



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Unit I

Supervised Learning Algorithm

Analytics is a collection of techniques and tools used for creating value from data. Techniques include concepts such as artificial intelligence (AI), machine learning (ML), and deep learning (DL) algorithms.

AI → Algorithm and systems that exhibit human like intelligence.

ML → Subset of AI that can learn to perform a task with extracted data and/or models.

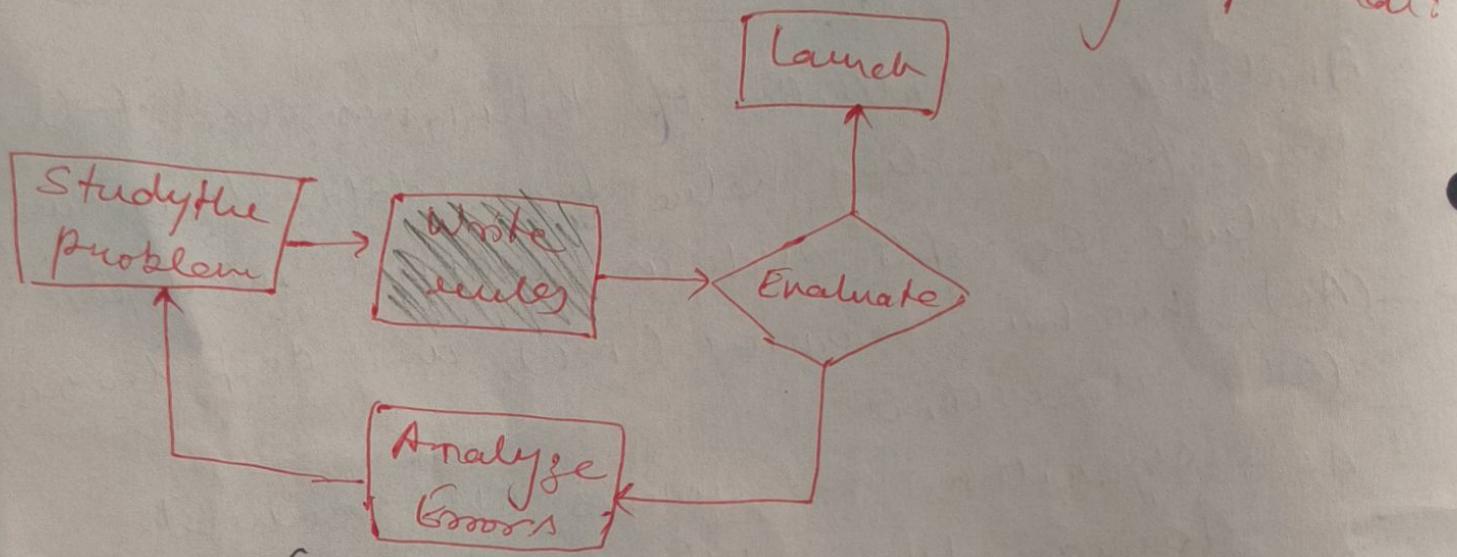
DL → Subset of ML that imitate the functioning of human brain to solve problems.



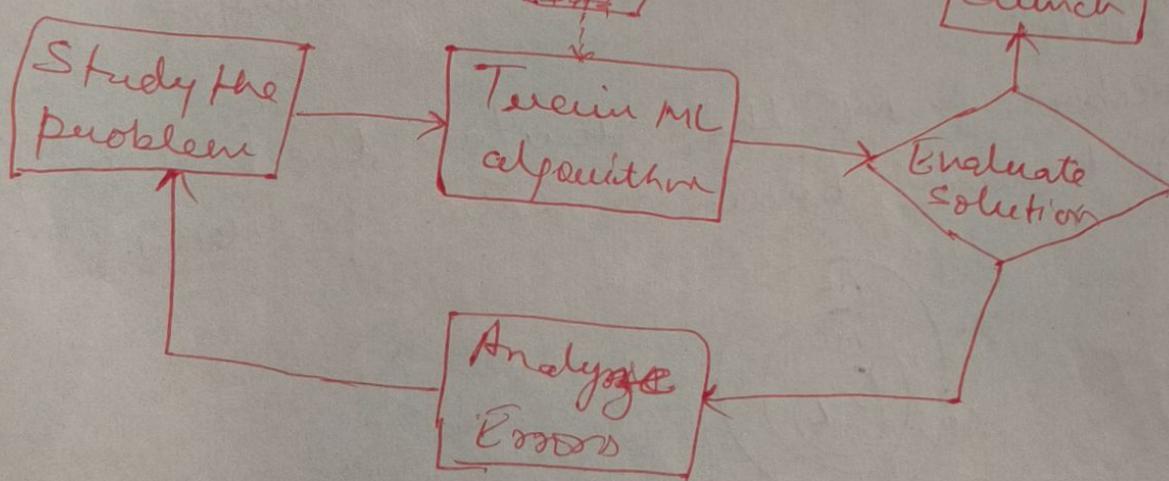
ML is a set of algorithm that have the capability to learn to perform tasks such as prediction and classification effectively using data.

★ An algorithm can be called a learning algorithm when it improves on a performance metric while performing a task, for e.g., accuracy of classification such as fraud, customer churn etc.

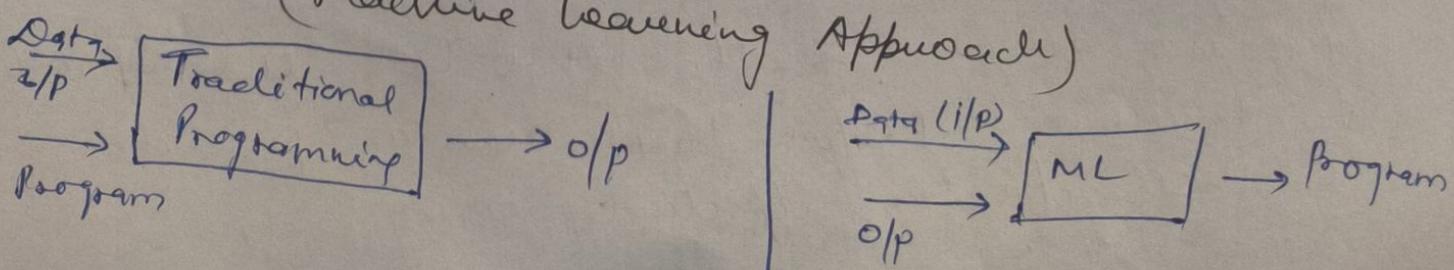
Differences between Traditional programming approach and Machine learning approach:



(Traditional programming Approach)



(Machine learning Approach)





→ Advantage of ML:-

PAGE NO.

1. Easily identifies trends and patterns : eg Amazon
2. No human intervention needed : eg anti-virus
They learn to filter new threats as they are recognized.
3. Continuous improvement eg: ~~the~~ weather forecasting
4. Handling multi-dimensional and multi-variety data

→ Disadvantage of ML:-

1. Data Acquisition
2. Time and Resources
3. Interpretation of Result
4. High error-susceptibility

→ Some important terms:-

- ↳ Labeled Data : LD has both the i/p and o/p Parameter in a complete machine-readable format. Used in Supervised learning.
- ↳ Unlabeled Data : is a piece of data that has not been tagged with labels for identifying the characteristics. It is raw data and used in basic unsupervised learning.

Machine learning algorithms are classified into four categories:

1. Supervised Machine learning
2. Unsupervised Machine learning
3. Semi Supervised learning
4. Reinforcement learning

Supervised learning

— Supervised learning is defined by its use of labeled datasets to train algorithm that to classify data or predict outcomes accurately.

Working
Supervised learning uses a training set to teach models to yield the desired o/p. This training dataset includes I/P and output o/p, which allow the model to learn over time,

The algorithm measures its accuracy through the loss function, adjusting until the error has been sufficiently minimized.

→ There are two varieties of supervised learning algorithms:

 [
 Classification Algorithm

 Regression Algorithm

 → Classification based supervised learning methods identify which category a set of data items belong to

 → Regression based supervised learning methods



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

try to predict output based on input variable.
 Examples of Supervised machine learning algorithms are as follows:

1. K- Nearest Neighbors ~~Algo~~ Algorithm
2. Linear Regression
3. Logistic Regression
4. Support Vector Machine
5. Decision Tree
6. Random Forest
7. Neural Network

Unsupervised learning Algo

Unsupervised learning is where we have only input data and no corresponding output variable. From that data, it discover pattern that help solve for ~~to~~ clustering or ~~assoc~~ association problem.

Eg: Grouping customers by their corresponding habit.

Unsupervised learning algo. divided in two categories:
 → Clustering: A clustering problem is where we want to discover the inherent grouping in the data. Such as based on purchasing behavior.

Ens of unsupervised learning:

* grouping the customers.

→ Association: An rule based method for finding relationships b/w variables in a given dataset. These method are frequently used for market basket analysis, allowing companies to better understand relationships b/w different products. Eg Amazon customer's

Eg. If unsupervised learning:

1. K-means for clustering problems.
2. Apriori algorithms for association rule learning problems.

Semi Supervised Learning:

Some algorithm deal with partially labeled training data along with a lot of unlabelled data.

Eg: Google photo (Photo archive)

Reinforcement learning: RL is a type of machine learning technique that enables an agent to learn in an interactive environment by trial and error using feedback from its own actions and experiences.

Eg: Google deep mind Chess engine.



Linear Regression Model:

Linear Regression is most well known algorithm in statistics and machine learning.

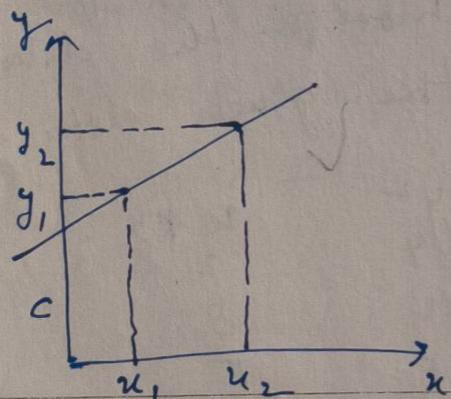
→ Linear regression is a linear model.

e.g. a model that assumes a linear relationship b/w the i/p variables (x) and the single o/p variables (y).

~~Note~~ The simplest form of a SLR equation with one dependent and one independent variable is represented by -

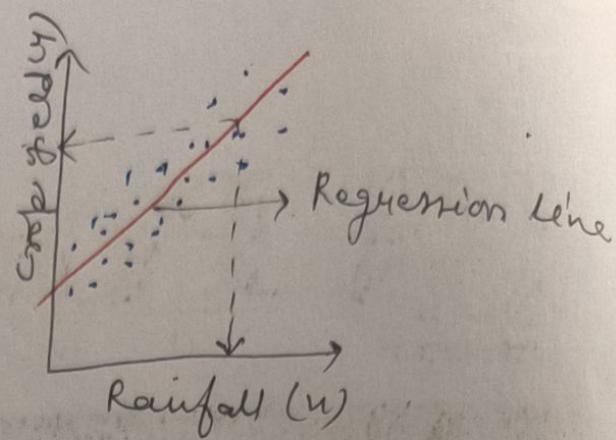
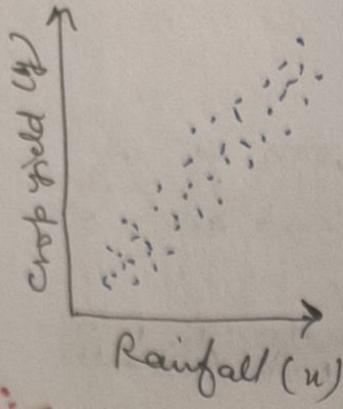
$$y = mx + c$$

where $m = \frac{y_2 - y_1}{x_2 - x_1}$



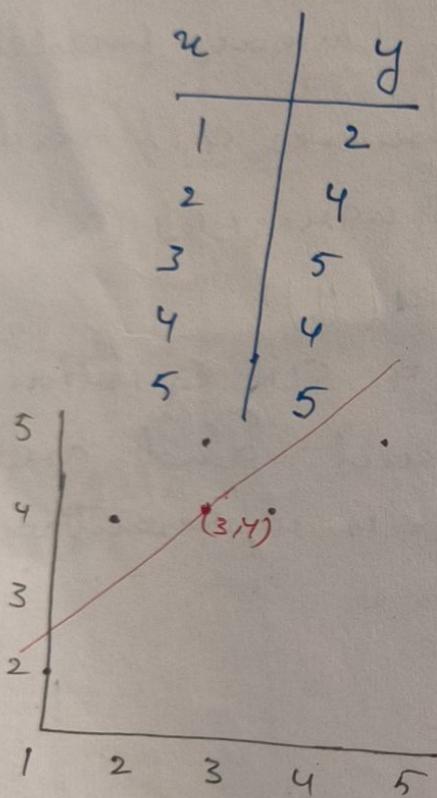
y - Dependant variable
 x - independent variable
 m → slope of the line
 c - coefficient of the line

For e.g.: let's plot the crop yield based on the amount of rainfall. Here, rainfall is the independent variable and crop yield is the dependent variable.



~~Eg:~~

Let's consider a sample data set with five rows and find out how to draw the regression line.



Now plot the x and y on the graph. Now we want to draw a regression line in ~~such a way~~ such a way that fits most of the data point in the given data set.

x	y	x^2
1	2	1
2	4	4
3	5	9
4	4	16
5	5	25

$$\Sigma = 15$$

$$\Sigma = 20$$

$$\Sigma = 55$$

$$\Sigma = 86$$

$$\Sigma = 66$$

y^2	$x \cdot y$
4	8
16	15
25	16
25	25



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Linear equation is represented by $y = mx + c$

$$m = \frac{((n * \Sigma(u * y)) - (\Sigma(n) * \Sigma(y)))}{((n * \Sigma(u^2)) - (\Sigma(u))^2)}$$

$$m = \frac{(5 * 66) - (15 * 20)}{(5 * 55) - (225)} = 0.6$$

$$c = \frac{((\Sigma(y) * \Sigma(u^2)) - (\Sigma(u) * \Sigma(u * y)))}{((n * \Sigma(u^2)) - (\Sigma(u))^2)}$$

$$c = 2.2$$

Line equation is $y = mx + c$

y_{pred} values are:

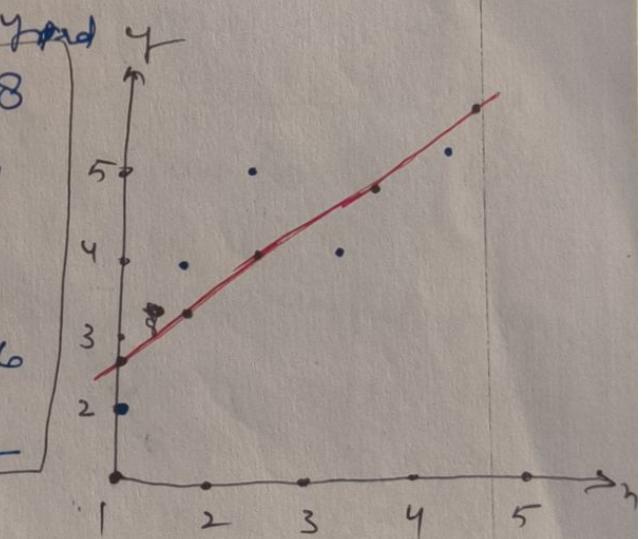
$$y = 0.6 * 1 + 2.2 = 2.8$$

$$y = 0.6 * 2 + 2.2 = 3.4$$

$$y = 0.6 * 3 + 2.2 = 4$$

$$y = 0.6 * 4 + 2.2 = 4.6$$

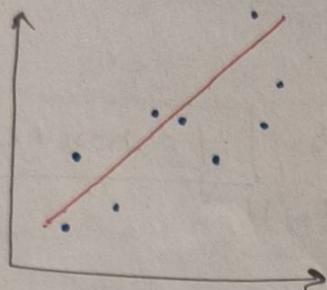
$$y = 0.6 * 5 + 2.2 = 5.2$$



Blue point represent actual value of y

Black point represent predicted value of y based on the model we created.

Distance between actual value and predicted values of y are known as errors or Residuals.
The best-fit line should have the lowest sum of squares of these errors, also known as 'e-square'.



We keep the line moving through the data points to make sure the best fit line has the least squared distance between the data points and the regression line.



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

Logistic Regression:

Logistic Regression is a mathematical model used in statistics to estimate the probability of an event occurring using some previous data.

Logistic Regression works with binary data, where the event either happens (1) or doesn't happen (0).

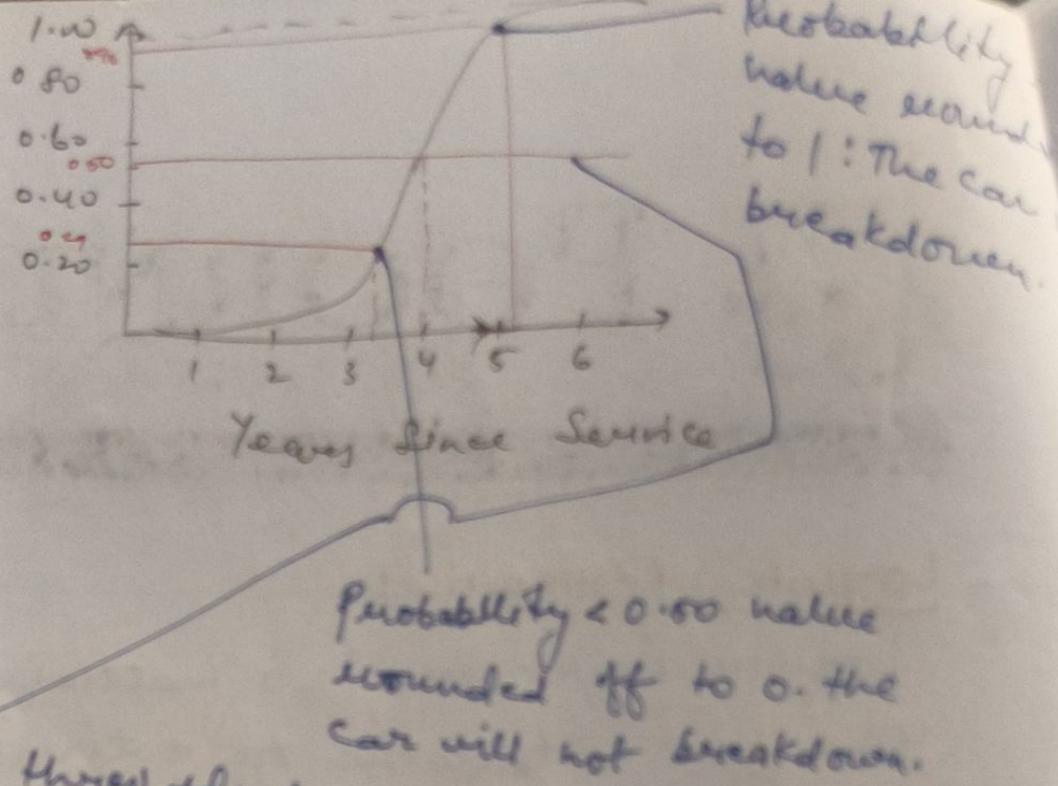
Eg: Employee got a promotion or not.

→ Logistic Regression is used to describe data and the relationship between one dependent variable and one or more independent variables. The independent variable can be nominal, ordinal or of interval type.

→ Logistic Regression name is derived from logistic function and this function is also called sigmoid function. Value of logistic function is lies between 0 and 1.

Eg:- The following is an example of a logistic function we can use to find the probability of a vehicle breaking down, depending on how many years it has been since it was service last.

Probability
of
breakdown



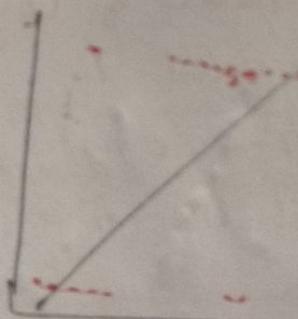
Here, the threshold value 0.50 indicates that the car is more likely to breakdown after 4 years of usage

→ This curve is called Sigmoid curve or S-curve, which mathematically represented like

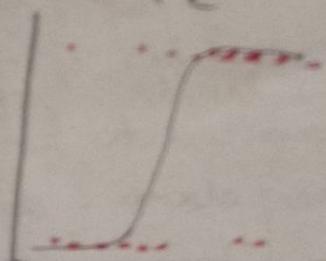
$$\sigma(y) = 1 / (1 + \exp(-y))$$

- The sigmoid function is also called squashing function as its domain, is the set of all real numbers and range is $(0, 1)$.
- If the input is either a very large negative number or a very large positive number, the σ_p is always between 0 and 1.

$$y = mx + b$$



$$y = \frac{1}{1 + e^{-(mx+b)}}$$





POORNIMA

COLLEGE OF ENGINEERING

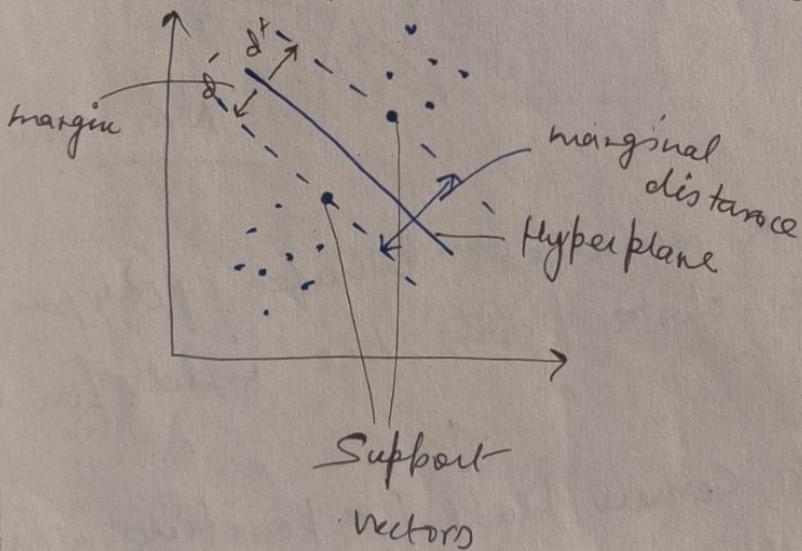
DETAILED LECTURE NOTES

PAGE NO.

Support Vector Machine :

Objective of SVM algorithm is to find a linear surface (hyper plane) in N dimensional space such that it distinctly classified the data points.

- A Linear surface is a line in 2D, plane in 3D, hyperplane in higher dimensions.
- Our objective is to find a higher plane such that it has maximum margin distance b/w the data points of the classes.



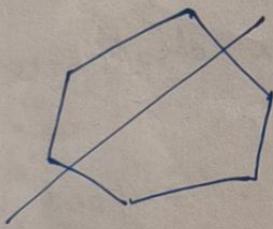
1. Support Vector
2. Hyperplane
3. Marginal Distance
4. Linear Separable
5. Non Linear Separable

- The higher the marginal distance, the more generalized our model.

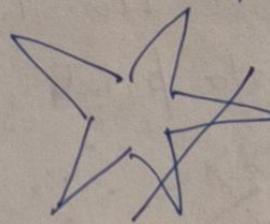
- The data points through which lines/ hyper
 (π^+) and (π^-) pass through are called as
Support Vector. Hence, the algorithm is called
as Support Vector Machine.
- Decision Boundary (Hyperplane)
- we always choose that hyper plane which is
having the largest margin (max width margin)
Maximal margin hyperplane should be selected
which decrease error rate and increase accuracy.

Geometric Intuition of SVM:-

Convex polygon: A shape is said to be a
convex polygon if a line drawn on that
shape intercept only in 2 points.



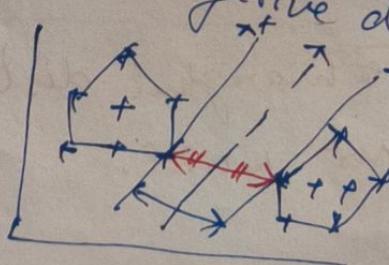
Convex.



Non-convex

Convex Hull: A smallest shape/polygon
where all the data points lie inside/on the
Polygon.

Step 1: Create a convex hull for positive data
points and negative data points.





POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

Step 2: Draw a shortest line that connects 1 data point of positive class and 1 data point of negative class (nearest point of each class).

Step 3: Draw a line/Hyperplane that bisects the shortest line drawn in Step 2. That line is Hyperplane.

PAGE NO.....

Equation of ~~Hyperline~~ Hyperplane:

Let

$$\pi: \omega^T n + b = 0$$

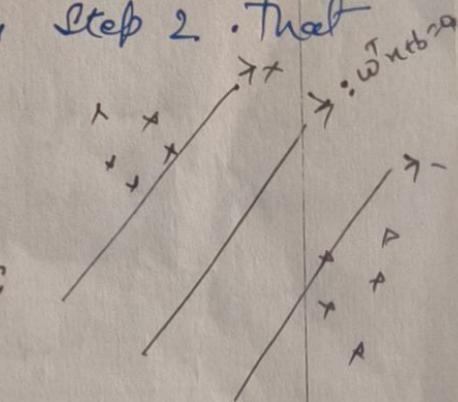
$$\pi^+: \omega^T n_1 + b = 1$$

$$\pi^-: \underline{\omega^T n_2 + b = -1}$$

$$\omega^T (n_2 - n_1) = 2$$

$$\frac{\omega^T}{\|\omega\|} (n_2 - n_1) = \boxed{\frac{2}{\|\omega\|}} \text{ optimization function}$$

$$\text{margin } (d) = \frac{2}{\|\omega\|}$$



$\omega^T \rightarrow$ Transpose of slope values

\rightarrow we can choose ω^T with norm of ω^T $\frac{2}{\|\omega\|}$

\rightarrow In ω^T , some direction involved

\rightarrow magnitude of ω^T off-

optimize
subject
but
condition
such that
 y_i

$$1 \quad \omega^T n_i + b \geq 1$$

$$-1 \quad \omega^T n_i + b \leq -1$$

$\Rightarrow y_i * \omega^T n_i + b_i \geq 1$ (for all n_i)
If y_i is not $>$ or $= 1$
then there is misclassification



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

67

Decision Tree in Machine Learning:

In decision tree each branch of the tree represents a possible action/decision occurrence or reaction. It is utilized for both regression and classification task.

Decision tree learning employs a divide and conquer strategy by conducting a greedy search to identify the optimal split points within a tree,

if else \Rightarrow Decision Tree

Important terms related to Decision Tree:

Entropy: Entropy is the measure of randomness or unpredictability in the dataset.

(Initially the dataset is having higher entropy).

\rightarrow Entropy is a metric used to train the decision tree. This metric measures the quality of split.

Information Entropy for a dataset with c classes are denoted as;

$$E = - \sum_{i=1}^c p_i \log p_i$$

~~log₂p_i~~ where p_i is the probability of randomly picking an element of class i .

Eg: Consider a dataset with 1 blue, 2 green, 3 red. Then

$$E = - (P_b \log_2 P_b + P_g \log_2 P_g + P_r \log_2 P_r)$$

$$P_b = \frac{1}{6}, P_g = \frac{2}{6}, P_r = \frac{3}{6}$$

$$E = - \left(\frac{1}{6} \log_2 \left(\frac{1}{6}\right) + \frac{2}{6} \log_2 \left(\frac{2}{6}\right) + \frac{3}{6} \log_2 \left(\frac{3}{6}\right) \right)$$

$$E = 1.46$$

Eg: Consider 3 blue only in a dataset

$$E = - \frac{3}{3} \log_2 \frac{3}{2,3}$$

$$= -1 \log_2 1$$

$$E = 0$$

Information Gain: how can we quantify the

IG = [How the features are selected] Quality of a Split?

~~Information gain is the measure of decrease in Entropy after the database is split.~~ Information gain is calculated for a split by subtracting weighted entropies of each branch from the original entropy.

→ Purity → Pure split

↳ Entropy

↳ Gini ~~coefficient~~ impurity

→ How the features are selected $\left[1 - \sum_{i=1}^n (P_i)^2\right]$

↳ Information Gain

Higher IG = more entropy removed

I^G is the measure of decrease in entropy after the dataset is split. (10)

[leaf node carries the classification labels that we are trying to predict]



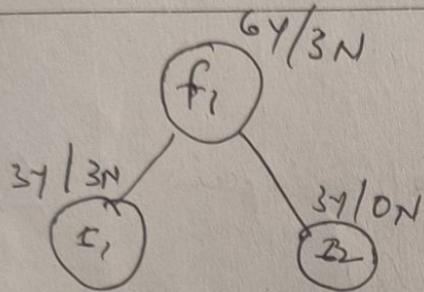
POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

The best-split is chosen by maximizing the I^G .

PAGE NO.

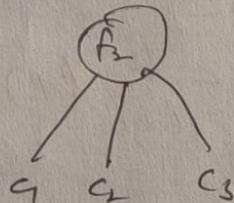
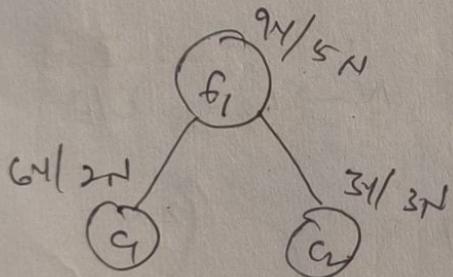


If Entropy \approx then
Please Split

Entropy ≈ 1 then
Impure split

Further division

→ which feature to take to split??



Information Gain :

$$\text{Gain}(S, f_i) = H(S) - \sum_{v \in V_i} \frac{|S_v|}{|S|} H(S_v)$$

on particular node total sample



$$H(S_{v_1}) \rightarrow c_1$$

$$H(S_{v_2}) \rightarrow c_2$$

$H(S)$ is Entropy of root node.

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

$$H(S) = 0.94$$

$$H(S_{VC_1}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$H(S_{VC_1}) = 0.81$$

$$H(S_{VC_2}) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$H(S_{VC_3}) = 1$$

$$\text{Gain}(S, f_1) = 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$= \underline{\underline{0.049}}$$

\Rightarrow Suppose for feature 2 information gain is $\underline{\underline{0.051}}$

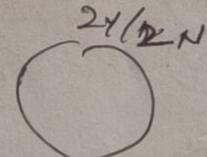
$$\text{Gain}(S, f_2) >> \text{Gain}(S, f_1)$$

\downarrow use feature 2, we will start the split.
Highest is ~~is~~ value feature is used for splitting

\rightarrow Gini Impurity

$$GI = 1 - \sum_{i=1}^N (P_i)^2 \quad N \rightarrow \text{no. of O/P} \quad \begin{cases} \text{Yes} \\ \text{No} \end{cases}$$

$$= 1 - [(P_+)^2 + (P_-)^2]$$



$$= 1 - \left[\left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2 \right]$$

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

$$= 0.5$$

\rightarrow It means complete impure split.

\rightarrow Gini calculates fast & in comparison to Entropy, Entropy use log function and Gini use simple math.



Poornima

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

Naïve Bayes: Naïve Bayes works on the principle of conditional probability as given by the Bayes theorem.

Bayes theorem gives the conditional probability of an event A given the another event B has occurred.

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

where :

$P(A)$ = The probability of A occurring

$P(B)$ = The probability of B occurring

$P(A, B)$ = The probability of A given B

$P(B, A)$ = The probability of B given A

Eg: Tossing two coins :

Sample space = { HH, HT, TH, TT }

$P(\text{first coin is tail when the second coin is head}) = ?$

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

$$P(A|B) = \frac{P(B|A) * \frac{1}{2}}{\frac{1}{2}}$$

$$= \frac{\frac{1}{4} * \frac{1}{2}}{\frac{1}{2}} = 0.25$$

Application \rightarrow Sentiment Analysis

Face Recognition

Weather Prediction

Medical Diagnosis

$$P(Y | x_1, x_2, \dots, x_n) = P(x_1|y) * P(x_2|y) * P(x_3|y) * \dots * P(x_n|y)$$

Sample (N/..)

Person	COVID (Yes/No)	Fever	
		P(x ₁)	P(x ₂)
1	Yes	Yes	No
2	No	Yes	Yes
3	Yes	Yes	Yes
4	No	No	Yes
5	Yes	No	No
6	No	No	Yes
7	Yes	No	Yes
8	Yes	Yes	No
9	No	Yes	Yes
10	No	Yes	No

Given Person (flu, COVID)

$$P(\text{Yes} | \text{flu, COVID}) = P(\text{flu/yes}) * P(\text{COVID/yes}) * P(\text{yes})$$

$$P(\text{No} | \text{flu, COVID}) = P(\text{flu/no}) * P(\text{COVID/no}) * P(\text{no})$$

Sol:

Step 1: Prior probability:

$$P(\text{fever} = \text{yes}) = 7/10$$

$$P(\text{fever} = \text{no}) = 3/10$$

Step 2: Conditional probability:

	Yes	No	Fever
Covid	4/7	2/3	
Flu	3/7	2/3	



POORNIMA

COLLEGE OF ENGINEERING

DETAILED LECTURE NOTES

PAGE NO.

(Y)

$$= \frac{9}{7} \times \frac{3}{7} \times \frac{7}{10}$$

$$= \frac{12}{17}$$

$$= 0.17$$

(N)

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

$$= \frac{4}{30}$$

$$= 0.13$$

Probability value of high for Yes.

So for given problem, value is Yes

fewer



K nearest neighbors (KNN) :- It uses the past data for learning and produces prediction.

KNN classifies a data point based on how its neighbors are classified. Stores all available cases that are being provided while training and then classifies new cases based on a similarity measures.

- The process of choosing K value is termed as parameter tuning and it is very important aspect that is going to determine the accuracy of model.
- ↳ Value of K can be chosen by applying following methods :
 - Hit and trial
 - \sqrt{n} , where n stands for total number of data samples in dataset.
 - Odd value of K is selected to avoid confusion between two classes of data.
- ↳ When we choose KNN then
 - Data-set should be properly labelled.
 - ~~can't~~ KNN cannot be used for huge data. It is lazy learner.
 - KNN provides higher accuracy for small datasets.

How does KNN works?

Euclidean distance of two point (A, B)

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

Eg:-

Weight	Height	Body type	<u>ED</u>
62	182	Obese	13
55	170	Obese	2
58	174	Non-Obese	4.1
57	173	Obese	3
65	172	Obese	8.2
58	169	Obese	1.4
51	167	Obese	6.7
69	176	Non-Obese	13.4
64	173	Obese	7.6

Testing Sample:-

57	170	?
----	-----	---

label

~~Euclidean Distance~~

Higher votes are for obese.

KNN Classifier in Machine Learning: KNN is supervised machine learning algorithm that can be used to solve the classification and regression problems.

-The k-nearest neighbors (KNN) algorithm is a simple, easy-to-implement, non-parametric, lazy learning, supervised machine learning algorithm that can be used to solve both classification and regression problems using feature similarity.

K- Nearest Neighbors is a

- Non parametric as it does not make an assumption about the underlying data distribution pattern

- Lazy algorithm as KNN does not have a training step. All data points will be used only at the time of prediction. With no training step, prediction step is costly.
- Supervised machine learning algorithm as target variable is known
- Used for both Classification and Regression
- Uses feature similarity/nearest neighbors to predict the cluster that the new point will fall into.

$1-KNN$ algorithm assumes that similar things exist in close proximity (means near to each other). Therefore, KNN captures the idea of similarity (distance, proximity or closeness) with some mathematics

(i.e. calculating the distance of a new data point with nearest point).

- KNN classified a data point based on how its neighbor's are classified.

• KNN algorithms steps:

1. Compute a distance value between the item to be predicted and every item in the training dataset.
2. Pick the k closest data points (the items with the k lowest distances)
3. Get the labels/values of the selected k neighbors

{ If regression, return the mean of the ' k ' neighbors
If classification, return the mode of the ' k ' neighbors.

Decision Tree (Numerical Example)

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Y
/
attribute
target

Step: 1 Find the overall entropy for the dataset denoted by S and then find the entropy for each attribute value for each and specific features.

Values (Outlook) = Sunny, Overcast, Rain

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

$$\text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

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

$$\text{Entropy}(S_{\text{Sunny}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$S_{\text{Overcast}} \leftarrow [4+, 0-]$$

$$\text{Entropy}(S_{\text{Overcast}}) = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = 0$$

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

$$\text{Entropy}(S_{\text{Rain}}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$$

$$\text{Gain}(S, \text{Outlook}) = \text{Entropy}(S) - \sum_{v \in \{\text{Sunny}, \text{Overcast}, \text{Rain}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Gain(S, Outlook)

$$= \text{Entropy}(S) - \frac{5}{14} \text{Entropy}(S_{\text{Sunny}}) - \frac{4}{14} \text{Entropy}(S_{\text{Overcast}})$$

$$- \frac{5}{14} \text{Entropy}(S_{\text{Rain}})$$

$$\text{Gain}(S, \text{Outlook}) = 0.94 - \frac{5}{14} 0.971 - \frac{4}{14} 0 - \frac{5}{14} 0.971 = 0.2464$$

Attribute (Temp)

Values (Temp) = Hot, Mild, Cool

$$S = [9+, 5-] \quad \text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

$$S_{Hot} \leftarrow [2+, 2-] \quad \text{Entropy}(S_{Hot}) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1.0$$

$$S_{Mild} \leftarrow [4+, 2-] \quad \text{Entropy}(S_{Mild}) = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.9183$$

$$S_{Cool} \leftarrow [3+, 1-] \quad \text{Entropy}(S_{Cool}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.8113$$

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

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

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

$$- \frac{4}{14} \text{Entropy}(S_{Cool})$$

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

$$= 0.289$$

$$\Rightarrow \underline{0.04}$$

Attribute (Humidity)

Values (Humidity) = High, Normal

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

$$\text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

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

$$\text{Entropy}(S_{High}) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.9852$$

$$S_{Normal} \leftarrow [6+, 1-]$$

$$\text{Entropy}(S_{Normal}) = -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.5916$$

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

Gain(S, Humidity)

$$= \text{Entropy}(S) - \frac{7}{14} \text{Entropy}(S_{High}) - \frac{7}{14} \text{Entropy}(S_{Normal})$$

$$\text{Gain}(S, \text{Humidity}) = 0.94 - \frac{7}{14} 0.9852 - \frac{7}{14} 0.5916 = 0.1516$$

Attribute (Wind)

Values (Wind) = Strong, Weak

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

$$\text{Entropy}(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

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

$$\text{Entropy}(S_{Strong}) = 1.0$$

0.8113

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

$$\text{Entropy}(S_{Weak}) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} = 0.3113$$

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

$$\text{Gain}(S, Wind) = \text{Entropy}(S) - \frac{6}{14} \text{Entropy}(S_{Strong}) - \frac{8}{14} \text{Entropy}(S_{Weak})$$

$$\text{Gain}(S, Wind) = 0.94 - \frac{6}{14} \cdot 1.0 - \frac{8}{14} \cdot 0.3113 = 0.6313$$

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

$$\text{Gain}(S, Wind) = \text{Entropy}(S) - \frac{6}{14} \text{Entropy}(S_{Strong}) - \frac{8}{14} \text{Entropy}(S_{Weak})$$

$$\text{Gain}(S, Wind) = 0.94 - \frac{6}{14} \cdot 1.0 - \frac{8}{14} \cdot 0.8113 = 0.0478$$

0.0478

Now Information gain for all attributes is:

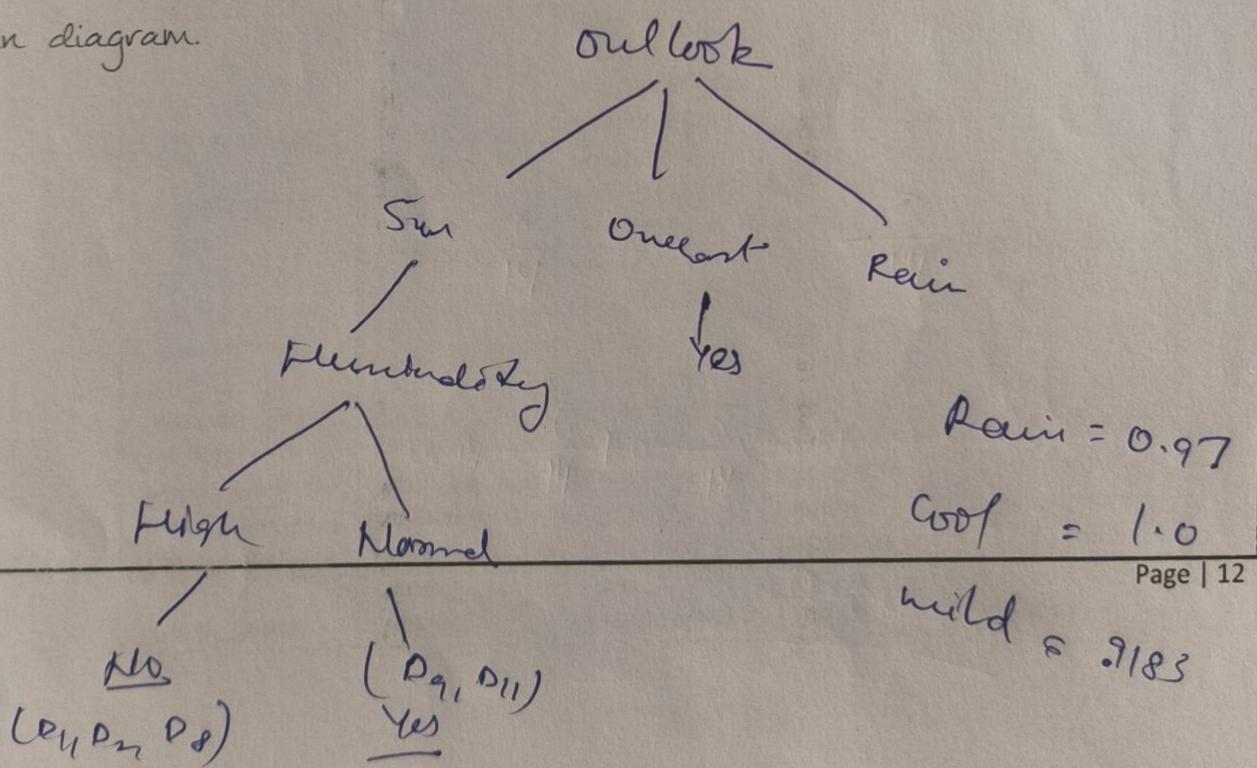
Day	Temp	Humidity	Wind	Play Tennis
D4	Mild	High	Weak	Yes
D5	Cool	Normal	Weak	Yes
D6	Cool	Normal	Strong	No
D10	Mild	Normal	Weak	Yes
D14	Mild	High	Strong	No

$$Gain(S_{Rain}, Temp) = 0.0192$$

$$Gain(S_{Rain}, Humidity) = 0.0192$$

$$Gain(S_{Rain}, Wind) = 0.97$$

Now choose wind as next attribute since it is having the highest gain. So updated and final decision tree look like this in the given diagram.



Now Information gain for all attributes is:

$$Gain(S, \text{Outlook}) = 0.2464$$

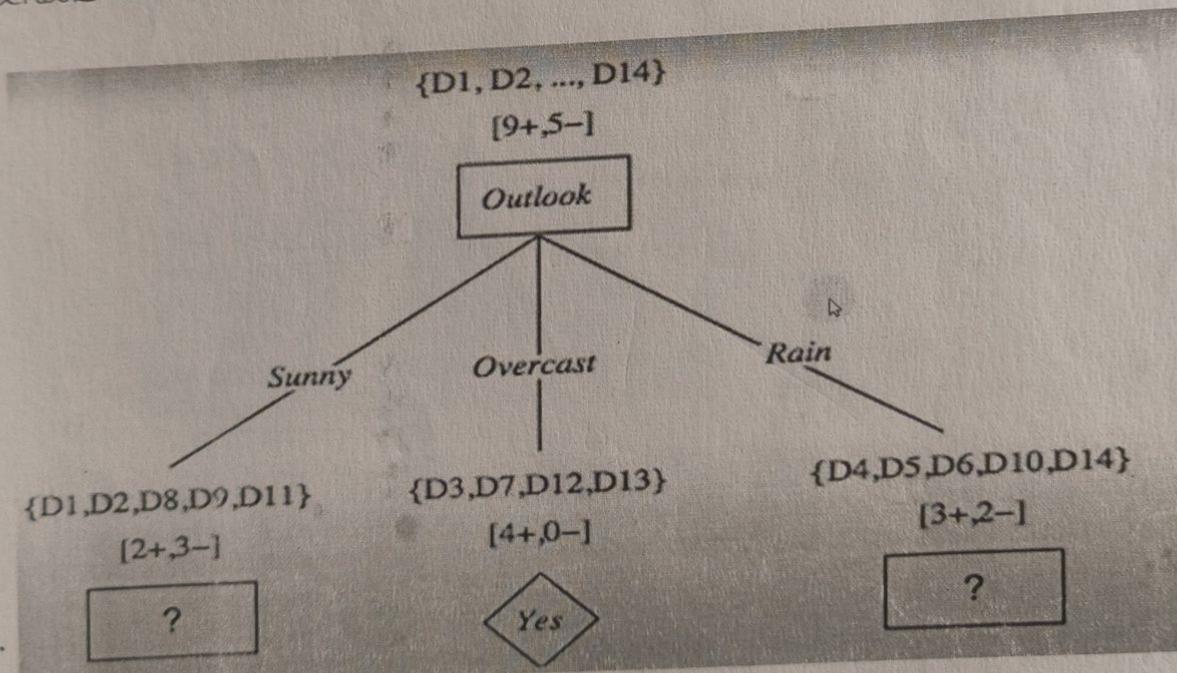


$$Gain(S, \text{Temp}) = 0.0289$$

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

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

So we choose Outlook variable as an root node. After choosing Outlook as an root node then the Decision tree look like this:



Ten, humidity

$$\text{Ten} = 0.97$$

~~Temp, humidity~~

$$\text{Temp} = 0.97$$

wind
0.97

Page | 6

$$Heat = 0$$

$$\text{High} = 0$$

$$\text{Strong} = 1.0$$

$$Wind = 1$$

$$\text{Normal} = 0$$

$$\text{Weak} = 0.97$$

$$Cool = 0.0$$

?

$$= 0.570$$

$$= 0.97$$

$$= 0.0142$$



Random Forest Algorithm:

Random forest is constructed using multiple decision trees and final decision is obtained by the majority votes of the decision trees.

→ If the problem is of classification type, decision tree algorithm takes the majority vote.

→ If the problem is of regression type, decision tree algorithm takes the means of decision.

Random forest or random decision forest is an ensemble learning method for classification, regression.
 ✓ [By this overfitting problem will be resolved].
 ✓ In classification gives more accurate result in comparison to regression.

for ex:-

Mango Type		Source/Sweetness (on scale of 5)	Diameter (cm)
A	M	4.5	7
A	M	4.5	6.8
D	UP	3.8	9
D	UP	3.6	7.9
K	M	2.7	6
K	M	4	6.5

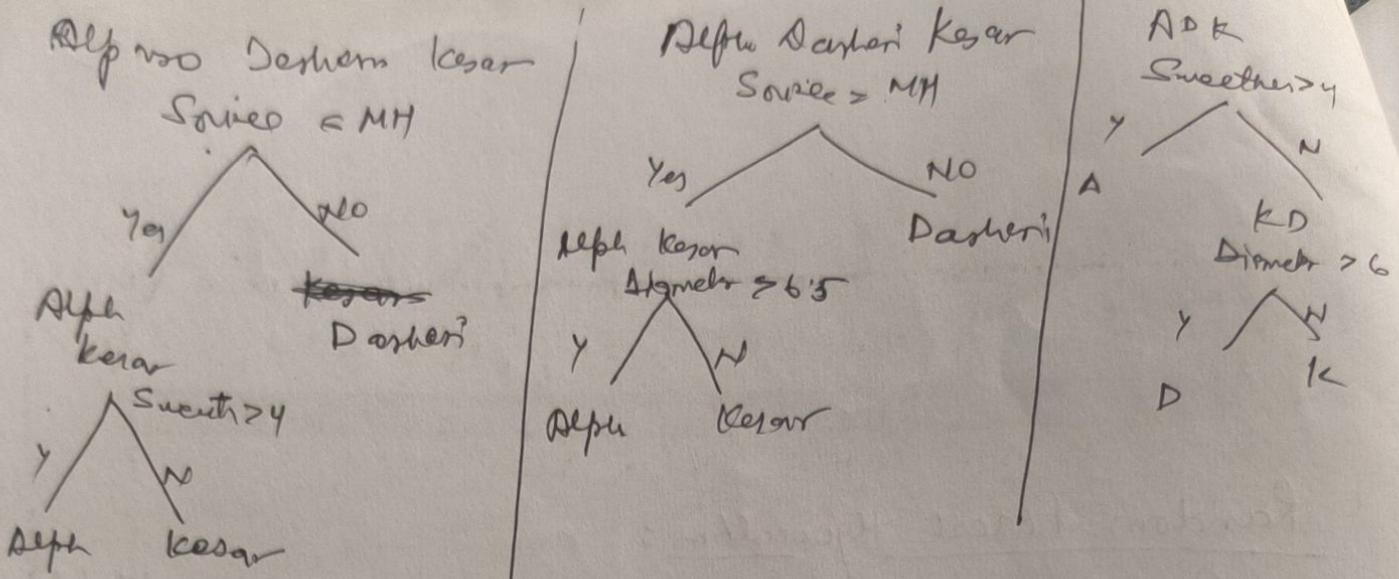
MP	Source	Sweetness		
			Mahan	Up
A	Mahan	4.5		
D	Up	3.8		
K	Mahan	3.7		
K	Up	4		

MP	Source	Sweetness		
			Mahan	Up
A	Mahan	4.5		
A	Up	4.5		
K	Mahan	3.7		
D	Up	3.8		

- The process of selecting features for bootstrap sample is known as Feature Selection.

- The process of creating trees for subsets is known as Random Sampling.

MP	Source	Sweetness		
			Mahan	Up
A	Mahan	4.5		
A	Up	4.5		
K	Mahan	3.7		
D	Up	3.8		



Feature Selection in Random Forest

- For classification, by default the feature selection is taken as:
= sqrt number of features in data
- For regression, by default the feature selection is taken as:
= total number of features in Dataset

Ensembling means Aggregating the results of decision trees and taking the majority votes in classification, and mean in case of regression.

New Sample :-

Sourie	Sweetness	Diameter	Class
Makorchi	3.9	6.4	Class of Mayo = ?
①	②	③	
Kesar	Kesar	D	Class is Kesar

Advantages

- low variance
- Reduced overfitting
- No need for Normalization
- High accuracy
- High Scalability

Variability
Handles noisy data

Random Forest Algorithm: A random forest is actually just a bunch of decision trees bundled together (that's why it is called as forest).

-- Random forest is a Supervised Machine Learning Algorithm that is used widely in Classification and Regression problems. It builds decision trees on different samples and takes their majority vote for classification and average in case of regression.

Consider the following dataset-

TRAINING DATASET

COLOR	DIAMETER	LABEL
RED	3	APPLE
YELLOW	3	LEMON
PURPLE	1	GRAPES
RED	3	APPLE
YELLOW	3	LEMON
PURPLE	1	GRAPES

CONDITIONS

COLOR == PURPLE?

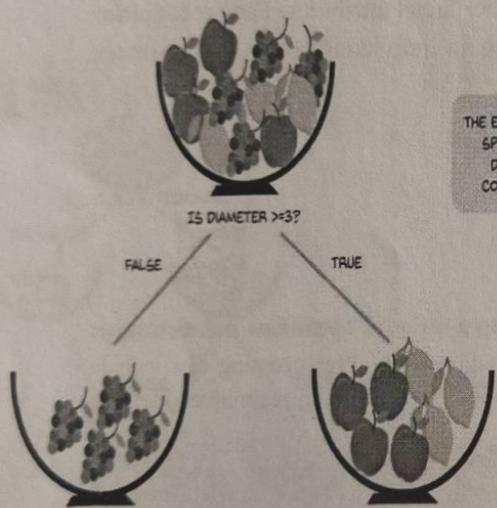
DIAMETER = 3

COLOR == YELLOW?

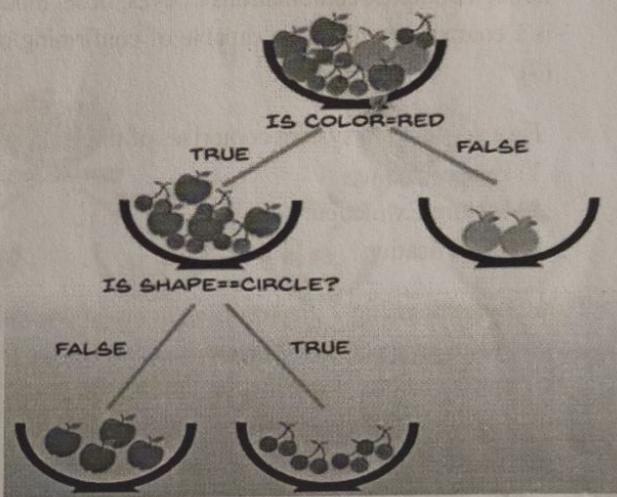
COLOR == RED?

DIAMETER = 1

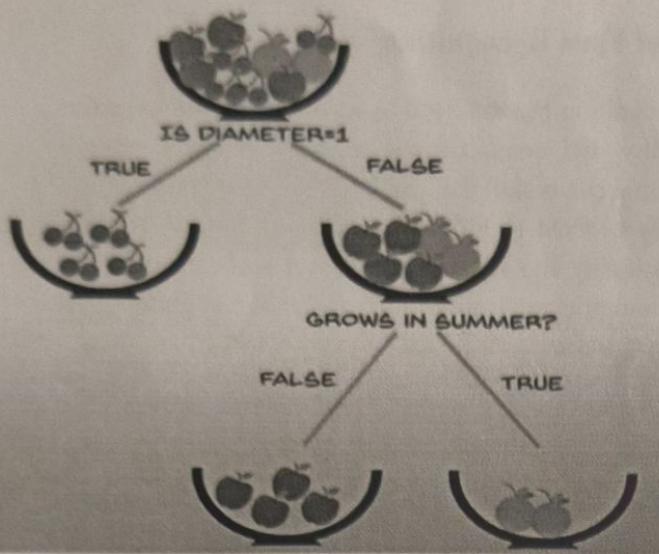
Decision Tree-1



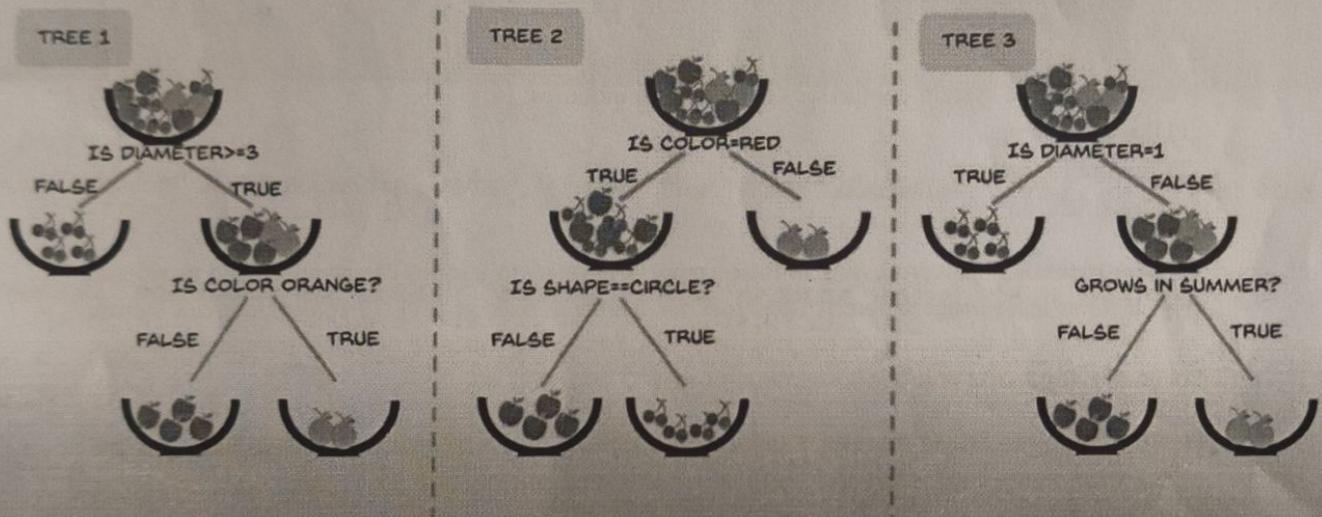
Decision Tree-2



Decision Tree-3



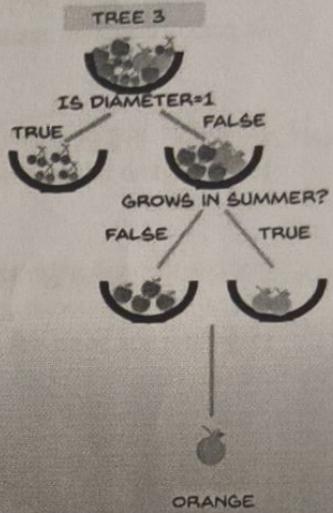
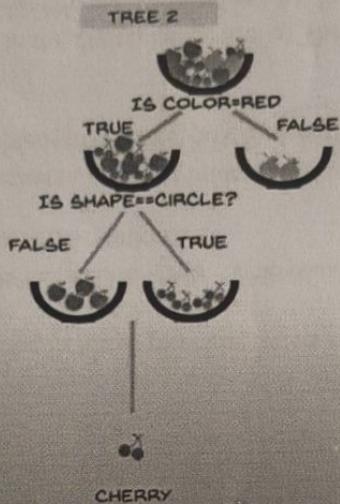
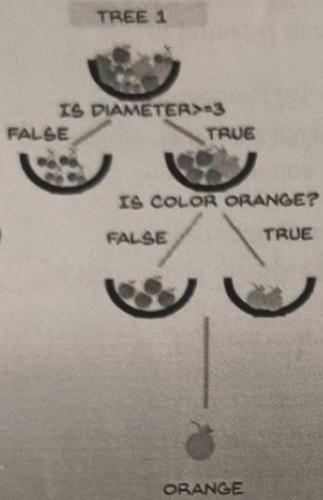
Now all decision trees are as-



Let's take a new fruit and find the name of the fruit.



DIAMETER = 3
COLOUR = ORANGE
GROWS IN SUMMER = YES
SHAPE = CIRCLE



Now take a majority vote, since the two-decision tree predict the orange, so random forest predict the orange label to classify the new fruit.

-- So a random forest or Random decision forest is a method that operates by constructing multiple decision tree during training phase.

--The decision of the majority of the trees is chosen by random forest as the final decision.

Random Forest Algorithm Steps:

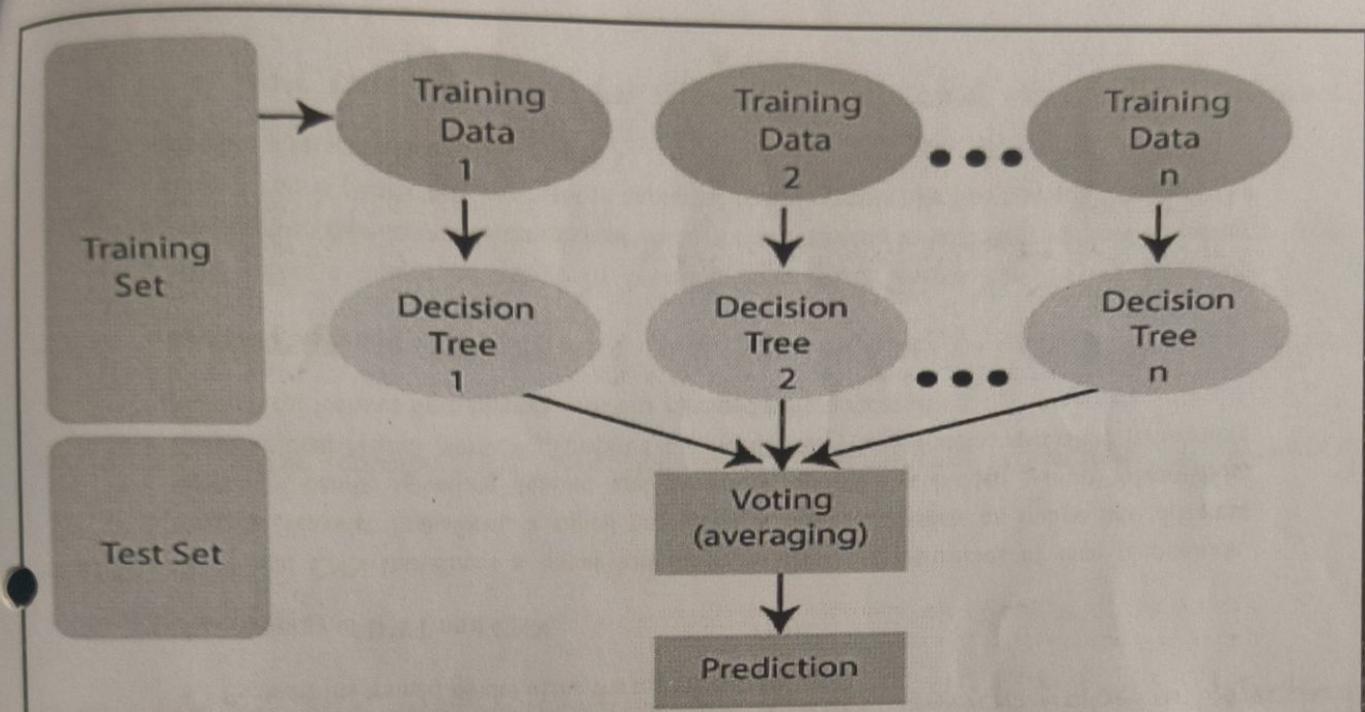
The following steps explain the working Random Forest Algorithm:

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: Finally, select the most voted prediction result as the final prediction result.



Advantages:

- Random forests is considered as a highly accurate and robust method because of number of decision trees participating in the process.
- It does not suffer from the overfitting problem. The main reason is that it takes the average of all the predictions, which cancels out the biases.

- The algorithm can be used for classification and regression problems.
- Random forests can also handle missing values. There are two ways to handle these. Using median values to replace continuous variables, and computing the proximity-weighted average of missing values.
- you can get the relative feature importance, which helps in selecting the most contributing features for the classifier.

Disadvantages:

- Random forests is slow in generating predictions because it has multiple decision trees. Whenever it makes a prediction, all the trees in the forest have to make a prediction for the same given input and then performing voting on it. This whole process is time consuming.

- The model is difficult to interpret as compared to a decision tree, where you can easily make a decision by following the path in the tree.
- Requires much computational power as well as resources as it builds numerous trees to combine their outputs.
- Requires much training time as it combines a lot of decision trees to determine the class.