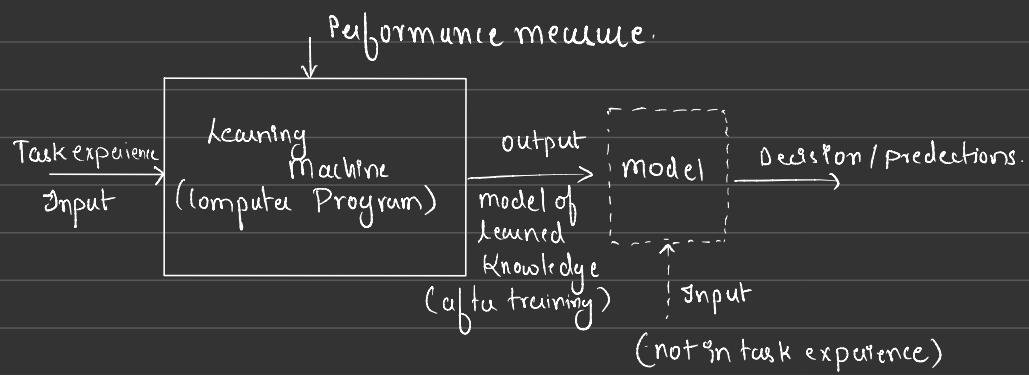


•> Metrics for Assessing Regression Accuracy.

⇒



•> Mean Square Error

→ The mean is obtained from the training data as arithmetic average, which measures the average square deviation of the predicted values from the true value.

$$MSE = \frac{1}{N} \sum_{i=1}^N (y^{(i)} - h(\omega, u^{(i)}))^2$$

•> Root mean Square error

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (y^{(i)} - h(\omega, u^{(i)}))^2}$$

•> Sum of Error Squares

⇒ Mathematical manipulation

$$= \sum_{i=1}^N (y^{(i)} - h(\omega, u^{(i)}))^2$$

•> Expected Mean square error

$$EMSE = E \left[\sum_{i=1}^N (y^{(i)} - \underbrace{h(\omega, u^{(i)})}_\text{hypothesis funct.})^2 \right]$$

statistical expectation operation

•> Mean absolute Error

$$MAE \Rightarrow \frac{1}{N} \sum_{i=1}^N |y^{(i)} - h(w, u^{(i)})|$$

•> Classification Accuracy

i) Miss classification Error

\Rightarrow Formula = $\frac{\text{Number of data points for which } (y^{(i)} - \hat{y}^{(i)}) \neq 0}{\text{Number of samples } (N)}$

•> Confusion Matrix :

		Predicted values	
		No	Yes
Actual values	No	TN	FP
	Yes	FN	TP

TN \rightarrow True negative

FN \rightarrow False negative

TP \rightarrow True positive

FP \rightarrow False negative

\Rightarrow The true +ve and true -ve are the accurate classifications.

\Rightarrow A false +ve takes place when the result is inaccurately predicted as +ve but its -ve in reality

\Rightarrow A false -ve is said to offer when the result is inaccurately predicted as -ve but in reality it is +ve

\Rightarrow An ideal algorithm should create a diagonal matrix which $FP = FN = 0$

→ Metrics for calculation of confusion matrix

1) Miss classification error or (Success Rate)

→ It has two parts

↓
Success Rate

$$\frac{TP + TN}{TP + TN + FP + FN}$$

↓
Miss classification Rate

$$\frac{FP + FN}{FP + FN + TP + TN}$$

True Positive Rate

$$\Rightarrow \frac{TP}{TP + FN}$$

True Negative Rate

$$\frac{TN}{FP + TN}$$

True Positive Rate tells how sensitive our decision method is in detecting the abnormal event

True Negative Rate tells how specific our decision method is in detecting the abnormal event

So a decision technique is said to be good if it is highly sensitive and highly specific.

(Q)

PV

		PV	
		Yes	No
AV	No	50	10
	Yes	5	100
		55	110

Calculate the success rate, failure rate

$$\text{Success} = 0.9091$$

$$\text{Failure} = 0.09$$

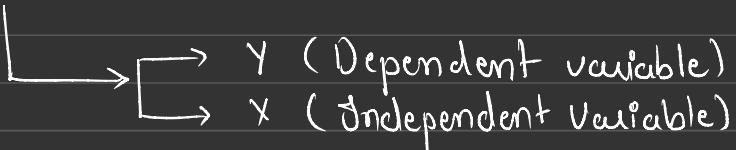
$$\text{TP rate} = 0.95$$

$$\text{TN rate} = 0.83$$

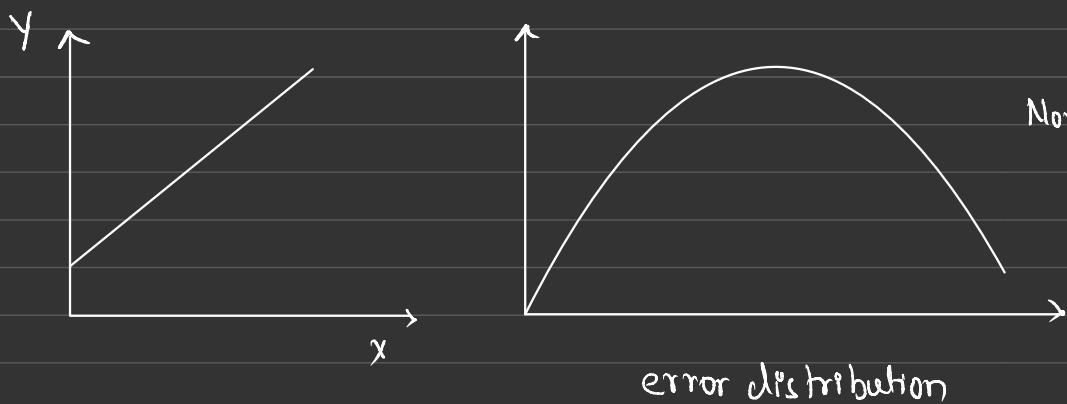
Unit - 7 : Generalized Linear Model (GLM)

- 1) Why GLM
- 2) Assumptions Made
- 3) Component of GLM

-) Linear Models :



→ When relation between Y & X are linear we will use linear model



Definition: Generalised Linear model allow to build a linear relationship between the response and the predictors. Even though their underlying relationship is not linear. This is made possible by means of link function or activation function.

→ Here the errors in the response variable are assumed to follow an exponential family of distribution.

-) Assumptions of GLM

→ There are 2-3 assumptions that we will make.

- 1) Data should be independent and random
- 2) Y or response does not need to be normally distributed but distribution is from exponential family (Binomial, Beta, and gamma distributions)
- 3) The original response variable need not have a linear relationship with X, but the transformed response variable should be linearly dependent on X.

•) Equation of Logistic Regression:

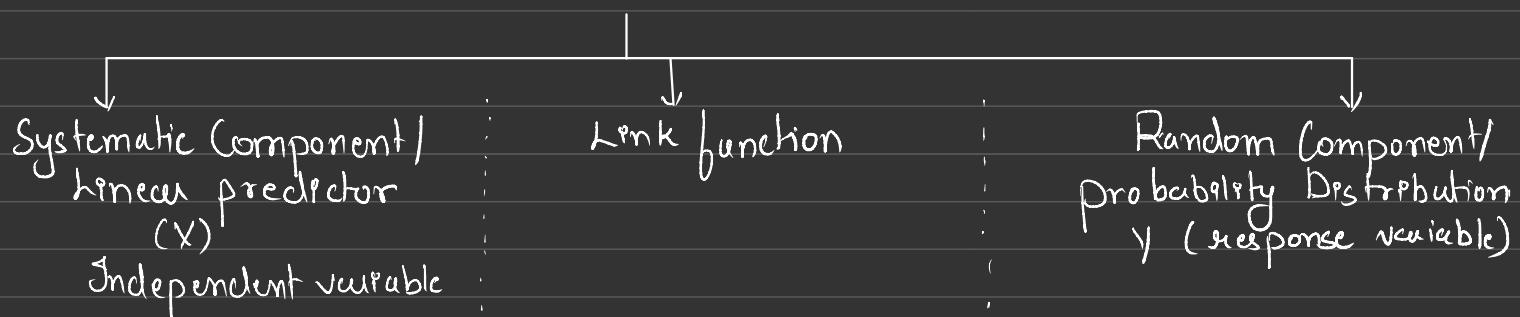
$$\text{log odds} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$$

Comparison between linear regression.

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots \quad (\text{only LHS part is changed})$$

•) Components of GLM

→ There are 3 components of GLM.



→ It is the linear combination of predictors and the regression parameters (β_1, β_2, \dots)

→ It specifies the link between random component & systematic component. It indicates how y relates to the linear combination of x .

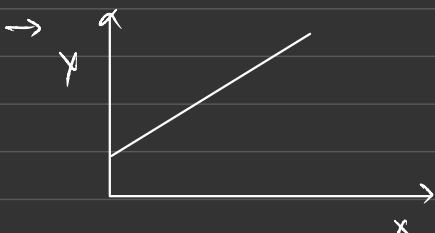
→ It refers to the probability distribution from the exponential family of distributions of y .

•) Response Variable Distribution Support Link name Link function

Only two possible	Bernoulli	Integers: {0, 1}	Logit	$X_B = \ln(u/(1-u))$
-------------------	-----------	------------------	-------	----------------------

Linear Regression

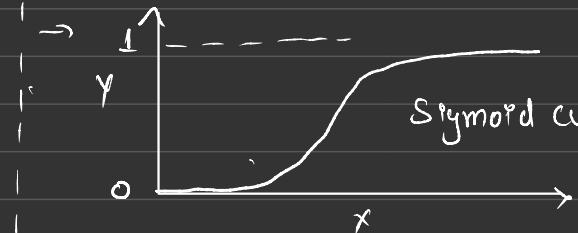
- Response : Continuous
- y & x linearly related



- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$
- Least Square

Logistic Regression

- Response : binary / categorical
- y & x relation is not linear.



- $\text{log odds} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$
- maximum likelihood estimation.

→ The odds: ratio of the proportions for the two possible outcomes is a number between 0 and infinity

• What is logistic function

$$\text{log(odds)} = \text{log} \frac{P(\text{class 1} | x)}{1 - P(\text{class 1} | x)} = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots$$

$$P(\text{class 1} | x) = \left(\frac{1}{1 + e^{-(\omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n)}} \right)$$

$$= \left(\frac{1}{1 + e^{-a}} \right) \rightarrow \text{sigmoid function}$$

• Derivation / Parameters of logistic Regression.

⇒ We are given a set of observed data

$$\begin{aligned} \overset{\text{Dataset}}{D} &= \{x^{(i)}, y^{(i)}\}_{i=1}^N \\ x &= [x_1, x_2, \dots, x_n]^T \in \mathbb{R}^n \\ y &\in \{0, 1\} \end{aligned}$$

∴ The proposed model

$$P(y=1|x) = \frac{1}{1 + e^{-(\omega_0 + \omega_1 x_1 + \omega_2 x_2 + \dots + \omega_n x_n)}}$$

$$= \frac{1}{1 + e^{-(\omega^T x + \omega_0)}} \quad \text{where } [\omega^T = \omega_0, \omega_1, \dots]$$

$$P(y=0|x) = 1 - P(y=1|x)$$

⇒ Exponential family of distribution

→ Bernoulli

→ Sigmoid function $\frac{1}{1 + e^{-a}}$

•> Maximum Likelihood Criteria / Estimation

→ It is used for estimating the parameters of logistic regression model.

→ The parameter values are found such that they maximize the probability of observing the dataset of a given model.

⇒

$$\mathcal{L}(\{\omega, \omega_0\}, D) = \mathcal{L}(\bar{\omega}, D) = \prod_{i=1}^N P(y^{(i)} | x^{(i)})$$

⇒ In the maximum likelihood problem, our goal is to find ω which will maximizes \mathcal{L} , i.e., we wish to find $\bar{\omega}^*$ where

$$\bar{\omega}^* = \arg \max_{\bar{\omega}} \mathcal{L}(\bar{\omega}, D)$$

⇒ Often we maximize $\log(\mathcal{L}(\omega, D))$ instead because it is analytically easier.

$$\begin{aligned} \log(\mathcal{L}(\bar{\omega}, D)) &= \log \left(\prod_{i=1}^N P(y^{(i)} | x^{(i)}) \right) \\ &= \sum_{i=1}^N \log P(y^{(i)} | x^{(i)}) \end{aligned}$$

We model log-likelihood as

$$P(Y=y) = P^y (1-P)^{1-y}; \quad y \in [0, 1]$$

$$\begin{aligned} \log(\mathcal{L}(\bar{\omega}, D)) &= \log \left(\prod_{i=1}^N P(y^{(i)} | x^{(i)})^{y^{(i)}} (1-P(y^{(i)} | x^{(i)}))^{1-y^{(i)}} \right) \\ &= \sum_{i=1}^N \log(g(x^{(i)}, \bar{\omega})^{y^{(i)}} (1-g(x^{(i)}, \bar{\omega}))^{1-y^{(i)}}) \end{aligned}$$

To maximize log-likelihood with respect to $\bar{\omega}$,

$$\begin{aligned} \log(\lambda_i(\bar{\omega}, D)) &= \log(g(x^{(i)}, \bar{\omega})^{y^{(i)}} (1-g(x^{(i)}, \bar{\omega}))^{1-y^{(i)}}) \\ &= y^{(i)} \log g(x^{(i)}, \bar{\omega}) + (1-y^{(i)}) \log(1-g(x^{(i)}, \bar{\omega})) \end{aligned}$$

Taking derivative from above eqn

$$\nabla_{\bar{\omega}} (\log \lambda_i(\bar{\omega})) = \frac{y^{(i)}}{g(x^{(i)}, \bar{\omega})} \nabla_{\bar{\omega}} (g(x^{(i)}, \bar{\omega})) - \frac{1-y^{(i)}}{1-g(x^{(i)}, \bar{\omega})} \nabla_{\bar{\omega}} (g(x^{(i)}, \bar{\omega}))$$

Now,

$$g(x^{(i)}, \bar{\omega}) = \frac{1}{1 + e^{-\alpha}} ; \quad \alpha = (\bar{\omega}^T \bar{x}^{(i)})$$

Therefore

$$\begin{aligned}\nabla_{\bar{\omega}} (g(x^{(i)}, \bar{\omega})) &= \frac{\partial}{\partial \alpha} \left(\frac{1}{1 + e^{-\alpha}} \right) \frac{\partial \alpha}{\partial \bar{\omega}} \\ &= \left(\frac{1}{1 + e^{-\alpha}} \right) \left(1 - \frac{1}{1 + e^{-\alpha}} \right) \bar{x}^{(i)} \\ &= g(x^{(i)}, \bar{\omega})(1 - g(x^{(i)}, \bar{\omega})) \bar{x}^{(i)}\end{aligned}$$

This gives

$$\nabla_{\bar{\omega}} (\log(\lambda_i(\bar{\omega}))) = (y^{(i)} - g(x^{(i)}, \bar{\omega})) x^{(i)}$$

∴ The gradient ascent function will be as follows.

$$\begin{aligned}\Rightarrow \omega_{\text{new}} &= \omega_{\text{old}} + \eta \frac{\partial (\text{likelihood function})}{\partial \omega} \\ &= \omega_{\text{old}} + \eta \sum_{i=1}^n (y^{(i)} - g(x^{(i)}, \bar{\omega})) x^{(i)}\end{aligned}$$

•> Difference Between linear model & Generalized Linear Model

Generalized Linear Model : Linear Model

1) Logistic Regression | 1) Linear Regression

2) Categorical output | 2) Output is Continuous.

3) Non Normal distribution | 3) Normal Distribution.

•> Logistic Regression.

- It is a Supervised machine learning algorithm that can be used to model the probability of a certain class or event.
- It is a non linear regression problem.
- > Difference Between Probability & Likelihood.
- Probability follows clear parameters and computations while a likelihood is based mainly on observed factors.

$$L(\mu, \sigma, \text{data}) = P(\text{data} ; \mu, \sigma)$$

$\downarrow \quad \downarrow$
mean standard deviation.

Probability	Likelihood
$\mu = \text{Mean}$ $\sigma = \text{S.D}$ } constant	μ } may vary σ

- $P(\text{data}; \mu, \sigma)$: It means "The probability density of observing the data with model parameters μ & σ ". We can generalize this to any number of parameters and any distribution.
- $L(\mu, \sigma; \text{data})$: It means "The likelihood of the parameters μ and σ taking certain values given that we have observed a bunch of data."

- Q) The data set of pass or fail in an exam for 5 students is given in the table below

Hours of studies	Result ($\frac{1}{0} = \text{Pass}$)
29	0
15	0
33	1
28	1
39	1

Use logistic regression as classifier where log of odds is given

$$\log(\text{odds}) = -6.4 + 2 \times \text{hours}$$

a) How to calculate the probability of pass for student who studied 33 hours.

b) Atleast how many hours the student should study that make sue will pass the course with the probability of more than 95 %.

⇒ Question on Calculating Probability with Logistic Regression

Solution :-

a) We will use the sigmoid function to calculate the probability

$$\text{Hours} = 33$$

$$= Z = -64 + 2 \times 33 = 2$$

$$P = \frac{1}{1 + e^{-Z}} = \frac{1}{1 + e^{-2}}$$

$$= 0.88 \Rightarrow 88\%$$

$$b) 0.95 = \frac{1}{1 + e^{-Z}}$$

$$\text{Solve for } z \Rightarrow Z = 2.944$$

$$\Rightarrow 2.944 = -64 + 2 \times \text{hours}$$

$$\Rightarrow 2.944 + 64 = 2 \times \text{hours}$$

$$= \text{Hours} = 33.44$$

• Tree based learner (Classification Problem) \hookrightarrow Supervised Learning technique.

→ Decision tree algorithm

→ Different ways to train a machine.

→ Decision tree → Common

→ random forest → many no of decision tree are combined randomly

→ naive bayes

→ KNN

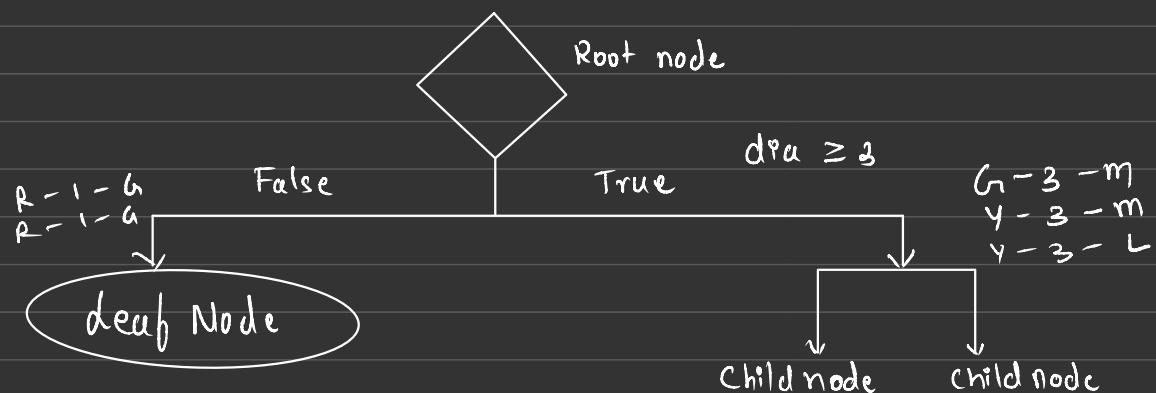


\Rightarrow Decision tree is a graphical representation of all possible solutions to a decision based on certain conditions.

\Rightarrow Random forest build multiple decision tree and merge them together to get a more accurate and stable prediction.

Example :-

Colour	diameter	label
Green	3	Mango
Yellow	3	Mango
Red	1	Grape
Red	1	Grape
Yellow	3	Lemon



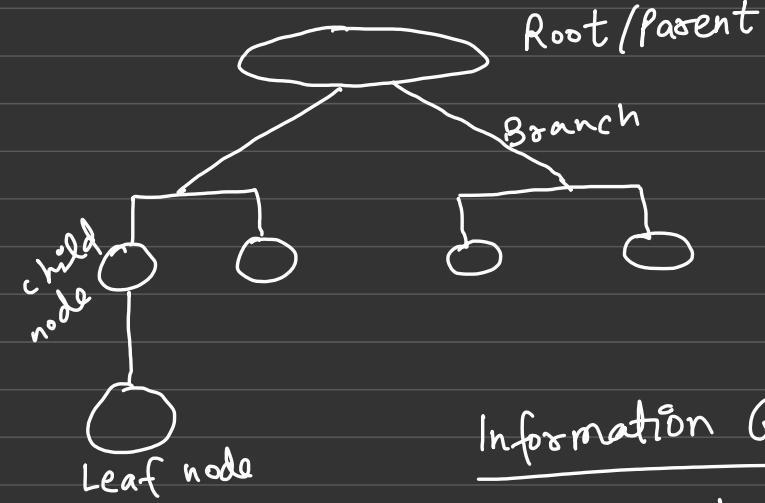
•> Different Terminologies

- 1) Root node :- It is a base node of a tree which represents the entire population or sample and further gets divided into two or more homogenous sets.
(Parent Node)
- 2) Leaf node :- It is one when we reach at the end of the tree it cannot further segregate into any other level.
- 3) Child node :- All root node is always the parent node while all other bottom node associated with that are called child node.
- 4) Splitting :- It is the dividing of root node into different subparts on the basis of some condition.
- 5) Branch/Subtree :- Branch gets formed when you split the tree.
- 6) Pruning :- We are removing the subnodes of a decision tree.
- 7) Gini Index :- It is a measure of impurity or purity used in building a decision tree.
- 8) Information Gain :- It is the decrease in "entropy" after a dataset is split on the basis of an attribute for constructing a decision tree we will select the nodes which have highest information gain.
- 9) Entropy :- $S = -P(yes) \log_2 P(yes) - P(No) \log_2 P(No)$

Unit - 9 : Decision Tree

Classification

Supervised ML



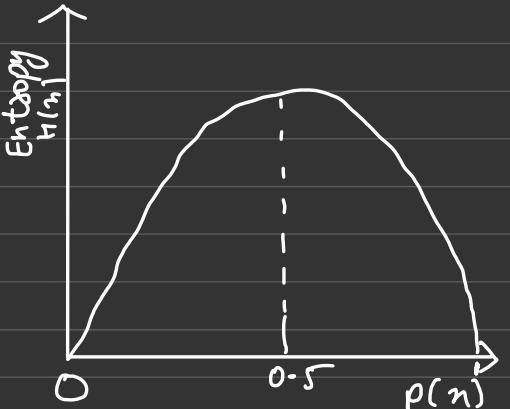
Information Gain :

Entropy

As the probability is zero or one then in that case entropy is zero.

When probability is 0.5 then the value of probability is maximum.

$$E(S) = -P_{\text{Yes}} \log_2 P(\text{Yes}) - P_{\text{No}} \log_2 P(\text{No})$$



- if no. of YES = no. of NO

$$\therefore P(S) = 0.5$$

$$E(S) = 1$$

- if it contains all YES or all NO

$$P(S) = 1 /, P(S) = 0$$

$$E(S) = 0$$

Problem

S: Sunny O: Overcast R: Rain N: Not C: Cold M: Mild

Instances	Outlook	Temperature	Humidity	Windy	Play Tennis
S	H	H	S	No	No
S	H	H	S	Yes	Yes
O	H	H	S	Yes	Yes
R	M	H	S	Yes	Yes
R	C	N	S	No	No
R	C	N	S	Yes	Yes
O	C	N	S	No	No
S	M	H	S	Yes	Yes
S	C	N	S	Yes	Yes
R	M	N	S	Yes	Yes
S	M	N	S	Yes	Yes
O	H	N	S	Yes	Yes

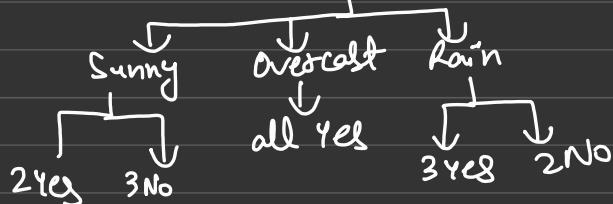
1 data missing

(i) Calculate the entropy for the dataset

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

$$= 0.94$$

(ii) Selection of Root Node from 4 features
 ① Outlook ② Temp ③ humidity ④ Windy



Calculate Entropy for each call.

$$E(\text{outlook/Sunny}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.971$$

$$E(\text{outlook/Overcast}) = 0$$

$$E(\text{outlook/Rain}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log \frac{2}{5} = 0.971$$

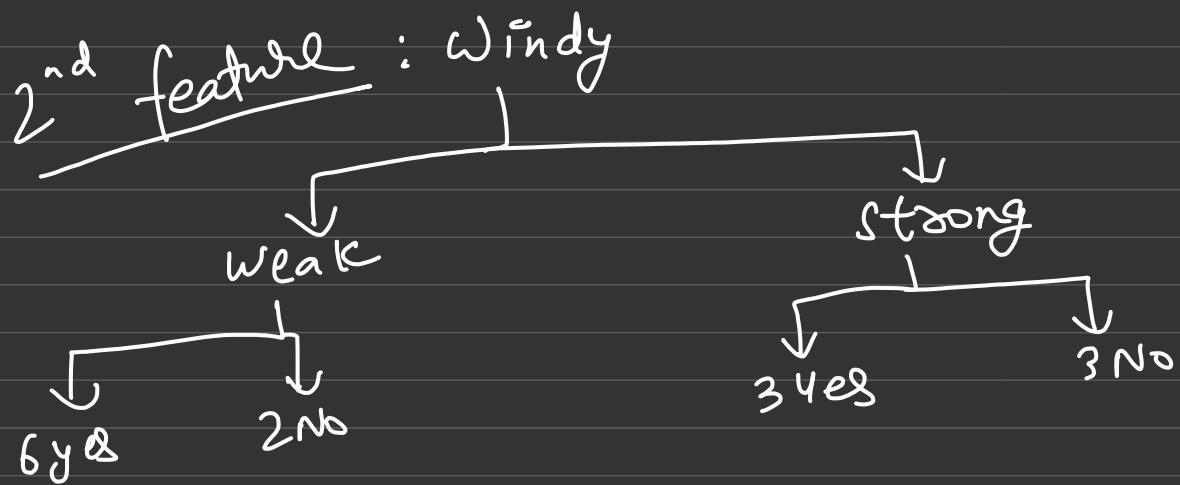
(iii) Calculate how much info you get from outlook

$$I(\text{outlook}) = [\text{Weighted average} \times \text{Entropy (for all features)}]$$

$$= \left[\frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 \right]$$

$$= 0.693$$

$$\begin{aligned}
 \text{Information gain(outlook)} &= E(S) - I(\text{outlook}) \\
 &= 0.94 - 0.693 \\
 &= 0.247
 \end{aligned}$$



$$E(\text{windy} | \text{weak}) = -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} = 0.811$$

$$E(\text{windy} | \text{strong}) = -\frac{3}{6} \log \frac{3}{6} - \frac{3}{6} \log \frac{3}{6} = 1$$

$$I(\text{windy}) = \left[\frac{8}{14} \times 0.811 + \frac{6}{14} \times 1 \right] = 0.892$$

$$\begin{aligned}
 \text{Information gain(Windy)} &= 0.94 - 0.892 \\
 &= 0.048
 \end{aligned}$$

Q)

DAY	Outlook	Temp	Humidity	Wind	Play Tennis
D ₁	Sunny	Hot	High	W	N
D ₂	Sunny	Hot	High	S	N
D ₃	Overcast	Hot	High	W	Y
D ₄	Rain	mild	High	W	Y
D ₅	Rain	cool	Normal	W	Y
D ₆	Rain	cool	Normal	S	N
D ₇	Overcast	cool	Normal	S	Y
D ₈	Sunny	mild	High	W	N
D ₉	Sunny	cool	Normal	W	Y
D ₁₀	Rain	mild	Normal	W	Y
D ₁₁	Sunny	mild	Normal	S	Y
D ₁₂	Overcast	mild	High	S	Y
D ₁₃	Overcast	Hot	Normal	W	Y
D ₁₄	Rain	mild	High	S	N

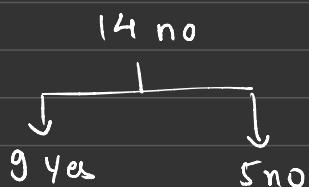
1st Method

$$\begin{aligned} I \cdot G_i &= E(S) - \text{Information} \\ &= E(S) - [\text{Weight} \times \text{Entropy (each feature coverage)}] \end{aligned}$$

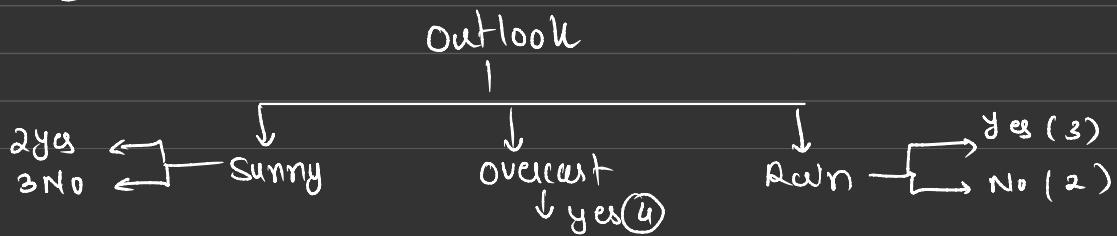
$$E(S) = -P_{\text{yes}} \log_2 P_{\text{yes}} - P_{\text{no}} \log_2 P_{\text{no}}$$

Step I: Calculate E(S)

$$E(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.54$$



Step II Selection of Root node. Calculate I.G. of each feature.

1st attribute

$$E(\text{outlook} | \text{sunny}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5}$$

$$= 0.971$$

$$E(\text{outlook} | \text{overcast}) = 0$$

$$E(\text{outlook} | \text{rain}) = 0.971$$

$$\text{Information (outlook)} = w_f \times \text{Entropy (each feature)}$$

$$= \frac{5}{14} \times 0.971 + \frac{5}{14} \times 0.971 = 0.693$$

$$\text{Information Gain} = E(s) - \text{Information}$$

$$= 0.94 - 0.693 = 0.247$$

The highest information gain of given feature will be the root node.

- Then go to the child nodes and repeat the steps all over again until the decision tree is completed.
- > Gini Index.
- It is a metric used to quantify the amount of uncertainty or impurity present at a single node.
- Opposite of information gain is gini index.
- The feature having lowest gini index is considered as root node.
- > Frame Decision tree using Gini Index.

•> Unit - 9 (Numericals & short questions) {Recap}

- What is Information gain?
- What is Gini index?
- What is Pruning?

Long Question

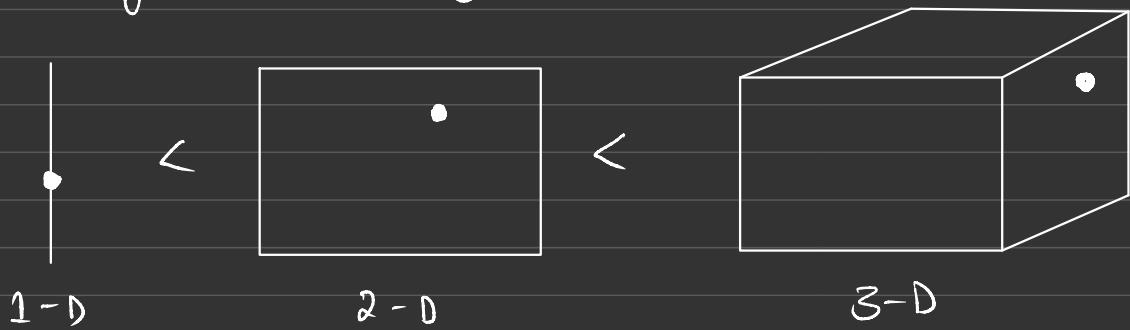
- Numericals (Form D.T.)
- Find 1st Root node.
- Information Gain
- highest → Root node

Decision tree is also known as ID₃ Algorithm.

- ID₃ stands for Iterative Dichotomiser Algorithm
 - ↓
 - It divides the nodes / Branches for better evaluation
- ID₃ is named such because the algorithm iteratively or repeatedly dichotomizes or divides the feature into two or more groups.
- First we are selecting the attribute / feature (Highest IG)
- Then we stop when leaf node appears.

•) Unit - 10 PCA (Principal Component Analysis) / Dimensionality Reduction technique

⇒ "Use of Dimensionality"

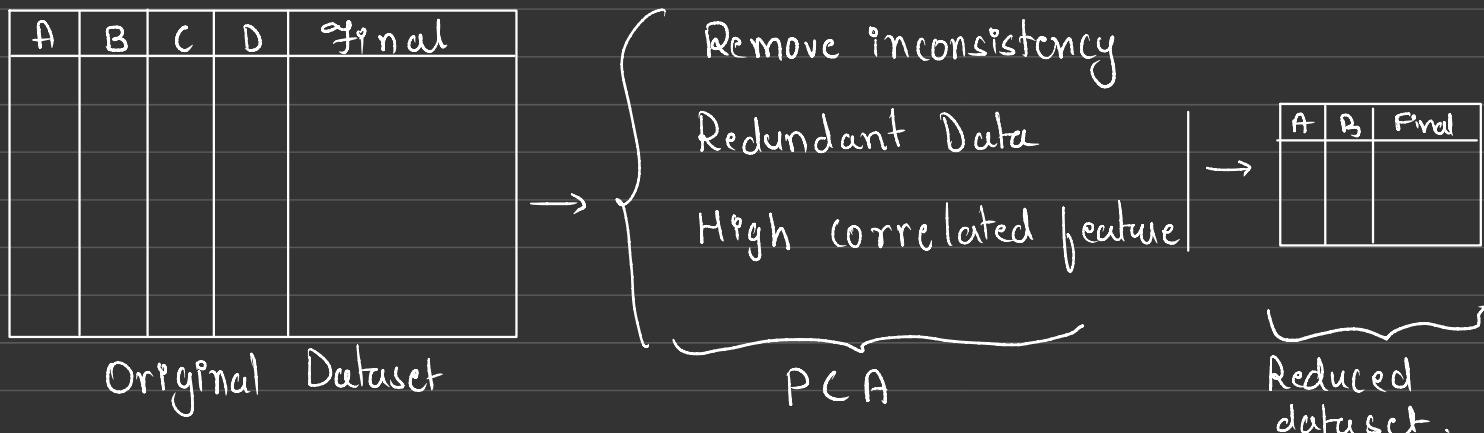


More the dimension difficult to evaluate the information.

⇒ What is PCA ?

→ PCA is a dimensionality reduction technique that helps to identify the correlations and patterns in a dataset so that it can be transformed into a dataset of significantly lower dimension without loss of any important information.

Example



Step 1 :- Computation of mean value of variable

Step 2 :- Computation of covariance of matrix

Step 3 :- Computation of Eigen value , vector, Normalized eigen vector,

Step 4 : Computation of Principal component

Step 5 : Write the reduced dataset.

Q)	Features	Example 1	Ex 2	Ex 3	Ex 4
x	4	8	13	7	
y	11	4	5	14	

⇒ Given the following data use PCA to reduce the dimension from 2 to 1.

Ans → Step 1: Calculate the mean of x and y

$$\Rightarrow \bar{x} = \frac{4+8+13+7}{4} = 8$$

$$\Rightarrow \bar{y} = \frac{11+4+5+14}{4} = 8.5$$

Step 2: Covariance matrix calculation.

Ordered pairs are (x, x) (x, y) (y, x) (y, y)

$$\text{Cov}(x, x) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 \quad N = \text{No of sample.}$$

$$= \frac{1}{4-1} [(4-8)^2 + (8-8)^2 + (13-8)^2 + (7-8)^2]$$

$$= 14$$

$$\text{Cov}(x, y) = \frac{1}{4-1} [(4-8)(11-8.5) + (8-8)(4-8.5) + (13-8)(5-8.5) + (7-8)(14-8.5)] \\ = -11$$

$$\text{Cov}(y, x) = \text{Cov}(x, y)$$

$$\text{Cov}(y, y) = \frac{1}{4-1} [(11-8.5)^2 + (4-8.5)^2 + (5-8.5)^2 + (14-8.5)^2]$$

$$= \frac{1}{3} (6.25 + 20.25 + 12.25 + 30.25)$$

$$= 23$$

Final Covariance matrix will be

$$S \Rightarrow \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix}_{2 \times 2}$$

Step 3: Calculation of Eigen value.

$$\det |(S - \lambda I)| = 0$$

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

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

$$S - \lambda I = \begin{bmatrix} 14 & -11 \\ -11 & 23 \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix}$$

$$= \begin{bmatrix} 14 - \lambda & -11 \\ -11 & 23 - \lambda \end{bmatrix} = 0$$

Simplify:

$$\Rightarrow (14 - \lambda)(23 - \lambda) - 121 = 0$$

$$\Rightarrow \lambda^2 - 37\lambda + 201 = 0$$

$$\Rightarrow \lambda_1, \lambda_2 = \frac{1}{2a} \sqrt{b^2 - 4ac} = \frac{1}{2} \sqrt{37^2 - 4 \times 201}$$

$$\lambda_1 = 30.3849$$

$$= \frac{1}{2} \sqrt{37^2 - 4 \times 201}$$

$$\lambda_2 = 6.6151$$

The value which is larger is considered as the first eigen value.

$\therefore \lambda_1 > \lambda_2$ we will have eigen values as 30.3849

Eigen vector for $\lambda_1 = 30.3849$.

Step 4

We will find out eigen vector of λ_1

$$(S - \lambda_1 I) \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$$

↑ eigen vector for λ_1

$$\Rightarrow \begin{bmatrix} 14 - \lambda_1 & -11 \\ -11 & 23 - \lambda_1 \end{bmatrix} \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} = 0$$

$$\begin{bmatrix} (14 - \lambda_1)v_1 - 11v_2 \\ -11v_1 + (23 - \lambda_1)v_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{aligned} (14 - \lambda_1)v_1 - 11v_2 &= 0 && \text{eq } ① \\ -11v_1 + (23 - \lambda_1)v_2 &= 0 && \text{eq } ② \end{aligned}$$

$$\text{From eqn } ① \quad \frac{v_1}{11} = \frac{v_2}{14 - \lambda_1} = t$$

$$\begin{aligned} \text{When } t = 1, \quad v_1 &= 11 \\ v_2 &= 14 - \lambda_1 \end{aligned}$$

$$\text{Eigen Vector } v_1 \text{ of } \lambda_1 = \begin{bmatrix} 11 \\ 14 - \lambda_1 \end{bmatrix}$$

$$\text{Putting the value of } \lambda_1 = \begin{bmatrix} 11 \\ -16.3849 \end{bmatrix}$$

c) Normalize eigen vector v_1 :

$$e_1 = \begin{bmatrix} 11 / \sqrt{(11)^2 + (-16.3849)^2} \\ -16.3849 / \sqrt{(11)^2 + (-16.3849)^2} \end{bmatrix} = \begin{bmatrix} 0.5574 \\ -0.8303 \end{bmatrix}$$

Similarly for $\lambda_2 = 6.6151$

$$e_2 = \begin{bmatrix} 0.8303 \\ 0.5574 \end{bmatrix}$$

Step 5:

first PC ₁	P ₁₁	P ₁₂	P ₁₃	P ₁₄
--------------------------	-----------------	-----------------	-----------------	-----------------

↪ Principal Component $x_1 - \bar{x}$

$$P_{11} = e_1^T \begin{bmatrix} 4-8 \\ 11-8.5 \\ y_1 - \bar{y} \end{bmatrix}$$

$$= [0.5574 \quad -0.8303] \begin{bmatrix} -4 \\ 2.5 \end{bmatrix} = -4.3652$$

$$P_{12} = 3.7361$$

$$P_{13} = 5.6928$$

$$P_{14} = -5.1238$$

PC ₁	-4.3652	3.7361	5.6928	-5.1238
-----------------	---------	--------	--------	---------

"Solution ends Here"

•> Naive Bayes Theorem / Classification

⇒ Test sample will be given & we have to find which class the test sample is using Naive Bayes.

$$\text{Naive Bayes theorem} = P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \rightarrow \begin{array}{l} \text{Likelihood} \\ \downarrow \\ \text{Posterior probability} \end{array} \quad \begin{array}{l} \text{Prior} \\ \downarrow \\ \text{evidence} \end{array}$$

$$\Rightarrow \text{Posterior} = \frac{\text{Prior} \times \text{Likelihood}}{\text{evidence}}$$

Q) dataset

Example No	Input Attribute			Target attribute
	Color	Type	Origin	Stolen
1				
2				
3				
4				
5				
6				
7				
8				
9				
10				

Step 1 : Calculate priori probability

$$\Rightarrow P(yes) = \frac{5}{10} = 0.5 \quad p(No) = \frac{5}{10} = 0.5$$

Step 2: Calculate Conditional probability

Color	Color		Type	Type		Origin	Origin	
	Yes	No		Yes	No		Yes	No
Red	3/5	2/5	Sport	4/5	2/5	Domestic	2/5	3/5
Yellow	2/5	3/5	SUV	1/5	3/5	Imported	3/5	2/5

\Rightarrow Step 3: New instance calculation.

$$\begin{aligned}
 P(\text{yes} | \text{new instance}) &= P(\text{yes}) * P(\text{Color} = \text{Red} | \text{yes}) * \\
 &\quad P(\text{Type} = \text{SUV} | \text{yes}) * P(\text{Origin} = \text{Domestic} | \text{yes}) \\
 &= \frac{5}{10} * \frac{3}{5} * \frac{1}{5} * \frac{2}{5} \\
 &\Rightarrow \frac{1}{5} * \frac{1}{5} * 3 * \frac{1}{5} = \frac{3}{125} = 0.024
 \end{aligned}$$

$$\begin{aligned}
 P(\text{No} \mid \text{New instance}) &= P(\text{No}) * P(\text{Color = Red} \mid \text{No}) * P(\text{Type = SUV} \mid \text{No}) \\
 &\quad * P(\text{Origin = Domestic} \mid \text{No}) \\
 &= \frac{1}{10} * \frac{2}{5} * \frac{1}{5} * \frac{3}{5} = \frac{9}{125} \\
 &= 0.072
 \end{aligned}$$

Compare these two posterior probability

$$P(\text{No} \mid \text{new instance}) > P(\text{yes} \mid \text{New instance})$$

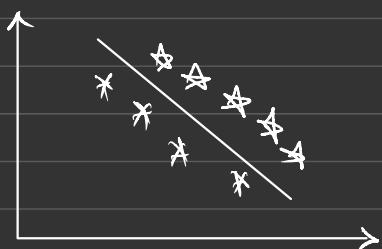
We can classify it as it belongs to No class.

•) Unit - II Support vector Machine (SVM)

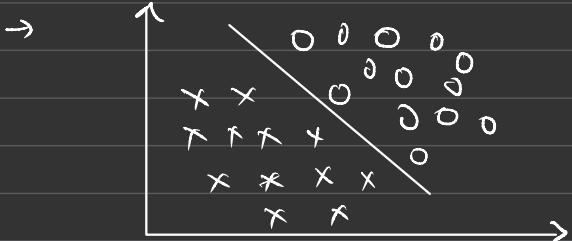
SVM Hypoplane is used (Objective is to find this hyperplane in n dimension.)

→ What is SVM?

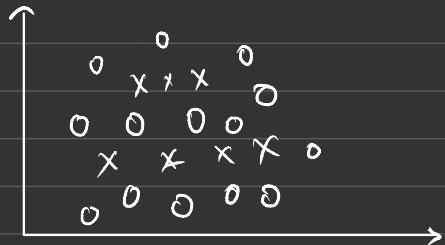
→ SVM is a discriminative classifier ie formally designed by a separative hyperplane such that the points of different categories are separated by a gap as wide as possible.



Linearly Separable Dataset



Non Linearly Separable Dataset



•> 1-D



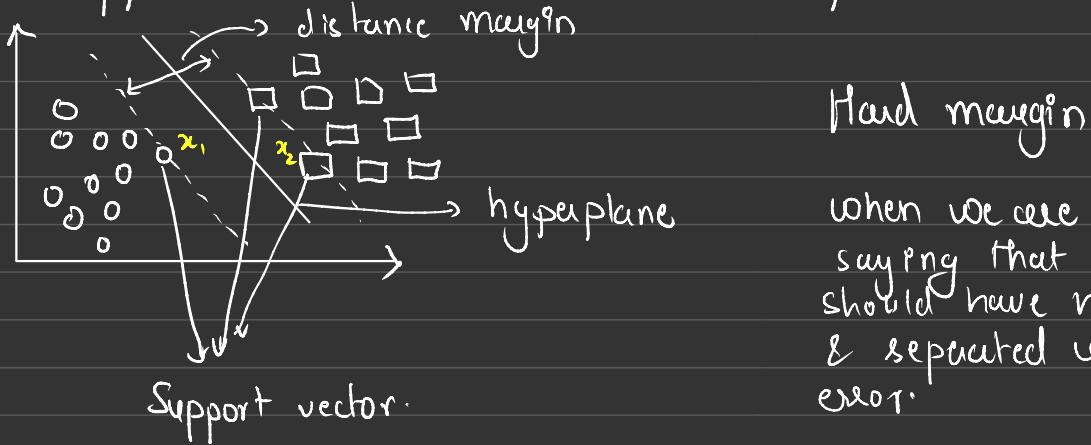
1-D plane to 2D plane

When there is no clear hyperplane it is necessary to move away from single dimensional view of the data to a 2-D view.

This process is called as Kerneling where it transforms low dimension data to high dimensional data to create a clear hyperplane.

•> Derivation of SVM

→ The Support Vectors are the extreme points in the dataset.



when we are strictly saying that classes should have no error & separated without any error.

$$r = \frac{g(n)}{\|\omega\|} \rightarrow \text{euclidean norm / distance of } \omega.$$

⇒ We want to compute the distance between x_2 and x_1
so the equation will be.

$$\begin{aligned} \omega^T x_2 + \omega_0 &= +1 \\ \omega^T x_1 + \omega_0 &= -1 \\ \xrightarrow{(-)} \quad \xrightarrow{(-)} \quad \xrightarrow{(+)} \\ \omega^T (x_2 - x_1) &= +2 \end{aligned}$$

$$\frac{\omega^T}{\|\omega\|} (x_2 - x_1) = \frac{2}{\|\omega\|} \Rightarrow (x_2 - x_1) = \frac{2}{\|\omega\|}$$

\Rightarrow Here we got our objective function, we need to update Objective function (ω^*, ω_0^*) so that $\frac{2}{\|\omega\|}$ should be maximized.

$$[\omega_1, \omega_n]$$

$$\Rightarrow \left[(\omega^{T*}, \omega_0^*) = \min \frac{\|\omega\|}{2} \right] \text{ eqn for hard margin}$$

$$(\omega^{T*}, \omega_0^*) = \min \frac{\|\omega\|}{2} + C_i \sum_{i=1}^n \xi_i$$

How many error
 model can consider
 eqn for
 soft margin

↓
 value of error.