25.7)^2 + (40 - 46)^2}$	18.88 3
$\sqrt{(43.6 - 23.3)^2 + (40 - 29)^2}$	23.09
$\sqrt{(43.6 - 31)^2 + (40 - 56)^2}$	20.37

→ Out of 3 cases '2' cases having sugar (mamony sugar)

→ out of '3' nearest neighbours 2 are one  
and one is zero, so the new test  
example is classified as ~~overweight~~  
sugar = ~~but it's normal too~~  
~~almost not~~

Ex:-	K=3	height (cm.)	weight (kg.)	Classification
		1 167	51	underweight
		2 182	62	Normal
		3 176	69	Normal
		4 173	64	Normal
		5 172	65	underweight
		6 174	56	Normal
		7 169	58	Normal
		8 173	55	Normal
		9 170	57	Normal
		10 170		

→ Range	Distance	Rank	Classification
169	1.46 - 0P	1	→ normal
170	2	2	→ normal
173	3	3	→ normal

$$E = 88.81 \quad (D_1 - 0P) + (R \cdot dE - d \cdot EP)$$

$$PO \cdot EG \quad (P_2 - 0P) + (R \cdot dE - d \cdot EP)$$

$$PE \cdot OG \quad (P_2 - 0P) + (E \cdot EG - d \cdot EP)$$

$$D_{avg} = \frac{D_1 + D_2}{2} \quad E = \frac{E_1 + E_2}{2}$$

maximum likelihood and least squared.

To find the maximum likelihood hypothesis in bayesian learning

$$h_{ML} = \arg \max_{h \in H} P(D|h)$$

probability density function

Let us take training instances  $(x_1, x_2, \dots, x_n)$  and consider target values  $D = (d_1, d_2, \dots, d_m)$

∴ we rewrite  $P(D|h)$  as product of  $P(d_i|h)$

$$h_{ML} = \arg \max_{h \in H} \prod_{i=1}^n P(d_i|h)$$

By assuming normal distribution

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

variable mean      standard deviation

$$h_{ML} = \arg \max_{h \in H} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i-h)^2}$$

$$= \arg \max_{h \in H} \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i-h(x_i))^2}$$

$\log e^{-x} = -x$

→ use log function

$$= \arg \max_{h \in H} \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} \left( -\frac{1}{2\sigma^2} (d_i-h(x_i))^2 \right)$$

$$= \arg \max_{h \in H} \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi\sigma^2}} \left( -\frac{1}{2\sigma^2} (d_i-h(x_i))^2 \right)$$

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^n -\frac{1}{\sqrt{2\pi\sigma^2}} (d_i - h(x_i))^2$$

y constant

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^n -\frac{1}{\sqrt{2\pi\sigma^2}} \frac{1}{2\sigma^2} (d_i - h(x_i))^2$$

y constant can remove

$$= \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^n (d_i - h(x_i))^2$$

square

$$h_{ML} = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^n (d_i - h(x_i))^2$$

minimum error

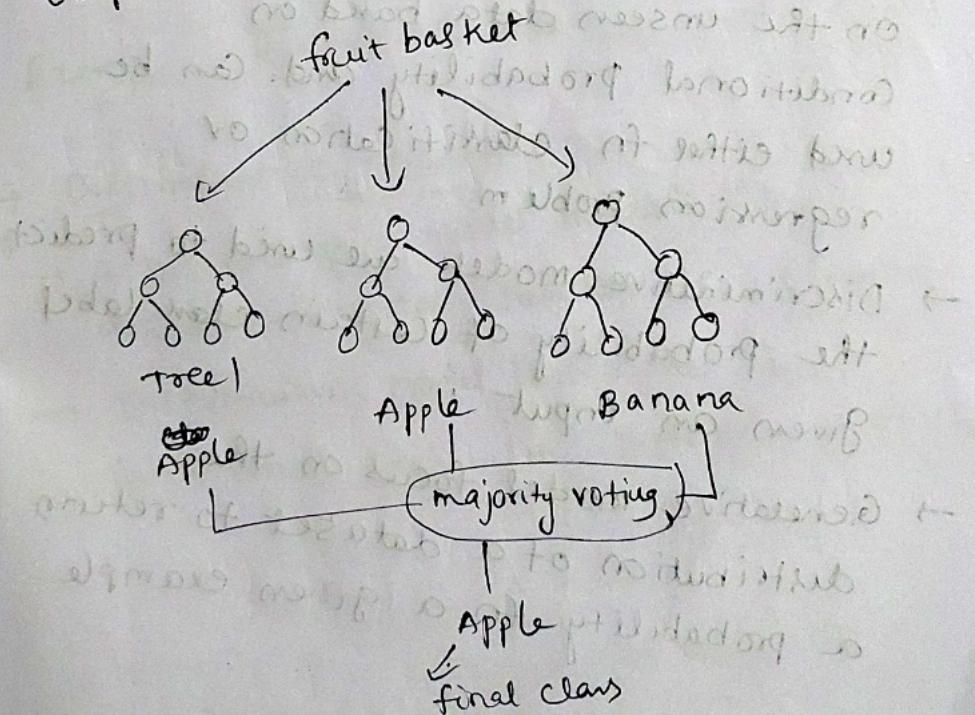
The maximum likelihood hypothesis which has minimum squared error between the hypothesis.

$$(d_i - h(x_i))^2$$

minimum error

## Random Forest model:-

- Random Forest is a supervised machine learning algorithm that is used widely in classification and Regression problem.
- It builds decision trees on different samples and takes their majority vote for classification.
- Ex:- Consider the fruit basket. ~~as the now 'n'~~ number of samples are taken from the fruit basket and individual decision tree is constructed for each sample. Each decision tree will generate an output. The final output is considered based on majority voting.



- Step 1: select random samples from a given data or training set
- Step 2: this algorithm will construct a decision tree for every training data
- Step 3: voting will take place by averaging the decision tree.
- Step 4: Final select the most voted prediction result as the final prediction result.

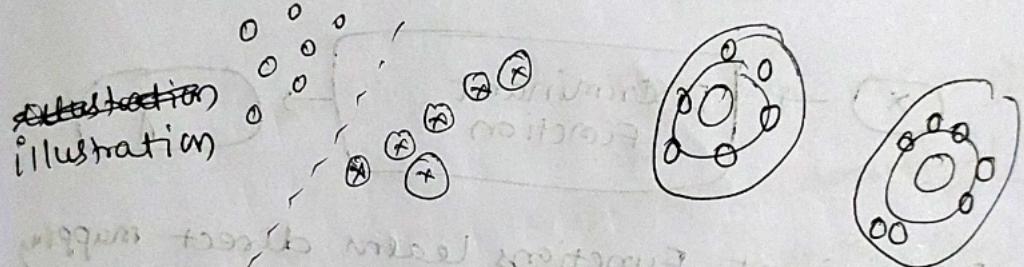
- Discriminative and Generative models
- machine learning models can be classified into two types of models.
- Discriminative model makes predictions on the unseen data based on conditional probability and can be used either for classification or regression problem.
- Discriminative models are used to predict the probability of certain class label given an input.
- Generative model, focus on the distribution of a dataset to return a probability for a given example

→ Discriminative models focus on predicting labels, while generative model focus on modeling the distribution of data.

→ ~~Generative vs~~ Discriminative vs Generative models

Goal	Discriminative Directly estimate $P(Y X)$	Generative model estimate $P(X Y)$ to then decide $P(Y X)$
------	----------------------------------------------	------------------------------------------------------------------

what's learned      Decision boundary      Probability of the data



~~Illustration~~  
illustration

Example      Regression, SVMs, Naive Bayes

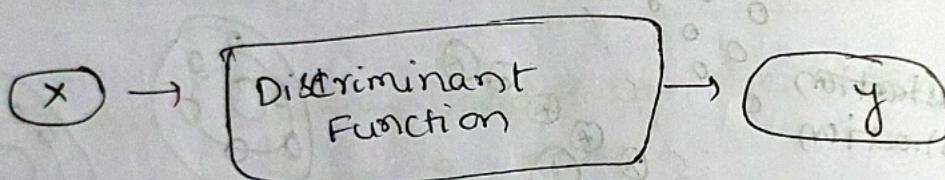
- Discriminative classifier
- A Discriminative model models the decision boundary between the classes
- Discriminative model learns the conditional probability distributions  $P(Y|X)$

$$\text{Posterior} = \frac{P(Y|X) \cdot P(X)}{\text{Prior}}$$

- Generative classifier
- A Generative model learns the joint probability distribution  $P(x,y)$
- It predicts the Conditional Probability with the help of Bayes theorem

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

- Discriminant Function:-



- Discriminant Functions learns direct mapping between features vector  $x$  and label  $y$ .
- In a binary classification set up with  $m$  features, the simplest discriminant function is very similar to the linear regression.

$$y = w_0 + w_1 x_1 + \dots + w_m x_m$$

$$J = w_0 + w^T x$$

label

$\downarrow$   
bias

Weight vector

$\swarrow$

feature vector

Geometrically, the simplest discriminant function  $y = w_0 + w^T x$  represents a hyperplane in  $m-1$  dimensional space where  $m$  is the number of features.

features ( $m$ )

1

Discriminant Function point

2

line

3

plane

4

Hyperplane in 3D space

:

$m$

Hyperplane in  $(m-1)D$  space

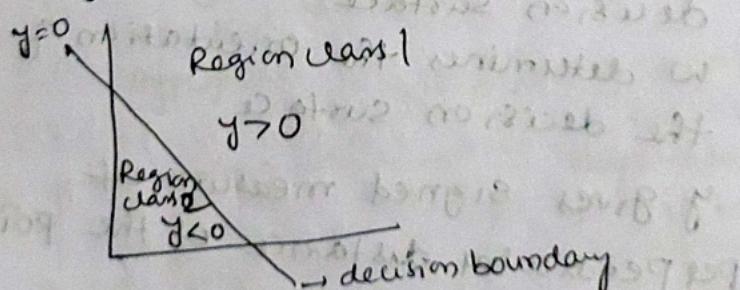
$\nearrow$   $(m-1)$  subspace in a  $m$  dimensional space.

$\rightarrow$   $(m-1)D$  hyperplane (which is line in this case) divides features space into two regions.

one for each class.

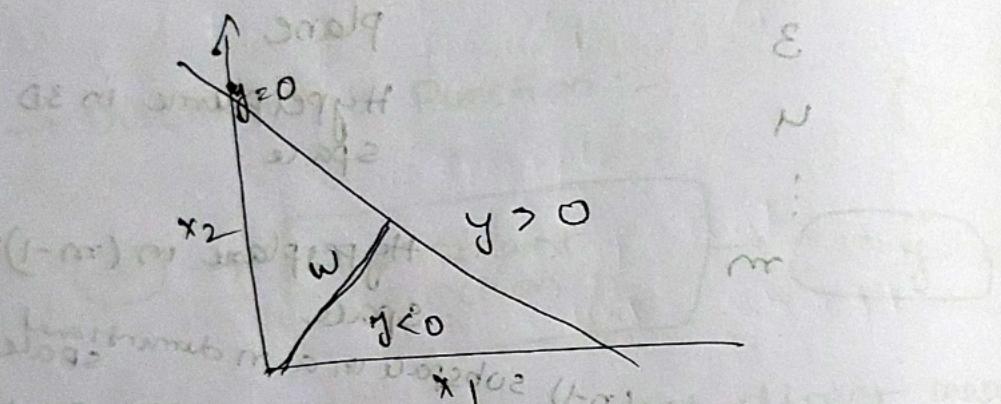
Region for Class 1, where  $y > 0$

Region for Class 2, where  $y < 0$



$\rightarrow$  The decision boundary between two classes is represented with  $(m-1)D$  hyperplane  $w_0 + w^T x = 0$ .

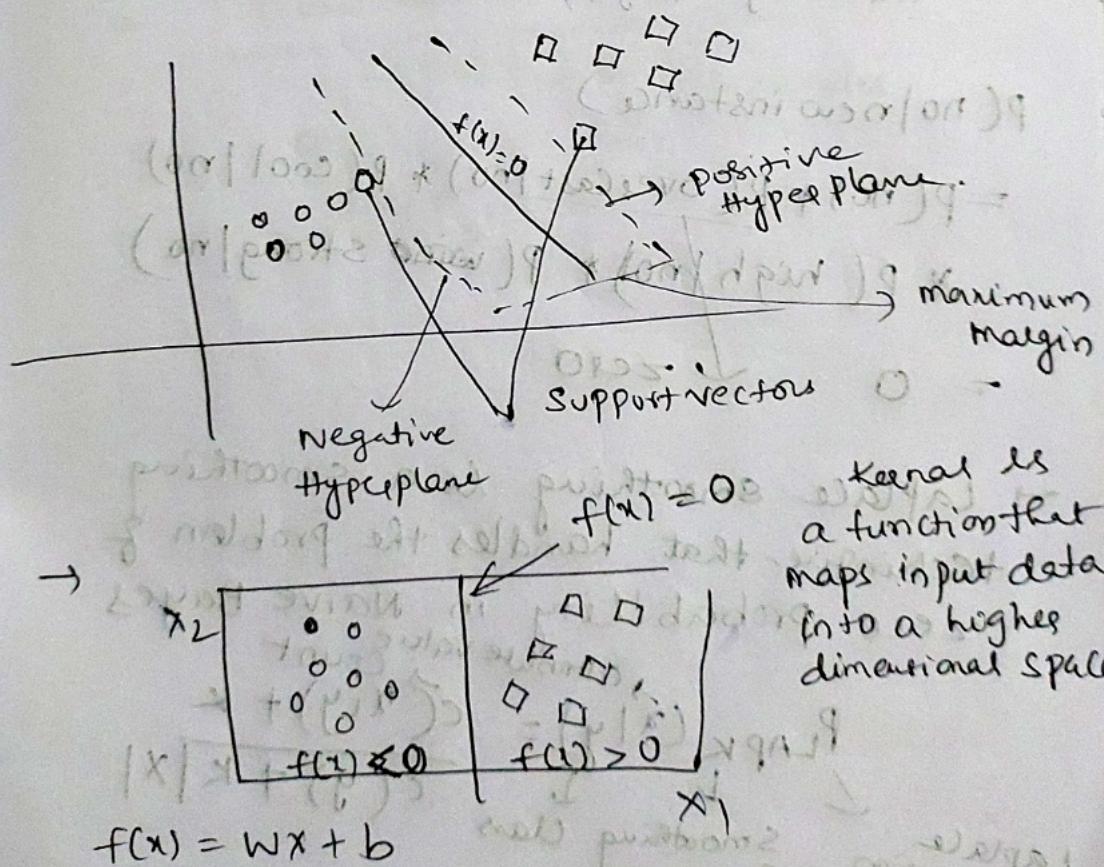
- Consider two points  $x^{(A)}$  and  $x^{(B)}$  on the decision surface, we will have,
- $$y^{(A)} = w_0 + w^T x^{(A)} = 0$$
- $$y^{(B)} = w_0 + w^T x^{(B)} = 0$$
- $$\Rightarrow w^T(x^{(A)} - x^{(B)}) = 0$$



- the vector  $w$  is orthogonal to every vector lying within the decision surface, hence it determines the orientation of the decision surface.
- $w_0$  determines the location of decision surface
- $w$  determines the orientation of the decision surface.
- $y$  gives signed measure of perpendicular distance of the point  $x$  from the decision surface
- Decision Surface divides feature space into two regions,

## Support vector machine

- SVM is machine learning algorithm which can be used for both classification or regression challenges.
- In this algorithm, we plot each data item as a point in  $n$ -dimensional space (where  $n$  is number of features you have) with the value of each feature being the value of a particular coordinate.
- Then, we perform classification by finding the hyper-plane that differentiate the two classes very well.



$w$  is the normal to the line,  $x$  is input vector and  $b$  is bias,  $w$  is known as the weight vector.

→ Laplace smoothing or

→ zero probability problem

outlook	-yes	NO
sunny	2/9	3/5
overcast	4/9	0/5
Rainy	3/9	3/5

→ clarity new example

(outlook = overcast, temp = cool, humidity = high, wind = strong)

$P(\text{no}|\text{new instance})$

$$= P(\text{no}) * P(\text{overcast}|\text{no}) * P(\text{cool}|\text{no})$$

$$* P(\text{high}|\text{no}) * P(\text{wind strong}|\text{no})$$

$$= 0 \quad \downarrow \text{zero}$$

→ Laplace smoothing is a smoothing technique that handles the problem of zero probability in Naive Bayes

$$P_{\text{LAP}_K}(x|y) = \frac{c(x|y) + k}{c(y) + k|x|}$$

Laplace Approximation Probability  
Smoothing Parameter  $k > 0$   
attribute value count  
class

$c(x|y)$  is the number of times  $x$  appears  
in  $\rightarrow$  examples

$c(y)$  is number of  $y$  examples

(suppose here  $y$  is no)

if  $y$  is yes number of yes examples

$x$  is number of values of that feature

Ex:- outlook ~~is~~ no of values

↓  
sunny, overcast, Rainy = 3

$$P_{LAP,K}(\text{outlook} = \text{overcast} | \text{no})$$

$$= \frac{c(\text{outlook} = \text{overcast} | \text{no}) + k}{c(\text{no}) + k * |x|} \rightarrow \text{modules}$$

$$= \frac{0 + 1}{5 + 1 * 3} = \frac{1}{8}$$