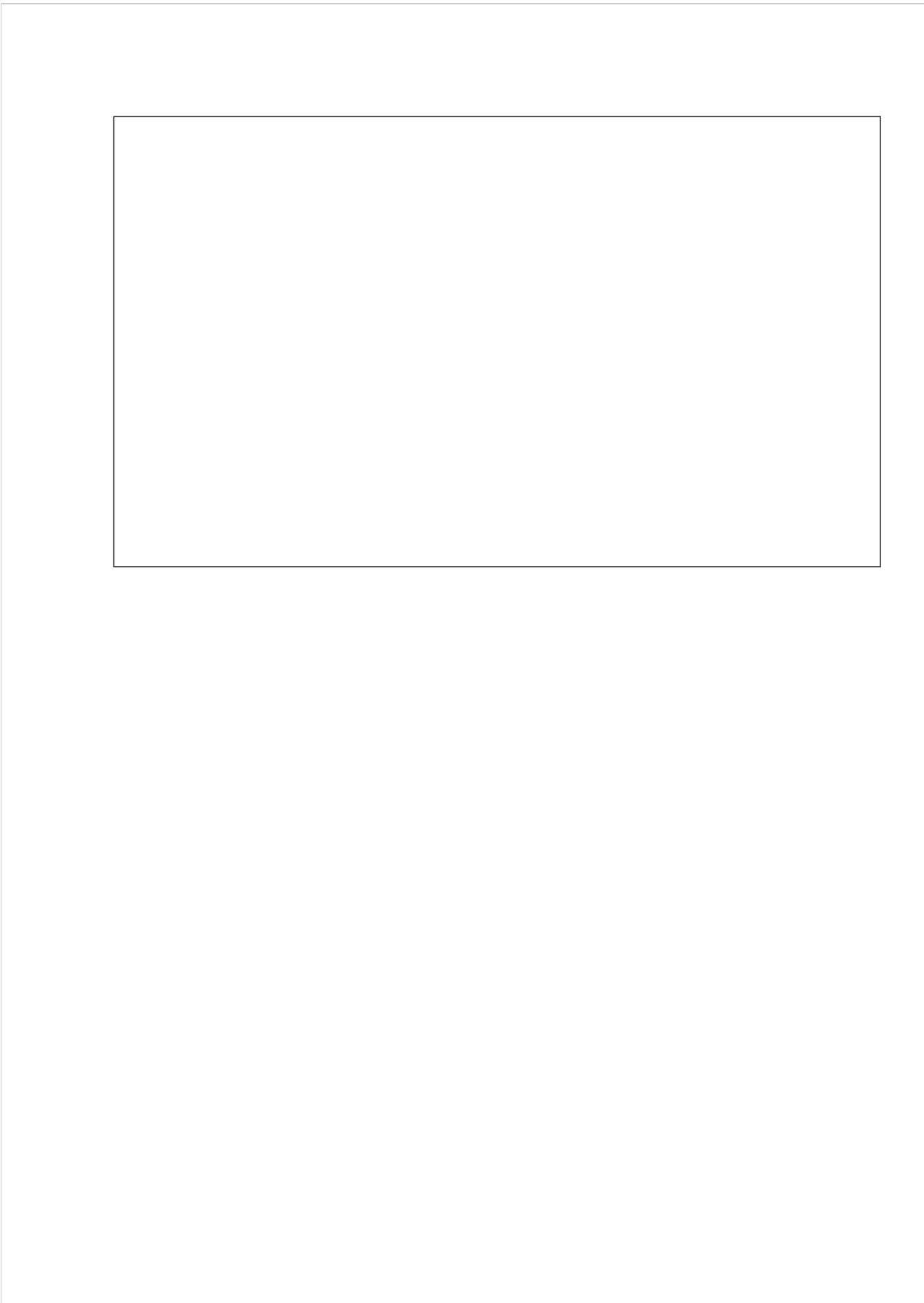


ML_QP2_Format_ESE_50 marks_R-2022

 VIT Vidyalankar Institute of Technology <small>Accredited A+ by NAAC</small> (Autonomous College Affiliated to University of Mumbai)	End Semester Examination (R-2022 Scheme) -(2024-25)																										
Date:	Branch: CMPN	Time: 2 Hrs.																									
Semester: VII	Subject: Machine Learning	Marks: 50																									
N.B. :- All Questions are Compulsory			CO																								
Q. 1)	Attempt any two (5 Marks Each) a) Compare and contrast supervised, unsupervised, and reinforcement learning. Provide an example for each type. b) Elaborate on Karl Pearson's coefficient of Correlation. c) Elaborate on following terms with respect to DBSCAN: (i) Directly Density Reachable (ii) Density Reachable																										
Q. 2)	Attempt any two. (5 Marks Each) a) Justify the need of cross validation. Explain k fold cross validation in details. b) <table border="1" style="margin-left: auto; margin-right: auto;"> <thead> <tr> <th></th> <th style="text-align: center;">INDEPENDENT</th> <th style="text-align: center;">DEPENDENT</th> </tr> <tr> <th style="text-align: center;">SR.NO</th> <th style="text-align: center;">X</th> <th style="text-align: center;">Y</th> </tr> </thead> <tbody> <tr><td style="text-align: center;">1</td><td style="text-align: center;">43</td><td style="text-align: center;">99</td></tr> <tr><td style="text-align: center;">2</td><td style="text-align: center;">21</td><td style="text-align: center;">65</td></tr> <tr><td style="text-align: center;">3</td><td style="text-align: center;">25</td><td style="text-align: center;">79</td></tr> <tr><td style="text-align: center;">4</td><td style="text-align: center;">42</td><td style="text-align: center;">75</td></tr> <tr><td style="text-align: center;">5</td><td style="text-align: center;">57</td><td style="text-align: center;">87</td></tr> <tr><td style="text-align: center;">6</td><td style="text-align: center;">59</td><td style="text-align: center;">81</td></tr> </tbody> </table> Tasks: Fit a linear regression model to the data and determine the equation of the regression line.				INDEPENDENT	DEPENDENT	SR.NO	X	Y	1	43	99	2	21	65	3	25	79	4	42	75	5	57	87	6	59	81
	INDEPENDENT	DEPENDENT																									
SR.NO	X	Y																									
1	43	99																									
2	21	65																									
3	25	79																									
4	42	75																									
5	57	87																									
6	59	81																									
c)	How does the AdaBoost algorithm iteratively adjust weights and combine weak classifiers to form a strong classifier, and what are the main steps involved in this process?																										
Q. 3)	Attempt any two. (5 Marks Each) a) Justify the need of pruning in XGBOOST Algorithm. b) How does a polynomial kernel work, and how does it help in modeling non-linear relationships in Support Vector Machines (SVMs)? c) How does graph-based clustering using a Minimum Spanning Tree (MST) work to identify clusters in a dataset, and how would you apply this method to a given graph example to demonstrate the process of separating clusters?																										

Q 4)	Attempt any One (10 Marks Each)																																																									
a)	<table border="1"> <thead> <tr> <th></th><th>Weekend</th><th>Weather</th><th>Parents</th><th>Money</th><th>Decision</th></tr> </thead> <tbody> <tr><td>W1</td><td>Sunny</td><td>Yes</td><td>Rich</td><td>Cinema</td></tr> <tr><td>W2</td><td>Sunny</td><td>No</td><td>Rich</td><td>Tennis</td></tr> <tr><td>W3</td><td>Windy</td><td>Yes</td><td>Rich</td><td>Cinema</td></tr> <tr><td>W4</td><td>Rainy</td><td>Yes</td><td>Poor</td><td>Cinema</td></tr> <tr><td>W5</td><td>Rainy</td><td>No</td><td>Rich</td><td>Stay In</td></tr> <tr><td>W6</td><td>Rainy</td><td>Yes</td><td>Poor</td><td>Cinema</td></tr> <tr><td>W7</td><td>Windy</td><td>No</td><td>Poor</td><td>Cinema</td></tr> <tr><td>W8</td><td>Windy</td><td>No</td><td>Rich</td><td>Shopping</td></tr> <tr><td>W9</td><td>Windy</td><td>Yes</td><td>Rich</td><td>Cinema</td></tr> <tr><td>W10</td><td>Sunny</td><td>No</td><td>Rich</td><td>Tennis</td></tr> </tbody> </table> <p>For the dataset given below, construct a decision tree using Gini Index, and determine which attribute is a root attribute.</p>		Weekend	Weather	Parents	Money	Decision	W1	Sunny	Yes	Rich	Cinema	W2	Sunny	No	Rich	Tennis	W3	Windy	Yes	Rich	Cinema	W4	Rainy	Yes	Poor	Cinema	W5	Rainy	No	Rich	Stay In	W6	Rainy	Yes	Poor	Cinema	W7	Windy	No	Poor	Cinema	W8	Windy	No	Rich	Shopping	W9	Windy	Yes	Rich	Cinema	W10	Sunny	No	Rich	Tennis	CO4
	Weekend	Weather	Parents	Money	Decision																																																					
W1	Sunny	Yes	Rich	Cinema																																																						
W2	Sunny	No	Rich	Tennis																																																						
W3	Windy	Yes	Rich	Cinema																																																						
W4	Rainy	Yes	Poor	Cinema																																																						
W5	Rainy	No	Rich	Stay In																																																						
W6	Rainy	Yes	Poor	Cinema																																																						
W7	Windy	No	Poor	Cinema																																																						
W8	Windy	No	Rich	Shopping																																																						
W9	Windy	Yes	Rich	Cinema																																																						
W10	Sunny	No	Rich	Tennis																																																						
b)	An insurance company wants to predict the likelihood of claims based on customer demographics and policy details. Suggest how a Random Forest model can be used for this task and explain how bagging improves the model's robustness and accuracy. Additionally, compare the performance of the Random Forest model with that of a single decision tree.	CO4																																																								
Q 5)	Attempt any One (10 Marks Each)																																																									
a)	How do we derive the loss function in a Support Vector Machine (SVM) that penalizes misclassification errors, and how does this loss influence the overall cost function for the classifier?	CO5																																																								
b)	Explain SVD as dimensionality reduction technique. Consider the matrix A as given below and find the Eigen values in SVD for ATA.	CO6																																																								
	$A = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}$																																																									



Q1.A Compare and contrast supervised, unsupervised, and reinforcement learning. Provide an example for each type.

Solution

① Supervised Learning

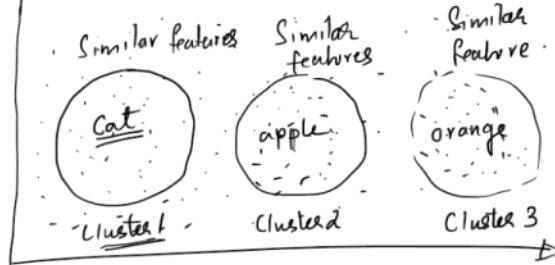
- Model is trained on a labeled dataset
- goal is to learn the mapping between input and outputs, enabling it to make predictions on unseen data.
- Model learns by comparing its predictions with the correct output and adjust its internal parameters to minimize the prediction error.

Ex → Image Classification → Given labeled images the model learns to classify new images.

Spam Detection: Training on emails labelled as Spam or not Spam, the model can filter new emails.

UnSupervised Learning

Given 10000 images and no label are specified.



- * Here the input is only data and no labels are associated with data.
- * Not Sure about type of output
- * Unsupervised Algo will work on 10000 images and will create clusters of images based on the similarities.
- ↗ It is "SELF LEARNING"

Semi Supervised

Text Document Classification

Ex 1000000 articles and need to classify them

into

- news
- literature
- Research Paper
- medical Report

Label is not provided

Manually labelling the 1000000 articles not possible.

→ We will label 10000 articles [Supervised Learning]

→ My model will be trained on these 10000 articles and will use the pattern identified to classify the remaining 990000 articles [Unsupervised]

- * Uses small amount of labelled data
- * and large amount of unlabelled data.
- * Benefits of Both labelled and unlabelled data.
- * Overcome the challenge of finding large amount of labelled data.

Eg ① Speech Analysis → we will label audio files
 which will need lot of human resource and
 thus we will use some technique that will help
 to improve traditional speech analysis tool.

② Web Content Classification → Organizing the
 knowledge available on billions of web pages
 will need advance processing, but the
 task needs human "Intervention" to classify

Q1.B Elaborate on Karl Pearson's coefficient of Correlation.

Solutions

① Karl Pearson's Coefficient of Correlation (r)

→ To calculate relationship bet two variable

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

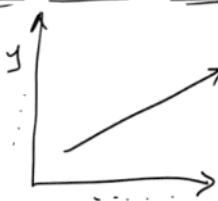
r quantifies the strength of relationship betⁿ two variables.

The value of r be between $\underline{\underline{+1}}$ and $\underline{\underline{-1}}$

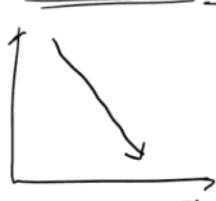
If $r = 1 \Rightarrow$ Total +ve correlation \Rightarrow If $x \uparrow$ then $y \uparrow$

If $r = -1 \Rightarrow$ Total -ve correlation \Rightarrow If $x \downarrow$ then $y \uparrow$

→ r gives strength (degree) & direction of correlation. or $x \uparrow$ then $y \downarrow$



+ve
Correlation



-ve
Correlation



Zero
Correlation

Q1.C Elaborate on following terms with respect to DBSCAN:

- (i) Directly Density Reachable
- (ii) Density Reachable

Solution >

Directly Density Reachable \rightarrow A point X is directly density reachable from point Y w.r.t epsilon & Minpt if

1. X belongs to neighbourhood of Y ($dist(X, Y) \leq \epsilon$)

2. Y is core point.

$$\underline{\text{Minpt}} = 4$$

① \circlearrowleft Y is a core point \Rightarrow

no of points in neighbourhood of $Y = 5 > \text{Minpt}$

$\therefore Y$ is core point

② X belongs to neighbourhood of Y .

$\therefore X$ is Directly Density Reachable from Y .

② Density Reachable \rightarrow

A point X is density reachable from Point Y w.r.t ϵ & Minpt

if there is a chain of points $P_1, P_2 \dots P_n$ and $P_1 = X$

& $P_n = Y$ such that P_{i+1} is directly density reachable from P_i

Minpt = 4 here X & Y are core points

Here P_2 is directly density reachable from X

P_3 is directly density reachable from P_2

Φ is directly density reachable from P_3

$\therefore X$ is Density Reachable from Y .

Q2.A Justify the need of cross validation. Explain k fold cross validation in details.

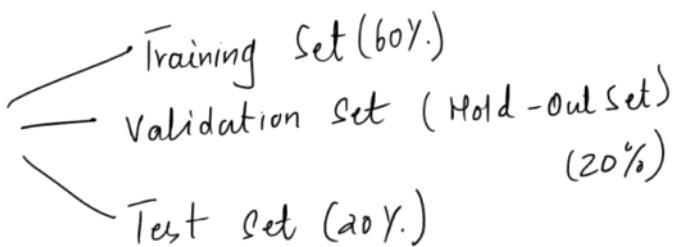
Solutions

Cross Validation →

- In Supervised M.L
- Train a Model on a Dataset
- Trained model is used to predict the target given new sample.
- How to know if model we have trained will produce effective and accurate result on new input

Cross Validation → It is process that ensures the model will perform well on new Data.

Split the Dataset



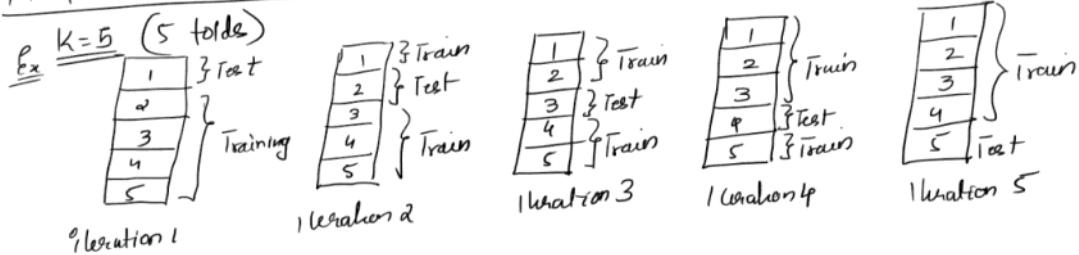
Training Set = part of data on which model is trained.
(This dataset will help to * build model)

Training Set = part of data on which model is trained.
 (This dataset will help to * build model)

* Validation Set ⇒ * Evaluate the Model

- Will help to chk if model overfits or Underfits.
- Update the parameters and again train the model.
- Repeat this until the model performs best on Validation set.

K-fold Cross Validation ⇒

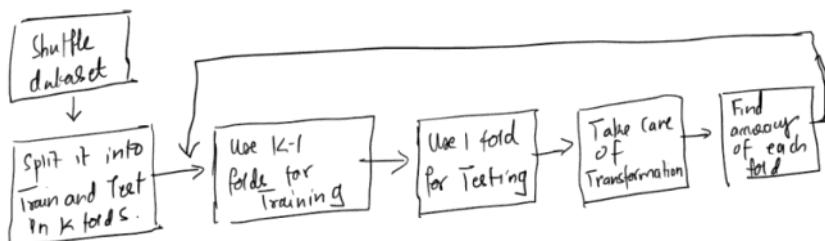


K fold ⇒ K fold helps us to build the model in generalized form.

→ To achieve this K fold Cross Validation splits the dataset into Training Testing & Validation.

→ Here Test and Train data will support building the model.

→ Life cycle of K-fold Cross Validation.



* The No of iterations ideally is K time

* Finding mean of accuracy score of each iteration will give the consistency of the Trained model.

Rules

① $K \geq 2$

if $K=2 \rightarrow$ just 2 iterations.

if $K = n \geq 2$ \Rightarrow $n-1$ for Training
1 for Testing

② most commonly used value of $K=10$

③ If K is very large then the running time of process will increase.

④ The value of K is inversely proportional to size of data i.e. if dataset size is small then number of folds can increase.

Q2.B INDEPENDENT DEPENDENT

SR.NO

X	Y
1	43
2	21
3	25
4	42
5	57
6	59
	99
	75
	87
	81

Tasks:

Fit a linear regression model to the data and determine the equation of the regression line.

Solution →

The linear regression line is represented as: $y = a + bx$

$$\text{where } b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$a = \frac{\sum y - b \sum x}{n}$$

$$\sum x = 247$$

$$\sum y = 486$$

$$\sum x^2 = 11409$$

$$\sum xy = 40022$$

$$\sum x y = 20485$$

$$\therefore y = 65.14 + 0.39x$$

Q2.C How does the AdaBoost algorithm iteratively adjust weights and combine weak classifiers to form a strong classifier, and what are the main steps involved in this process?

Solution: →

Summary of AdaBoost

- ① Assign Sample weight (Initial Equal)
- ② Create stump for each feature.
- ③ Use $Gini$ index to identify first stump.

- ④ For first stump:
 - ① Calculate Total Error
 - ② Amt of Say

- ⑤ Update the sample weights for each sample

$$\text{For Incorrectly classified } \frac{\text{updated weight}}{\text{old weight}} = e^{+\text{Amt of Say}}$$

$$\text{For Correctly classified } \frac{\text{updated weight}}{\text{old weight}} = e^{-\text{Amt of Say}}$$

- ⑥ Now Specify updated weight for each sample & Normalize the weight to generate New sample weight.

- ⑦ Identify new stump based on New sample weight & Repeat step 3 to 7

- ⑧ For Classification (Testing).

8.1 \Rightarrow Run the test sample through all the stumps.

8.2 \Rightarrow Calculate Amt of Say for sample classifying Yes & No

8.3 \Rightarrow Classify the test sample based on largest sum of Amount of Say

Idea Behind AdaBoost (Boosting Technique \Rightarrow reduce Bias)

1. Adaboost combines lot of weak learners to make classification
2. The weak learners are almost always Stump.
3. Some Stump will have more say in classification than other Stump.
4. Each Stump is made by taking previous Stump mistakes into account

Q3.A Justify the need of pruning in XGBOOST Algorithm.

Solution+

Pruning: Kyon chahiye

- ① Right distribution of Tree
- ② Ignoring Unnecessary branches
- ③ Decrease Computation
- ④ Less Time
- ⑤ Avoid Overfitting

Pruning of trees in XGBoost is purely done on gain

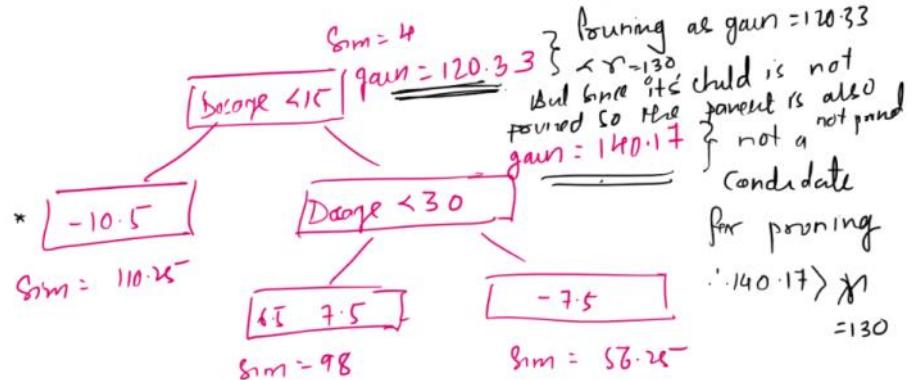
Let us Assume a Threshold gain = γ [gamma].

Rule for Pruning * If at branch $\frac{\text{gain} - \gamma}{\text{sum}}$ < 0 then prune and go above

* If $\text{gain} - \gamma \geq 0$ then do not Prune

* If child is not pruned then parent also will not be pruned]

Tree $\lambda = 0$
Let $\gamma = 130$



* Note
If α is sufficiently large then it may result into complete tree pruning
till root \rightarrow Extreme Pruning

* Regardless of value of $\underline{\lambda}$ and $\underline{\eta}$ let us assume the tree as:

Q3.B How does a polynomial kernel work, and how does it help in modeling non-linear relationships in Support Vector Machines (SVMs)?

Solution:

Polynomial Kernel

* It calculates higher Dimensional Relationship between observations (data points).

* The kernel that transforms $\underline{\underline{1D}}$ to $\underline{\underline{dD}}$ is Polynomial Kernel.

* It may look like $(\underline{\underline{a}} \times \underline{\underline{b}} + \gamma)^d$

where a and b are two different observations in dataset
(any two data points).

γ = coefficient of Polynomial.

d = degree of Polynomial.

* Here we use SVM with Polynomial Kernel to compute relationship betⁿ observations in higher dimension and then find good classifier.

* Let $\underline{\underline{a}}$ and $\underline{\underline{b}}$ be two observations.

* Let $\gamma = 1/2$ & $d = 2$ } Note \Rightarrow value of γ and d is determined by cross validation.

$$\begin{aligned}
 \Rightarrow (\underline{\underline{a}} \times \underline{\underline{b}} + \gamma)^d &= (\underline{\underline{a}} \times \underline{\underline{b}} + 1/2)^2 \\
 &= (\underline{\underline{a}} \times \underline{\underline{b}} + 1/2) \cdot (\underline{\underline{a}} \times \underline{\underline{b}} + 1/2) \\
 &= a^2 b^2 + \frac{1}{2} ab + \frac{1}{2} ab + \frac{1}{4} \\
 &= \underline{\underline{a}} \cdot \underline{\underline{b}} + a^2 b^2 + \frac{1}{4} \quad \Rightarrow \text{Polynomial.}
 \end{aligned}$$

= can be written as dot product of

$$\Rightarrow = \left(\underline{\underline{a}}, \underline{\underline{a}}^2, \frac{1}{2} \right) \cdot \left(\underline{\underline{b}}, \underline{\underline{b}}^2, \frac{1}{2} \right)$$

For data point a

original value $\underline{\underline{a}}$
Higher dimension value of a

original value $\underline{\underline{b}}$
value of b in higher dimension

Dot product is sum of 1st term multiplied and terms multiplied and so on]

Note: $\frac{1}{2}$ is third axis but value same so ignore.

So $(a+b)^d$ is used to get higher dimension "all" between two data points $\underline{\underline{a}}$ & $\underline{\underline{b}}$.

Q3.C How does graph-based clustering using a Minimum Spanning Tree (MST) work to identify clusters in a dataset, and how would you apply this method to a given graph example to demonstrate the process of separating clusters?

Solution =

- * Graph Based Clustering Using Minimum Spanning Tree →
 - Spanning Tree
 - Subgraph of given graph
 - Some No of vertex
 - and no cycle
 - Subgraph is connected (All vertex)

Algo

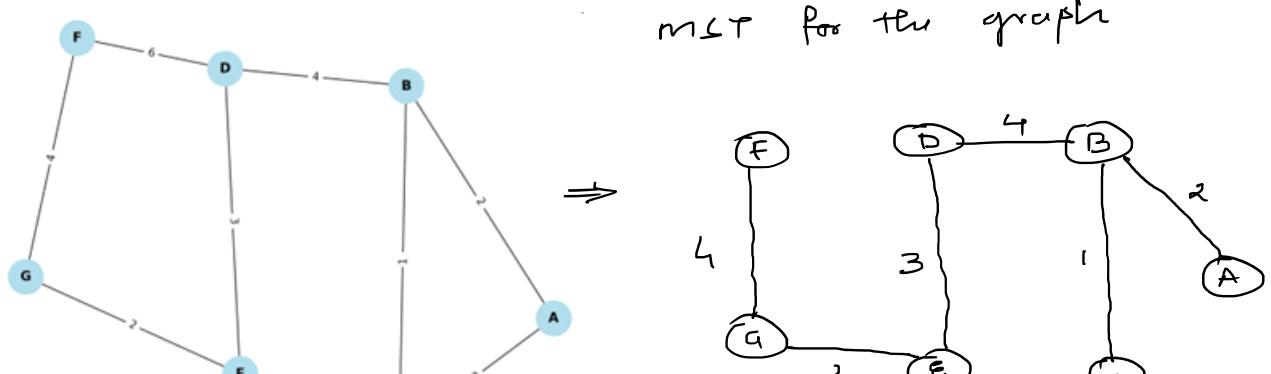
1. Determine MST of given graph.
2. Delete branch iteratively.
3. Each connected component is a cluster

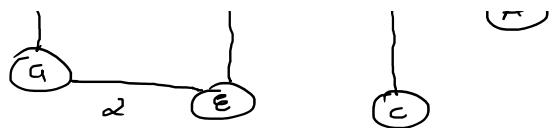
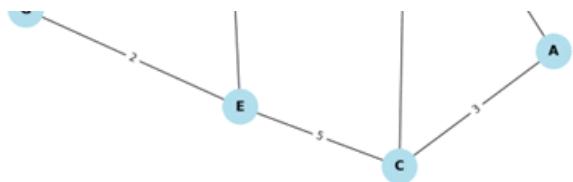
Different Strategies to Delete Branches.

- ① Delete Branch with Max weight ↴.
- ② Delete Inconsistent Branch →
- ③ Delete by Analysis of weight.

Delete Inconsistent Branches.

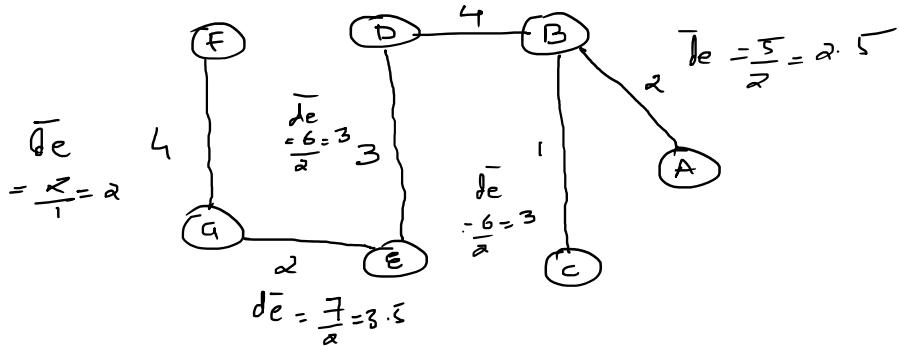
- * A branch is inconsistent if the corresponding weight (d_e) is much larger than the reference value \bar{d}_e
- * The reference value \bar{d}_e can be defined by the average weight of all the branches adjacent to edge e



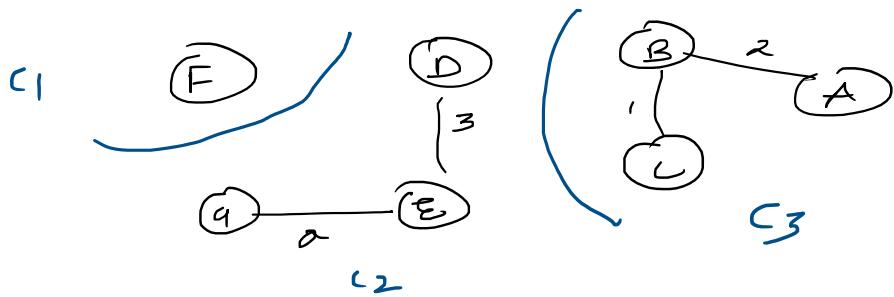


Calculating \bar{d}_e for every edge

$$\bar{d}_e = \frac{6}{3} = 2$$



can delete edge with case 4 i.e FG. and DB

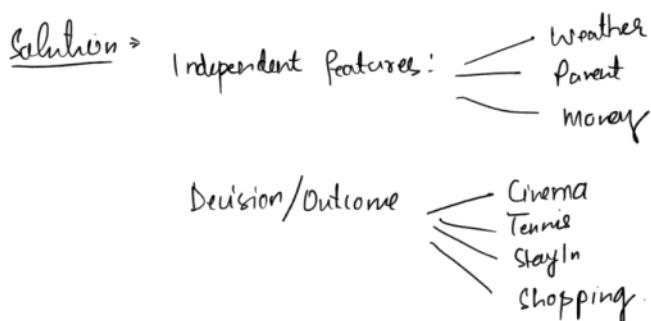


we have 3 clusters.

Q4.A a) Weekend Weather Parents Money Decision

W1	Sunny	Yes	Rich	Cinema
W2	Sunny	No	Rich	Tennis
W3	Windy	Yes	Rich	Cinema
W4	Rainy	Yes	Poor	Cinema
W5	Rainy	No	Rich	Stay In
W6	Rainy	Yes	Poor	Cinema
W7	Windy	No	Poor	Cinema
W8	Windy	No	Rich	Shopping
W9	Windy	Yes	Rich	Cinema
W10	Sunny	No	Rich	Tennis

For the dataset given below, construct a decision tree using Gini Index, and determine which attribute is a root attribute.

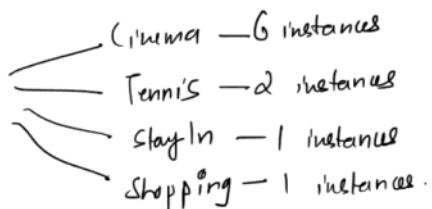


Step 1

We will calculate Gini Index for Overall collection of Outcomes

of Training Examples -

These are 4 possible outcomes for decision

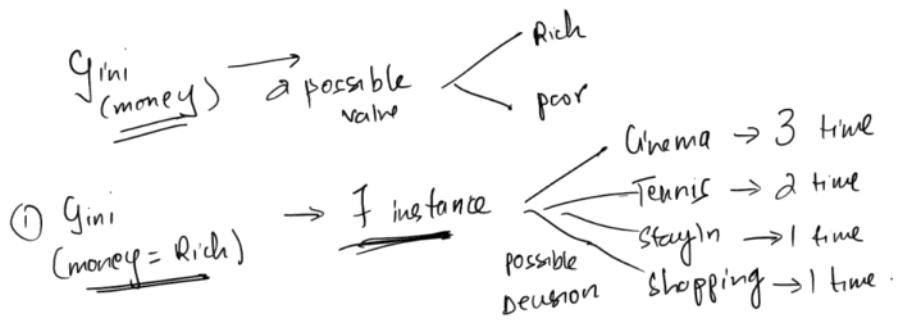


$$\begin{aligned}
 G_{\text{ini}} &= 1 - \left(\underline{\underline{(6/10)^2}} + \underline{\underline{(2/10)^2}} + \underline{\underline{(1/10)^2}} + \underline{\underline{(1/10)^2}} \right) \\
 (\text{devision}) &= 1 - \left(\frac{42}{100} \right) = \underline{\underline{0.58}}
 \end{aligned}$$

Note \rightarrow In Machine Learning, Gini index/coefficient is utilized as an Impurity measures in decision tree for classification.

$$G_{\text{ini}} = 1 - \sum_{i=1}^n (P_i)^2 \quad \text{where } P_i \text{ probability of outcome of specific class}$$

Step 2: To find Gini Index for Money



$$\begin{aligned}
 Gini_{(money=Rich)} &= 1 - \left(\left(\frac{3}{7}\right)^2 + \left(\frac{2}{7}\right)^2 + \left(\frac{1}{7}\right)^2 + \left(\frac{1}{7}\right)^2 \right) \\
 &= \underline{0.694}
 \end{aligned}$$

② $Gini_{(money=poor)}$ → 3 instance
Decision

$$= 1 - \left(\left(\frac{3}{3}\right)^2 \right) = 0_{11}$$

$$\begin{aligned}
 \text{Weighted Average } Gini_{(money)} &= \left(Gini_{(money=Rich)} * \text{proportion of rich} \right) + \left(Gini_{(money=poor)} * \text{proportion of poor} \right) \\
 &= \left(0.694 * \frac{7}{10} \right) + \left(0 * \frac{3}{10} \right)
 \end{aligned}$$

$$\boxed{Gini_{(money)} = \underline{0.485}}$$

Step 3 Gini Index on Parent

For Parent feature $\xrightarrow[\text{value}]{\text{possible}}$ Yes
No

$$\begin{aligned} \underline{\text{Gini}}_{(\text{parent}=\text{Yes})} &= 5 \text{ instances} \xrightarrow[\text{possible decision}]{\text{Cinema}} \\ &= 1 - ((5/5)^2) = 0.00 \end{aligned}$$

$$\begin{aligned} \underline{\text{Gini}}_{(\text{parent}=\text{No})} &= 5 \text{ instances} \xrightarrow[\text{possible decision}]{\text{Tennis} \rightarrow 2 \text{ times}, \text{StayIn} \rightarrow 1 \text{ times}, \text{Shopping} \rightarrow 1 \text{ time}, \text{Cinema} \rightarrow 1 \text{ time}} \\ &= 1 - \left(\frac{2}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 = 0.72 \end{aligned}$$

$$\begin{aligned} &= 1 - \left(\frac{2}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 + \left(\frac{1}{5} \right)^2 \\ &= 1 - \left(\frac{2}{5} \right) = 0.72 \end{aligned}$$

Weighted Average of Gini $= (0 * 5/10) + (0.72 * 5/10) = \underline{0.36}$

$$\boxed{\underline{\text{Gini}_{(\text{parent})} = 0.36}}$$

Step 4 Gini Index for Weather

Weather $\xrightarrow[\text{value}]{\text{possible}}$ Sunny = 3 instances
Windy = 4 instances
Rainy = 3 instances.

G_{ini} (weather = sunny) \Rightarrow 3 instances possible outcomes

$$= 1 - \left(\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right) = \underline{\underline{0.444}}$$

Cinema \rightarrow 1 time

Tennis \rightarrow 2 time

G_{ini} (weather = Windy) \Rightarrow 4 instances possible outcomes

$$= 1 - \left(\left(\frac{3}{4}\right)^2 + \left(\frac{1}{4}\right)^2 \right) = \underline{\underline{0.375}}$$

3 time Cinema

1 time Shopping

G_{ini} (weather = Rainy) \Rightarrow 3 instances possible outcomes

$$= 1 - \left(\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right) = \underline{\underline{0.444}}$$

2 cinema

1 Stay In

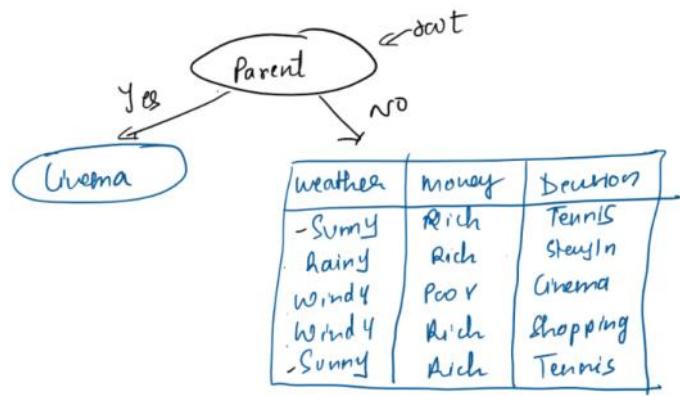
Weighted Average $G_{ini}(\text{Weather}) = \frac{(0.44 * 3/10) + (0.375 * 4/10) + (0.44 * 3/10)}{10} = \underline{\underline{0.414}}$

$G_{ini}(\text{money}) = 0.486$
$G_{ini}(\text{parent}) = 0.36$
$G_{ini}(\text{weather}) = 0.416$

Minimum $G_{ini} \rightarrow$ Minimum Impurity in Decision.

Here Minimum G_{ini} Value = $G_{ini}(\text{parent}) = 0.36$

So the root of Decision is Parent



Q4.B An insurance company wants to predict the likelihood of claims based on customer demographics and policy details. Suggest how a Random Forest model can be used for this task and explain how bagging improves the model's robustness and accuracy. Additionally, compare the performance of the Random Forest model with that of a single decision tree.

Solution >

Using a Random Forest Model in Claim Prediction

→ A Random Forest model is an ensemble method that combines multiple decision trees to make predictions.

Feature Selection

→ Customer demographics (ex: age, gender, location) and policy details (ex: policy type, coverage amount, previous claims)

are treated as input features for the model

→ Label → the likelihood of a claim (ex: claim or no claim, probability of claim) can be target label.

How Bagging Improves Robustness and Accuracy

• Bootstrap → In bagging multiple samples are drawn with replacement from the original dataset to create multiple bootstrapped datasets. Each decision tree in the Random Forest is trained on a different bootstrapped dataset.

• Diversity → Since each tree is trained on a unique sample and may use a random subset of features, the trees become diverse, capturing different patterns in data.

→ Avoided Overfitting → By averaging the results of multiple trees, bagging reduced the likelihood of overfitting which is common problem with individual decision trees.

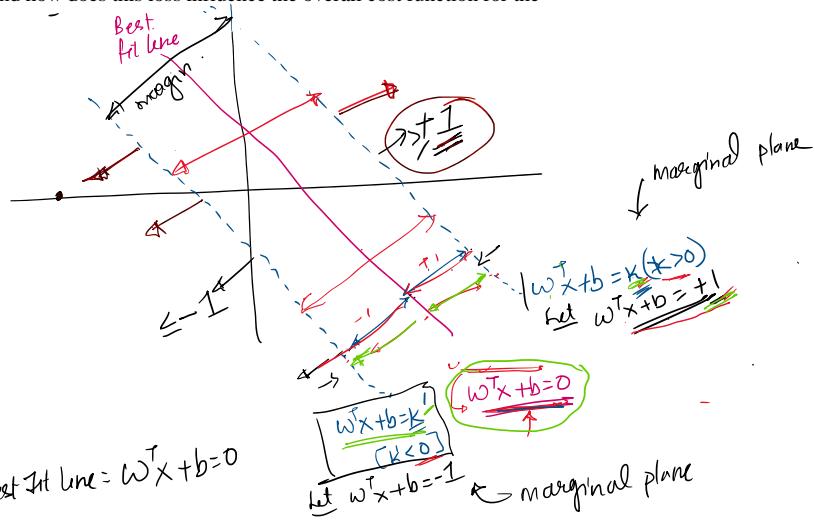
⇒ Increased Robustness → The Random Forest model is more robust to outlier and noise, as the effect of errors in individual trees are minimized by combining results across many trees.

Comparison :

Feature	Random Forest	Single Decision Tree
Accuracy	Higher, due to ensemble effect	lower, prone to overfitting
Variance	Low	High
Interpretability	Lower	Higher
Overshooting	Less likely	More likely
Prediction Stability	High (due to averaging results)	Low (depends on only one model)

Q5.A How do we derive the loss function in a Support Vector Machine (SVM) that penalizes misclassification errors, and how does this loss influence the overall cost function for the classifier?

Solution



Let the Marginal plane be $w^T x + b = +1$ [on +ve side]
And $w^T x + b = -1$ [on -ve side]

Our aim is to draw two marginal planes (+ve & -ve side)
and need to ensure the distance (Margin) is maximum.

* We want to find distance bet'n the Marginal Planes

$$\text{lets find difference: } \begin{aligned} w^T x_1 + b &= +1 \\ w^T x_2 + b &= -1 \\ \hline w^T(x_1 - x_2) &= 2 \end{aligned}$$

$\therefore w^T(x_1 - x_2) = 2$.
Here w = slope (coefficient) magnitude
 direction } Vector.

To convert w^T into Vector ie \vec{w}^T

divide by $\|w\|$

$$\begin{aligned} \therefore \frac{w^T(x_1 - x_2)}{\|w\|} &= \frac{2}{\|w\|} \\ &\quad \uparrow \text{difference} \\ &\quad \uparrow \text{value} \end{aligned} = \boxed{\frac{\vec{w}^T(x_1 - x_2)}{\|w\|} = \frac{2}{\|w\|}}$$

\therefore To Maximize the Margin \rightarrow

$$\boxed{\text{Maximize}_{(w,b)} = \frac{2}{\|w\|}}$$

Constraints \Rightarrow $x_i + 1 \geq w^T x_i + b \geq -1$

(w, b)

Constraints \Rightarrow

$$y_i \begin{cases} +1 & \text{when } w^T x + b \geq 1 \\ -1 & \text{when } w^T x + b \leq -1 \end{cases}$$

(point lies outside of $w^T x + b = \pm 1$)

(point lies below of $w^T x + b = -1$)

The constraints are for Correctly Classified Data point

∴ Final Constraints for S.V Classifier.

$$y_i * (w^T x + b) \geq 1$$

↑ ↑
observed (actual) predicted value

For Correctly Classified Data point

∴ To Maximize $= \frac{2}{\|w\|}$, Subjected to $y_i * (w^T x + b) \geq 1$

Also can be written as

To Minimize $= \frac{\|w\|^2}{2}$, Subjected to $y_i * (w^T x + b) \geq 1$
only for Correctly Classified Data points.

Here we have not considered for Misclassification
But in real world there will be always Misclassification.

∴ To consider for Misclassification we have to use

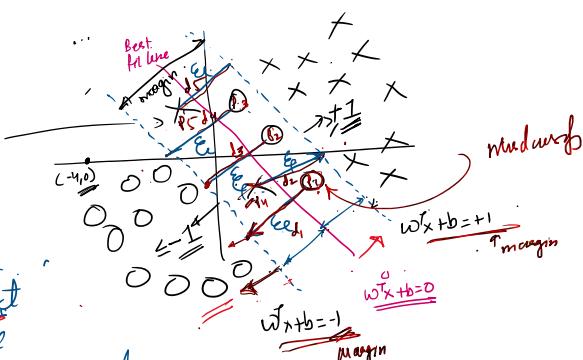
"Hyper parameters"

We allow Misclassification

Now we will use
two Hyperparameters

$\epsilon \Rightarrow$ distance between the misclassified point and correct Marginal plane

$C_i \Rightarrow$ No of Allowed Misclassified Points



Now Final Cost F^n (For All data points) $\left[\begin{array}{l} \text{Correctly Classified} \\ \text{Incorrectly Classified} \end{array} \right]$

↓ Min. Opt.

Now Final Cost \bar{F}^n (For All data points) [Incorrectly classified]

* Objective:

$$\text{To minimize}_{(\omega, b)} \frac{\|\omega\|}{\alpha} + C_1 \sum_{i=1}^{C_0} \epsilon_i \quad \text{constraint } y_i \cdot (\omega^T x_i + b) \geq 1$$

ϵ_i Hyperparameters.
The total Cost \bar{F}^n for SVC used for Classification
[For Allowed Misclassification].

= Sum of ϵ_i values
of all the misclassified points.

↙

accipit

Q5.B Explain SVD as dimensionality reduction technique. Consider the matrix A as given below and find the Eigen values in SVD for ATA.

$$A = \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \end{bmatrix}$$

Solution \rightarrow To find U $\therefore \underline{AA^T} = \begin{bmatrix} 3 & -1 \\ 1 & 3 \end{bmatrix} * \begin{bmatrix} 3 & 1 \\ 1 & 3 \end{bmatrix} = \begin{bmatrix} 11 & 1 \\ 1 & 11 \end{bmatrix}$

$$\Rightarrow AA^T - \lambda I = 0$$

for Eigen Values $\Rightarrow \begin{vmatrix} 11-\lambda & 1 \\ 1 & 11-\lambda \end{vmatrix} = 0 \quad \text{--- (I)}$

$$\lambda^2 - 22\lambda + 120 = 0$$

$$\boxed{\lambda_1 = 10} \quad \boxed{\lambda_2 = 12}$$

For Eigen Vector $\rightarrow \begin{bmatrix} 11-\lambda & 1 \\ 1 & 11-\lambda \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

(case 1) $\boxed{\text{Substitute } \lambda = \lambda_1 = 10}$

$$\begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{x_1}{1} = \frac{-x_2}{1} = 1$$

$$\boxed{x_1 = 1 \quad x_2 = -1}$$

Eigen Vector $x_1 = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$

(and) Substitute $\lambda = \lambda_2 = 10$

$$\begin{bmatrix} 11-\lambda & 1 \\ 1 & 11-\lambda \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\text{Eigen Vector } x_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

* In U the Eigen Vector generated by larger Eigen value will be the first column.

$$\therefore U = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$$

↑ ↑
 Eigen vector Eigen vector
 of $\lambda=12$ of $\lambda=10$

Now we need to Normalize the matrix \Rightarrow divide by length of respective vector.

$$U = \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{bmatrix}$$

↑ ↑
 Length of Length of vector
 vector $\lambda=12$ $\lambda=10$
 $\sqrt{1^2+1^2}$ $\sqrt{1^2+(-1)^2} = \sqrt{2}$
 $= \sqrt{2}$

Slop 2 To find V

$$A^T A = \begin{bmatrix} 3 & -1 & 1 \\ 1 & 3 & 1 \\ 1 & 1 & 3 \end{bmatrix} \begin{bmatrix} 3 & 1 & 1 \\ -1 & 3 & 1 \\ 1 & 1 & 3 \end{bmatrix} = \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 12 \end{bmatrix}_{3 \times 3}$$

there for Eigen values

$$A^T A - \lambda I = 0$$

$$\lambda^3 - 22\lambda^2 + 132\lambda = 0$$

$$\begin{array}{l} \therefore \boxed{\lambda_1 = 0} \\ \lambda_2 = 10 \\ \lambda_3 = 12 \end{array}$$

We need to find Eigen Vectors for the 3 Eigen values.

(one) Eigen Vector for $\lambda_1 = 0$

$$\rightarrow \begin{bmatrix} 10 & 0 & 2 \\ 0 & 10 & 4 \\ 2 & 4 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = 0$$

By Cramers Rule

$$\frac{x_1}{\begin{vmatrix} 10 & 2 \\ 0 & 4 \end{vmatrix}} = \frac{x_2}{\begin{vmatrix} 10 & 0 \\ 0 & 4 \end{vmatrix}} = \frac{x_3}{\begin{vmatrix} 10 & 0 \\ 2 & 10 \end{vmatrix}}$$

$$\frac{x_1}{-20} = -\frac{x_2}{40} = \frac{x_3}{100} = -\frac{1}{20}$$

$$\boxed{x_1 = 1 \quad x_2 = 2 \quad x_3 = -5}$$

$$\text{Eigen Vector } x_1 = \begin{bmatrix} 1 \\ 2 \\ -5 \end{bmatrix}$$

(case 2) Eigen Vector for $\lambda = 10$

$$\begin{bmatrix} 0 & 0 & 2 \\ 0 & 0 & 4 \\ 2 & 4 & -8 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

By Cramers Rule

$$\frac{x_1}{\begin{vmatrix} 0 & 4 \\ 4 & -8 \end{vmatrix}} = \frac{-x_2}{\begin{vmatrix} 0 & 4 \\ 2 & -8 \end{vmatrix}} = \frac{x_3}{\begin{vmatrix} 0 & 0 \\ 2 & 4 \end{vmatrix}} \Rightarrow \frac{x_1}{-16} = \frac{-x_2}{-8} = \frac{x_3}{0} = \frac{-1}{8}$$

$$\therefore x_1 = 2 \quad x_2 = -1 \quad x_3 = 0$$

$$\text{Eigen Vector } x_2 = \begin{bmatrix} 2 \\ -1 \\ 0 \end{bmatrix}$$

(case 3) Eigen Vector for $\lambda_3 = 12$

$$\begin{bmatrix} -2 & 0 & 2 \\ 0 & -2 & 4 \\ 2 & 4 & -10 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

By Cramers Rule

$$\frac{x_1}{\begin{vmatrix} 0 & 2 \\ -2 & 4 \end{vmatrix}} = \frac{-x_2}{\begin{vmatrix} -2 & 2 \\ 0 & 4 \end{vmatrix}} = \frac{x_3}{\begin{vmatrix} -2 & 0 \\ 0 & -2 \end{vmatrix}} \Rightarrow \frac{x_1}{4} = \frac{-x_2}{-8} = \frac{x_3}{4} = \frac{1}{4}$$

$$\therefore x_1 = 1 \quad x_2 = 2 \quad x_3 = 1$$

Eigen Vector $X_3 = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix}$

$$\text{Now } V = \begin{bmatrix} 1 & 2 & 1 \\ 2 & -1 & 2 \\ 1 & 0 & 5 \end{bmatrix}$$

↑ Eigen Vector for $\lambda=12$ ↑ Eigen Vector for $\lambda=10$ ↑ Eigen Vector for $\lambda=0$.

Now Normalizing the matrix \rightarrow Divide by length of vector.

$$V = \begin{bmatrix} \sqrt{56} & 2/\sqrt{56} & 1/\sqrt{30} \\ 2/\sqrt{56} & -1/\sqrt{56} & 2/\sqrt{30} \\ 1/\sqrt{56} & 0 & -5/\sqrt{30} \end{bmatrix} //$$

↑ length of vector $\sqrt{56}$ ↑ length of vector $\sqrt{10}$ ↑ length of vector $\sqrt{30}$

$$V^T = \begin{bmatrix} 1/\sqrt{56} & 2/\sqrt{56} & 1/\sqrt{56} \\ 2/\sqrt{56} & -1/\sqrt{56} & 0 \\ 1/\sqrt{30} & 2/\sqrt{30} & -5/\sqrt{30} \end{bmatrix} //$$

Step 3 To find Σ (or D)

$$\Sigma = \begin{vmatrix} \sqrt{12} & 0 & 0 \\ 0 & \sqrt{10} & 0 \\ 0 & 0 & \sqrt{0} \end{vmatrix}_{3 \times 3}$$

↑ diag matrix \Rightarrow The diag elements are Eigen values in Square root of decreasing order.

$$\therefore A = U \Sigma V^T$$

$$U = \begin{bmatrix} \sqrt{\frac{1}{2}} & \sqrt{\frac{1}{2}} \\ \sqrt{\frac{1}{2}} & -\sqrt{\frac{1}{2}} \end{bmatrix}, \Sigma = \begin{bmatrix} 5 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 0 \end{bmatrix}, V = \begin{bmatrix} \sqrt{\frac{1}{3}} & \sqrt{\frac{1}{3}} & \sqrt{\frac{1}{3}} \\ \sqrt{\frac{1}{3}} & -\sqrt{\frac{1}{3}} & 0 \\ \sqrt{\frac{1}{3}} & 0 & -\sqrt{\frac{1}{3}} \end{bmatrix}$$