

NAIVE BAYESIAN PROBLEM

Consider the training data in the following table where play is a class attribute. In the table, the humidity attribute has values "L" (for Low) or "H" (for High), sunny has values "N" (for No) or "Y" (for Yes) wind has values "S" (for Strong) or "W" (for weak) and play has "Yes" or "No"

Humidity	Sunny	Wind	play
L	N	S	No
H	N	W	Yes
H	Y	S	Yes
H	N	W	Yes
L	Y	S	No

What is the class label for the following day (Humidity=L, Sunny=N, Wind=S)

GIVEN, No. of datasets = 5

Evidence (Humidity=L, Sunny=N, Wind=S)

$$P(\text{Play} = \text{"Yes"}) = 3/5 \quad P(\text{Play} = \text{"Yes"}, \text{Humidity} = \text{"L"}) = 0$$

$$P(\text{play} = \text{"No"}) = 2/5 \quad P(\text{play} = \text{"No"}, \text{Humidity} = \text{"L"}) = 2/2 =$$

$$P(\text{play} = \text{"Yes"}, \text{Sunny} = \text{"N"}) = 2/3 = 0.66$$

$$P(\text{play} = \text{"No"}, \text{Sunny} = \text{"N"}) = 1/2 = 0.5$$

$$P(\text{play} = \text{"Yes"}, \text{Wind} = \text{"S"}) = 1/3 = 0.33$$

$$P(\text{play} = \text{"No"}, \text{Wind} = \text{"S"}) = 2/2 = 1$$

GIVEN Test set:-

$$P(\text{Yes} / \text{Hum=L, Sunny=N, Wind=S}) = 0 \times 0.66 \times 0.33 \\ \leq 0.$$

$$P(\text{No} / \text{Hum=L, Sunny=N, Wind=S}) = 1 \times 0.5 \times 1 \\ = 0.5$$

Since $P(\text{No}) > P(\text{Yes})$ play = "No"

DECISION TREE CLASSIFIER.

tree like predictive model used for classification and regression.

Supervised learning model.

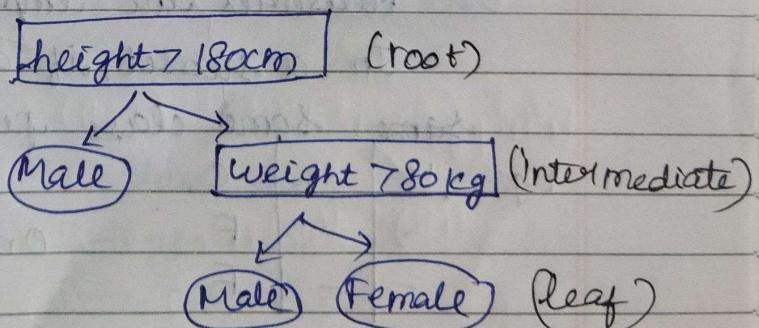
ALGORITHMS USED:

CART, ID3

Regression Classification

Height (cm)	Weight (kg)	Class
170	75	Female
170	90	Male
185	85	Male
155	50	Female
175	70	Female

STRUCTURE:



In this DT classifier class label is Male / Female based on height and weight.

H | w | class

> 180 < 180

H	w	C
185	85	Male

H	w	C
170	75	F
170	90	M
155	50	F
175	70	F

H	w	C
170	90	M

H	w	C
170	75	F
155	50	F
175	78	F

At each node of dt construction choose a feature as a splitting attribute.

Random attribute \Rightarrow Not a good idea less accuracy.

Therefore we need to find best splitting attribute.

FEATURE SELECTION:

Information gain

Gini index

INFORMATION GAIN:

Identify entropy (finding impurities)

ENTROPY (S):

Measures the impurity in the dataset

In a subset, if all the samples are in

same some class pure node.

H	w	C
-	-	F
-	-	F

H	w	C
-	-	M
-	-	M

All the samples in the above are in the same class.

Eg: H W C

170 75 F

170 90 M

155 50 F

175 78 F

NO. of samples = 4

2 classes - F, M

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

$$= -\frac{1}{4} \log_2 \left(-\frac{1}{4} \right) - \frac{3}{4} \log_2 \left(\frac{3}{4} \right)$$

$$= -\frac{1}{4} \times (-2) - \frac{3}{4} \times (-0.415)$$

$$= 0.81125$$

If the samples are in equal class then
subset is impure.

Entropy of pure node = 0

Impurity is high \Rightarrow entropy = 1.

Value of entropy range 0-1.

Dataset contains 2 class then the
entropy will be 0 or 1.

Dataset contains "n" classes the entropy
is $0 - \log_2 n$.

$$\text{Entropy} = \sum_{i=1}^c -p_i \log_2 p_i$$

C - Class label

$p_i \rightarrow$ probability of ith class

S \rightarrow subset.

Eg:-

Day	Weather	Temp	Windy	play
1.	Rainy	Mild	Weak	No
2.	Normal	Hot	Weak	Yes
3.	Wind	Mild	Strong	Yes
4.	Normal	Cool	weak	No
5.	Rain	hot	Strong	No

$$\text{Entropy} = \sum_{i=0}^C -p_i \log_2 p_i$$

$$\text{Entropy}[+2, -3] = \frac{-2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5}$$

$$\begin{aligned} \log_2 \frac{2}{5} &= \frac{\log 2/5}{\log 2} \\ 2 \text{ Yes} &= \frac{-2}{5} (-1.32) + \frac{3}{5} (0.736) \\ 3 \text{ No} &= \frac{2.64 + 2.2}{5} = \frac{9.84}{5} = 0.96 \end{aligned}$$

$$\text{Information Gain } G_{\text{rain}} = \text{Entropy}(S) - \frac{1}{|S|} \sum_{i=1}^{|S|} \text{Entropy}(S_i)$$

$$G_{\text{rain}}(\text{weather}) = \{ \text{Rainy, Normal, Windy} \}$$

$$\begin{aligned} \text{Entropy}(\text{Rain}) &\Rightarrow \text{Entropy}(0, -2) = \frac{0}{2} \log_2 \frac{0}{2} - \frac{2}{2} \log_2 \frac{2}{2} \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Normal}) &\Rightarrow (1, 1) = \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{Wind}) &\Rightarrow (1, 0) = \frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} \\ &= 0 \end{aligned}$$

$$\text{Gain(Weather)} = \text{Entropy} - \left(\frac{2}{5} \times 0 \right) - \left(\frac{2}{5} \times 1 \right) - \left(\frac{1}{5} \times 0 \right)$$

$$= 0.568$$

$\text{Gain(Temp)} = \{\text{Mild, Hot, Cool}\}$

$$\text{Entropy(Mild)} = (1, -1) = \frac{-1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 1$$

$$\text{Entropy(Hot)} = (1, -1) = \frac{-1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 1$$

$$\text{Entropy(Cool)} = (0, -1) = \frac{-0}{1} \log_2 \frac{0}{1} - \frac{1}{1} \log_2 \frac{1}{1}$$

$$= 0$$

$$\text{Gain(Temp)} = 0.97 - \left(\frac{2}{5} \times 1 \right) - \left(\frac{2}{5} \times 1 \right) - \left(\frac{1}{5} \times 0 \right) = 0.17$$

$\text{Gain(windy)} = \{\text{weak, strong}\}$

$$\text{Entropy(weak)} = (1, -2) = \frac{-1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.5283 + 0.3899$$

$$= 0.9182$$

$$\text{Entropy(Strong)} = (1, -1) = \frac{-1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 1$$

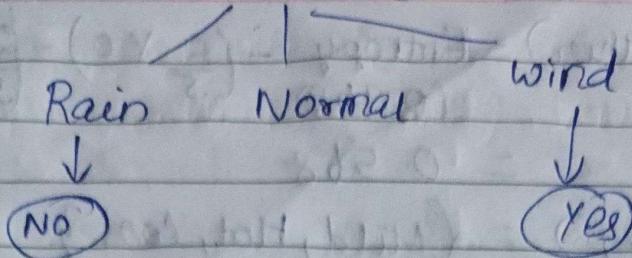
$$\text{Gain(windy)} = 0.97 - \left(\frac{3}{5} \times 0.9182 \right) - \left(\frac{2}{5} \times 1 \right)$$

$$= 0.97 - 0.55 - 0.4$$

$$= 0.02$$

Weather

DATE: / /



Day	Temp	wind	play
2	HOT	weak	Yes
4	Cool	weak	No

$$\text{Gain(Temp)} = \{ \text{HOT}, \text{COOL} \}$$

$$\text{Entropy(Hot)} = (1, 0) = \frac{-1}{2} \log_2 \frac{1}{1} - \frac{0}{1} \log_2 \frac{0}{1} = 0$$

$$\text{Entropy(Cool)} = (0, -1) = \frac{-0}{1} \log_2 \frac{0}{1} - \frac{1}{1} \log_2 \frac{1}{1} = 0$$

$$\text{Gain(Temp)} = 0.97 - \frac{1}{2}(0) - \frac{1}{2}(0)$$

RULE BASED ALGORITHM

Rule based classifier are just another type of classifier, which makes the class decision depending by various "IF ELSE" rules.

IMPORTANT KEY WORDS:

If, And, then

SET OF IF THEN RULES FOR CLASSIFICATION:

If Condition then conclusion.

"If" part of the rule is called rule Antecedent precondition.

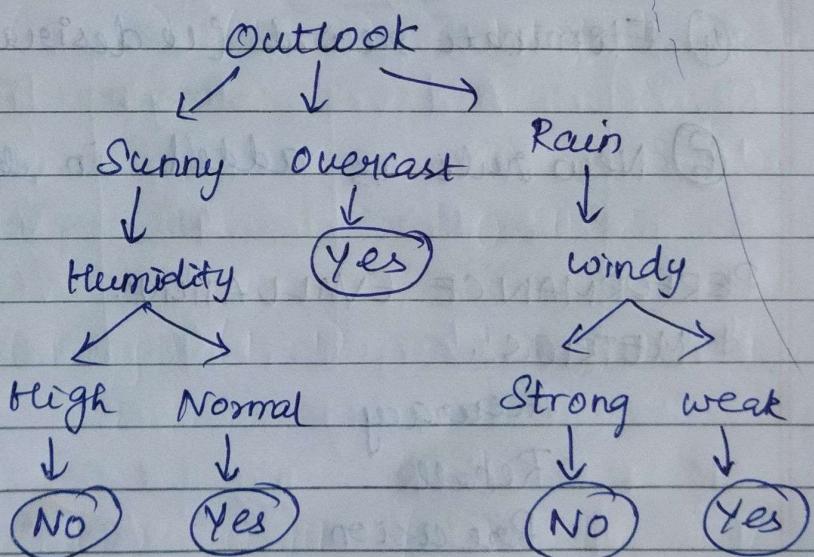
"Then" part of the rule is called rule consequent

Consequent part consists of class prediction.
↳ Pre-condition

e.g.: If Outlook = "Sunny" AND Humidity = "High"
then play = "No" → Rule consequent

To extract the rule from decision tree
One rule is created for each path
from the root to the leaf node.
Each splitting condition is logically
AND

The leaf node holds the class prediction,
forming the rule consequent.

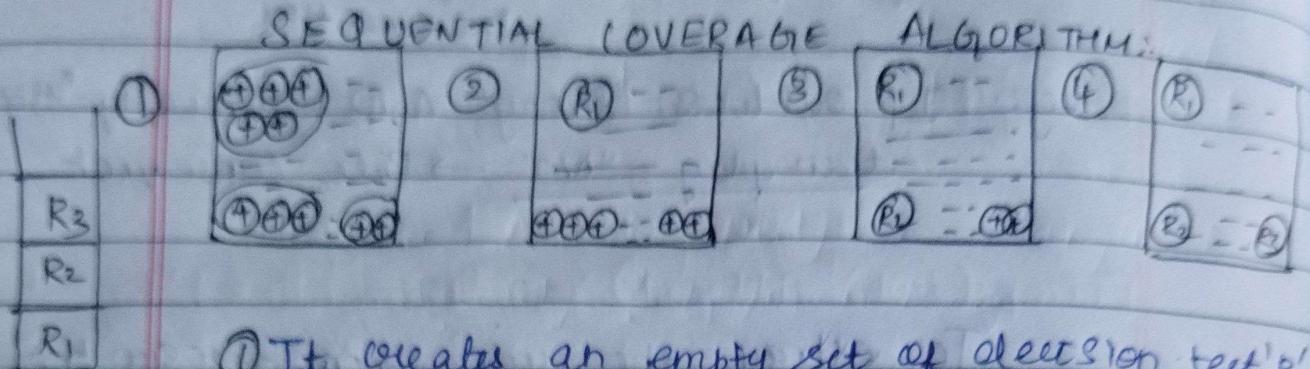


If Outlook = "Sunny" AND Humidity = "High"
then PLAY = "No"

If Outlook = "Sunny" AND Humidity = "Normal"
then Play = "Yes"

If Outlook = "Overcast", then play = "Yes"
 If Outlook = "Rain" AND Windy = "Strong"

RULE BASED



- ① It creates an empty set of decision test 'P'
- ② A function called "Learn of rule" function is used. It extract the best rule for class "y"
 - If all the training records & class y =>
 - If all the training records & class y =>
- ③ Get desirable value (Only true)
- ④ Eliminate records (ie desirable one)
- ⑤ New rule is added in the bottom of the P

PERFORMANCE EVALUATION:

METRICS:

Accuracy

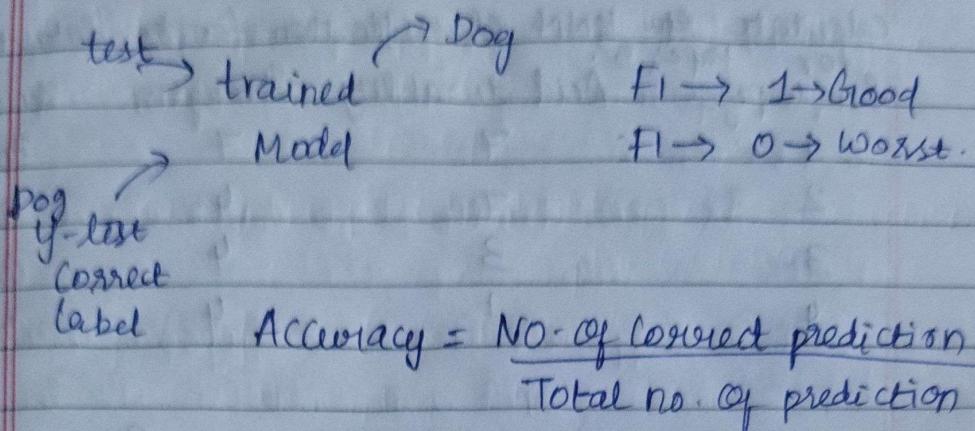
Recall

Precision

F1 Score

Eg: Binary Images.

100	Cat v. Not train	30% test
Image	dog	

TEST IMAGE:

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

$$\text{Recall} = \frac{\text{All the relevant datapoint}}{\text{Total datapoint}} = \frac{TP}{TP + FN}$$

$$\text{Precision} \Rightarrow \frac{\text{only the relevant datapoint}}{\text{Total datapoint}} = \frac{TP}{TP + FP}$$

$$\text{F1 Score} \Rightarrow \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

CONFUSION MATRIX:

N=165		Predicted	
		No	Yes
Actual No	TN=50	FP=10	60
	Actual Yes FN=5	TP=100	105
	55	110	

$$\text{Accuracy} = \frac{150}{165} = 0.90$$

$$\text{F1 Score} = \frac{2 \times 0.85}{1.85}$$

$$\text{Recall} = \frac{100}{105} = 0.9$$

$$= \frac{1.7}{1.85} = 0.91$$

$$\text{Precision} = \frac{100}{110} = 0.9$$

K - NEAREST NEIGHBOR

calculate the KNN classification for the dataset of Predict for Paper 5.

Sample paper	Acid durability	Strength	Quality
1	7	7	Bad
2	7	4	Bad
3	3	4	Good
4	1	4	Good

5th Data: Acid durability = 3 Strength = 7
Quality = ? given k=3

Identify k.

STEP-1: k = 3

Calculate the distance b/w the query instance and all the training samples

∴ Given: instance (3, 7) & calculate using Euclidean distance

$$d(x_i, y_i) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

SP	A.d	Strength	Distance
1	7	7	$\sqrt{(3-7)^2 + (7-7)^2} = 4$
2	7	4	$\sqrt{(3-7)^2 + (7-4)^2} = 5$
3	3	4	$\sqrt{(3-3)^2 + (7-4)^2} = 3$
4	1	4	$\sqrt{(3-1)^2 + (7-4)^2} = \sqrt{13} = 3$

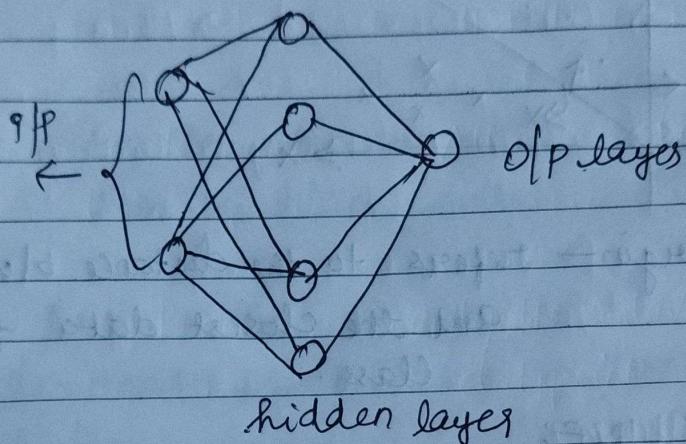
SP ED Rank Included in 3 neighbor

1	4	3	Yes
2	5	4	No
3	3	1	Yes
4	3.6	2	Yes

Dataset: 1, 3, 4

SP	A:D	Strength	Quality
1	2	7	Bad
3	3	4	Good
4	1	4	Good
5	3	7	Good

CLASSIFICATION BY BACK PROPAGATION:-



FORWARD PASS:-

- ① Σ
- ② Activation function.

BACKWARD PASS:-

- ① Error comparison
- ② weight updation
- ③ Bias updation.

SVM CLASSIFICATION:-

SVC:-

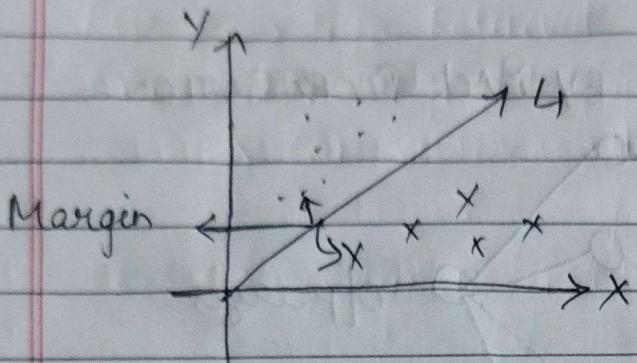
classification method for both linear & non-linear data.

An SVM classifies the data by

finding the best hyperplane that separates all data points of one class from those of the other class

It can handle high dimensional data.

Support Vector → datapoint in a dataset that are closest to the hyperplane



Margin - refers to the distance b/w hyperplane and the closest data from each class.

HYPER LINE:

Decision boundary that separates a data points into different classes.

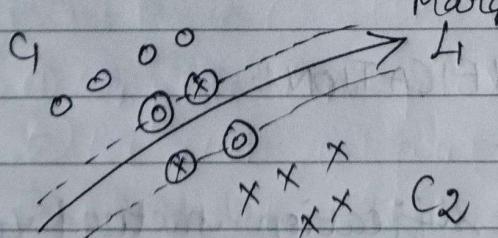
MAIN GOAL OF SVM:

to find the hyperline that maximise the margin b/w 2 classes.

MAXIMUM MARGIN CLASSIFIER:

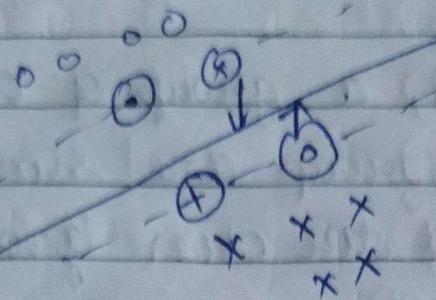
uses hard margin

Margin



MMC will not allow misclassification. It won't allow for & whether it may belong to C2

SVM allows misclassification datapoints
uses soft margin



MARGIN: Cross Validation.

Used to determine the best soft margin.

Maximum Margin

Min-Mis classification.

LAZY

→ instance based learning → Given a set of training
it is,
→ less time in training
more time in predicting

Eg:- kNN, lazy NB

EAGER

people, construct a
classification model
& receives a new data
to classify

Eg: Decision Tree, SVM,
ANN

MODULE - IV

CLUSTERING

CLUSTERING:-

In general, a group of objects such that objects in the group are similar and dissimilar with objects in the other groups.

- APP:-
 - Info. Retrieval
 - Earthquake prediction
 - Can also used → Data preprocessing for data summarisation.
 - Dimensionality reduction.

TYPES:- * partitional clustering
* Hierarchical "

PARTITIONAL:-

a division data objects subset (clusters) such that each data object is in exactly one subset

Algorithms: KNN, CLARA, PAM

HIERARCHICAL:-

a set of nested cluster organised as a hierarchical tree.

STEPS:-

- ① Complete the proximity Matrix
- ② let each datapoint be a cluster
- ③ Merge the two closest cluster
- ④ update the proximity Matrix until only one cluster is defined-

TYPES:-

- ① Agglomerative cluster
- ② Division cluster.

KNN CLUSTER ANALYSIS:

Consider the dataset (D) given below that records the programming skill rating of students in the scale of 0-10. Use k-Means algorithm and euclidean distance measure to cluster the students into two clusters.

$$C(K) = 2 \quad \text{Initial: } D_1(1,3), D_2(3,2)$$

Student Id	Java	Python	
D_i h=2 data object	1	3	similarity
	3	2	Measurin-
	6	5	ED
	7	9	
	6	8	

STEP-1: Find the k and fix it. $k=2$.

STEP-2: Randomly choose 2 data points from D and assign it as cluster center or Centroids.

Here, D_1 & D_2
 $(1,3)$ $(3,2)$
 cluster-1 cluster-2.

STEP-3: Calculate the Distance $\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$
 b/w datapoints & cluster center.

Object	Data point	cluster center-1	CC 3 (3,2)	Assign Center C1
D ₁	(1,3)	0	$\sqrt{5}$	C ₁
D ₂	(3,2)	$\sqrt{5}$	$\sqrt{10}$	C ₂
D ₃	(6,5)	$\sqrt{29}$	$\sqrt{18}$	C ₂
D ₄	(7,9)	$\sqrt{72}$	$\sqrt{65}$	C ₂
D ₅	(6,8)	$\sqrt{50}$	$\sqrt{45}$	C ₂

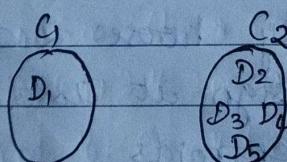
$$(6,5) (1,3) = \sqrt{(1-6)^2 + (3-5)^2} \\ = \sqrt{25+4} = \sqrt{29}$$

$$(6,5) (3,2) = \sqrt{(3-6)^2 + (2-5)^2} \\ = \sqrt{9+9} = \sqrt{18}$$

$$(7,9) (1,3) = \sqrt{(1-7)^2 + (3-9)^2} \\ = \sqrt{36+36} = \sqrt{72}$$

$$(7,9) (3,2) = \sqrt{(3-7)^2 + (2-9)^2} \\ = \sqrt{16+49} = \sqrt{65}$$

STEP-4: Find the cluster center by calculating mean of each cluster.



$$\text{cluster Mean } C_1 = (1,3)$$

$$\text{cluster Mean } C_2 = (5.5, 6) \quad \frac{3+6+9+6}{4}, \frac{2+5+9+8}{4}$$

Iteration-2:

Object	Data point	C ₁ (1,3)	C ₂ (5.5, 6)	cluster
D ₁	(1,3)	0	$\sqrt{29.25}$	C ₁
D ₂	(3,2)	$\sqrt{5}$	$\sqrt{22.5}$	C ₁
D ₃	(6,5)	$\sqrt{29}$	$\sqrt{1.25}$	C ₂
D ₄	(7,9)	$\sqrt{72}$	$\sqrt{15.25}$	C ₂
D ₅	(6,8)	$\sqrt{50}$	$\sqrt{6.25}$	C ₂

$$(1, 3) \quad (5.5, 6) = \sqrt{(5.5-1)^2 + (6-3)^2}$$

$$= \sqrt{(4.5)^2 + 9} = \frac{5.5}{\cancel{10}}$$

$$(7, 9) \quad (5.5, 6) = \sqrt{(5.5-7)^2 + (6-9)^2} = \frac{45 \times 45}{225}$$

$$= \sqrt{15.25} = \frac{180x}{20.25}$$

c_1
(1, 3)
D₁ D₂

c_2
(5.5, 6)
D₃, D₄
D₅

$$\begin{aligned} D_1 &= (1, 3) & D_3 &= (6, 5) \\ D_2 &= (3, 2) & D_4 &= (7, 9) \\ \text{Mean} &= \frac{1+3}{2}, \frac{3+2}{2} & D_5 &= (6, 8) \\ &= (2, 2.5) & & \frac{6+7+6}{3}, \frac{5+9+8}{3} \\ & & & = \frac{29}{3}, \frac{22}{3} \\ & & & = (6.33, 7.33) \end{aligned}$$

Iteration 3:

Object	Datapoint	$c_1(2, 2.5)$	$c_2(6.33, 7.33)$	cluster
D ₁	(1, 3)	$\sqrt{3.25}$	$\sqrt{46.58}$	C ₁
D ₂	(3, 2)	$\sqrt{1.25}$	$\sqrt{38.98}$	C ₁
D ₃	(6, 5)	$\sqrt{22.25}$	$\sqrt{5.38}$	C ₂
D ₄	(7, 9)	$\sqrt{22.25}$	$\sqrt{5.78}$	C ₂
D ₅	(6, 8)	$\sqrt{46.25}$	$\sqrt{0.79}$	C ₂

$$(1, 3) \quad (6.33, 7.33) = \sqrt{(6.33-1)^2 + (7.33-3)^2}$$

$$= \sqrt{(5.33)^2 + (4.33)^2}$$

$$= \sqrt{28.4 + 18.74}$$

$$= \sqrt{47.14} = \sqrt{46.58}$$

$$(6, 8) \quad (6.33, 7.33) = \sqrt{(6-6.33)^2 + (8-7.33)^2}$$

$$= \sqrt{(0.3)^2 + (0.7)^2}$$

$$= \sqrt{0.09 + 0.49}$$

$$= \sqrt{0.58}$$

Apply the DB Scan algorithm to the given data points and create the cluster with min points = 4 and $\epsilon = 1.9$ (Epsilon)

Data points:

$P_1: (3, 7)$ $P_2: (4, 6)$ $P_3: (5, 5)$
 $P_4: (6, 4)$ $P_5: (7, 3)$ $P_6: (6, 2)$
 $P_7: (7, 2)$ $P_8: (8, 4)$ $P_9: (3, 3)$
 $P_{10}: (2, 6)$ $P_{11}: (3, 5)$ $P_{12}: (2, 4)$

Min point: 4 $\epsilon = 1.9$

	P_1	P_2	P_3	P_4	P_5	P_6	P_7	P_8	P_9	P_{10}	P_{11}	P_{12}
$(3, 7) P_1$	0	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41
$(4, 6) P_2$	1.41	0	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41
$(5, 5) P_3$	1.41	1.41	0	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41
$(6, 4) P_4$	1.41	1.41	1.41	0	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41
$(7, 3) P_5$	1.41	1.41	1.41	1.41	0	1.41	1.41	1.41	1.41	1.41	1.41	1.41
$(6, 2) P_6$	1.41	1.41	1.41	1.41	1.41	0	1.41	1.41	1.41	1.41	1.41	1.41
$(7, 2) P_7$	1.41	1.41	1.41	1.41	1.41	1.41	0	1.41	1.41	1.41	1.41	1.41
$(8, 4) P_8$	1.41	1.41	1.41	1.41	1.41	1.41	1.41	0	1.41	1.41	1.41	1.41
$(3, 3) P_9$	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	0	1.41	1.41	1.41
$(2, 6) P_{10}$	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	0	1.41	1.41
$(3, 5) P_{11}$	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	0	1.41
$(2, 4) P_{12}$	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	1.41	0

$$P_1 \rightarrow P_2, P_0 \quad P_5 \rightarrow P_4, P_6, P_8, P_7$$

$$P_2 \rightarrow P_3, P_1, P_{11} \quad P_6 \rightarrow P_5, P_7$$

$$P_3 \rightarrow P_2, P_4 \quad P_7 \rightarrow P_5, P_6$$

$$P_4 \rightarrow P_3, P_5 \quad P_8 \rightarrow P_5$$

$$P_9 \rightarrow P_{12}$$

$$P_{10} \rightarrow P_1, P_{11}$$

$$\begin{aligned} P_{11} &\rightarrow P_2, P_{10}, P_{12} \\ P_{12} &\rightarrow P_9, P_{11} \end{aligned}$$

Points	Noise/Core	Border / Not
P ₁	Noise	Border
P ₂	Core	Border
P ₃	Noise	Border
P ₄	Noise	Border
P ₅	Core	Border
P ₆	Noise	Border
P ₇	Noise	Border
P ₈	Noise	Border
P ₉	Noise	Border
P ₁₀	Noise	Border
P ₁₁	Core	Border
P ₁₂	Noise	Border

