

Machine Learning.

Unit :- 2

Date _____
Page _____

* Types of descriptive statistics :-

- 1) distribution
- 2) central tendency
- 3) variability

* Measures of central tendency :-

↳ mean

↳ median

↳ mode

Mean :-

Dataset : 15, 3, 12, 0, 24, 3.

$$\begin{aligned} \text{sum} &= 15 + 3 + 12 + 0 + 24 + 3 \\ &= 57 \end{aligned}$$

Total no. of responses $N = 6$.

∴ Mean = $\frac{57}{6} = \boxed{9.5}$

* Measures of variability

↳ Range

↳ Standard deviation

↳ Variance.

Range :- Ordered dataset : 0, 3, 3, 12, 15, 24

$$\begin{aligned} \text{Range} &= 24 - 0 \\ &= 24. \end{aligned}$$

Standard deviation :- [S or SD]

It is the avg. amount of variability in the dataset.

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{N}}$$

Date _____
Page _____

Ex: Given:-

Rawdata: 15, 3, 12, 0, 24, 3.

Sd n.

Rawdata	Deviation from Mean	Squared deviation
15	$15 - 9.5 = 5.5$	30.25
3	$3 - 9.5 = -6.5$	42.25
12	$12 - 9.5 = 2.5$	6.25
0	$0 - 9.5 = -9.5$	90.25
24	$24 - 9.5 = 14.5$	210.25
3	$3 - 9.5 = -6.4$	42.25
<hr/>		421.5
$\bar{x} = \frac{9.5}{6}$		

Mean.

$$\rightarrow \text{Mean} = \frac{15 + 3 + 12 + 0 + 24 + 3}{6} = 9.5$$

$$\begin{aligned} \rightarrow S.D &= \sqrt{\frac{\text{sum of squared deviation}}{N-1}} \\ &= \sqrt{\frac{421.5}{5}} \\ &= \sqrt{84.3} \end{aligned}$$

$$\therefore S = 9.18$$

Variance : Square of standard deviation

$$\Leftrightarrow S = 9.18$$

$$\therefore S^2 = 84.3 \leftarrow \text{Variance of given dataset}$$

* Types of descriptive statistics :-

- 1) Univariate 2) Bivariate 3) Multivariate
- focuses on a one variable at a time → simultaneously study frequency & variability of two variables.
- concerns bivariate but with more than two variables.

* Mode :-

→ The value that has higher frequency.

For Ungrouped data :-

Match no.	1	2	3	4	5	6	7	8	9	10
No. of wickets	2	1	1	3	2	3	2	2	4	1

2 wickets were taken by the batter frequently.

∴ mode is 2

For grouped data :-

$$\text{Mode} = l + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h.$$

l = lower limit of the modal class

h = size of class interval

f_1 = frequency of modal class

f_0 = freq. of class preceding the modal class.

f_2 = freq. of class succeeding the modal class.

Ex:- In a class of 30 students marks obtained by students in mathematics out of 50 is tabulated as below.
Calculate the mode of given data.

Marks obtained	Number of students	
10 - 20	5 $\leftarrow f_0$	$h = 30 - 20$
$l \rightarrow 20 - 30$	12 $\leftarrow f_1$	$= 10$
30 - 40	8 $\leftarrow f_2$	
40 - 50	5.	

\Rightarrow Maximum class freq. = 12

\therefore corresponding modal class = 20 - 30

$$\therefore \text{Mode} = l + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h$$

$\therefore l$ = lower limit of modal class = 20

$\therefore h$ = size of class interval = $30 - 20$
 $= 10$

$\therefore f_1$ = freq. of modal class = f_{12}

$\therefore f_0$ = freq. of the class preceding
the modal class = 5

$\therefore f_2$ = freq. of the class succeeding
the modal class = 8

$$\begin{aligned}
 \text{Mode} &= l + \left(\frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times h \\
 &= 20 + \left(\frac{12 - 5}{24 - 5 - 8} \right) \times 10 \\
 &= 20 + \boxed{26.364}
 \end{aligned}$$

* Median :

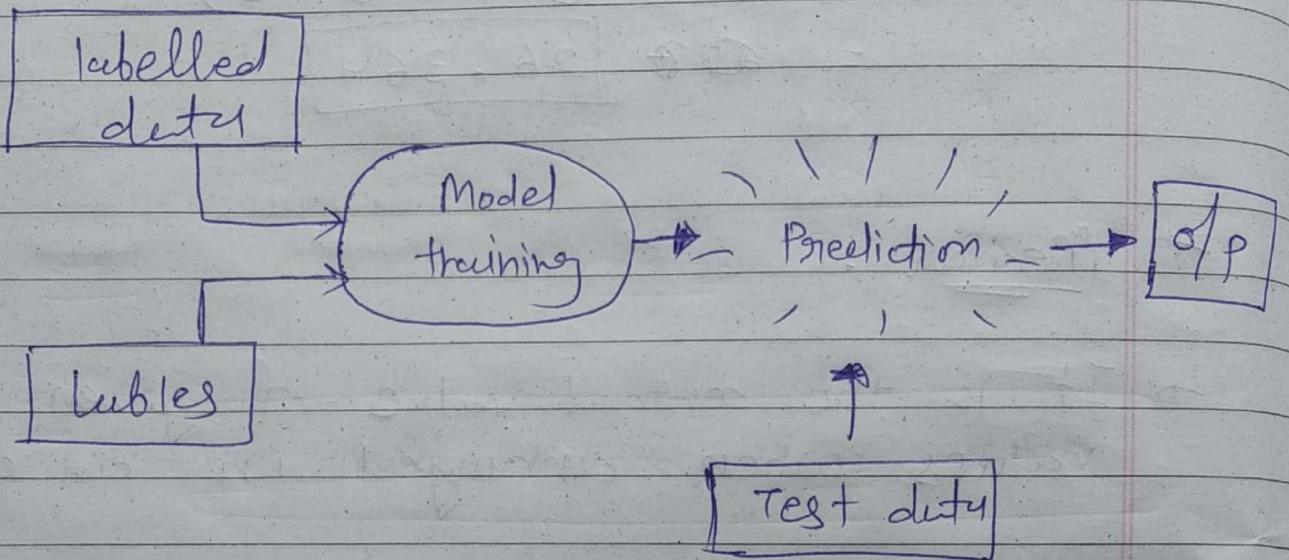
⇒ It is the central value of given set of values when arranged in an order.

* Formula :

$$\text{Mode} = 3 \text{ median} - 2 \text{ mean}$$

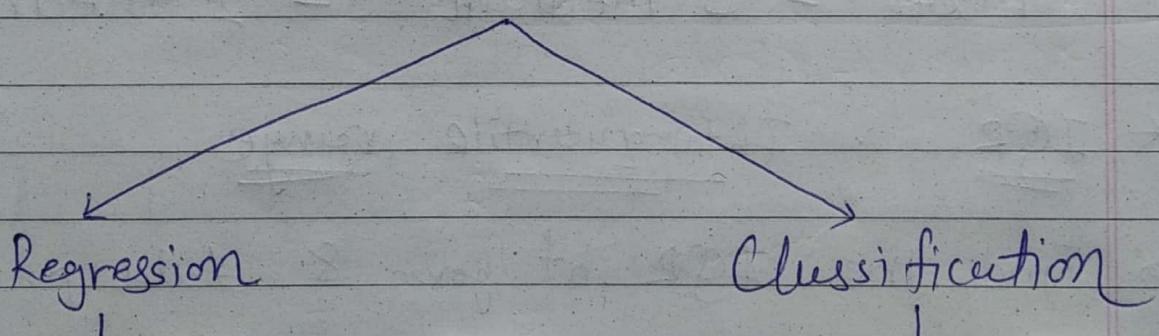
Supervised learning

* Diagram :-



* Types of supervised ML Algorithms :-

knowledge
purpose



Linear Regression

Bayesian linear regression

Regression Trees

Non-linear regression

Polynomial regression.

Random forest

Decision Trees

Logistic regression

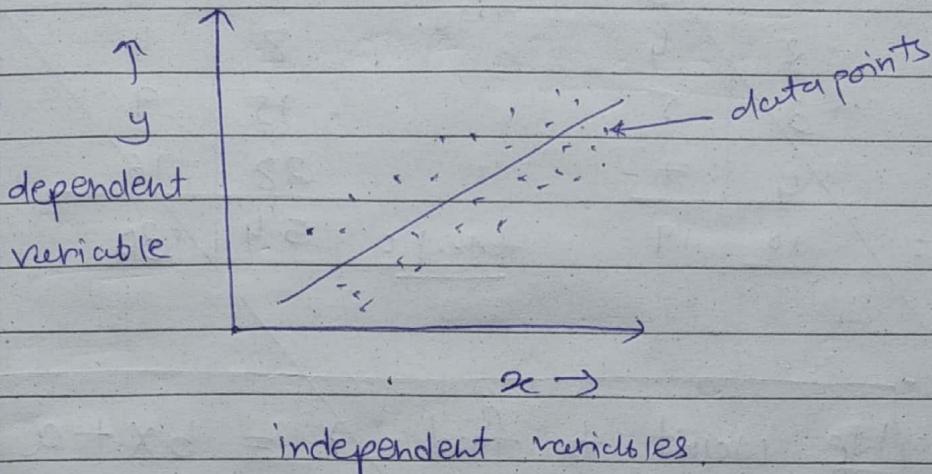
Support vector machine

Naive Bayes Classifier

K-Nearest Neighbours.

* Linear Regression :-

→ Dependent variable is continuous in nature.



1) Simple linear regression :

→ one dependent and only one independent variable exists.

$$y = \alpha_0 + \alpha_1 x_1$$

↑ ↑
dependent variable independent variable
(to be found) co-efficient of regression.

2) Multiple linear regression :-

→ more than one independent variable and only one dependent variable exists.

$$y = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_m x_m$$

α_i = Regression co-efficient

x_i = Independent variable

y = Dependent variable.

$$* b = \bar{y} - a\bar{x}$$

* Ex: Number of TR Ads Number of cars sold

Linear Regression	(x)	(y)
	1	14
	3	24
	2	18
	1	17
	3	27
	10	100

Date _____
Page _____

OLS method
Ordinary least square method.

$$\rightarrow \bar{x} = \frac{10}{5} = 2 \quad | \quad \bar{y} = \frac{100}{5} = 20$$

$$\rightarrow \text{Regression eqn: } \boxed{y = ax + b}$$

$$\checkmark \quad a = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{slope}$$

$$\checkmark \quad b = \bar{y} - a\bar{x} \quad \text{intercept}$$

x	y	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$	$(x_i - \bar{x})^2$
1	14	-1	-6	6	1
3	24	1	4	4	1
2	18	0	-2	0	0
1	17	-1	-3	3	1
3	27	1	7	7	1
10	100			20	4

$$\therefore a = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$
$$= \frac{20}{4}$$

$$\therefore \boxed{a = 5}$$

and $b = \bar{y} - a \bar{x}$

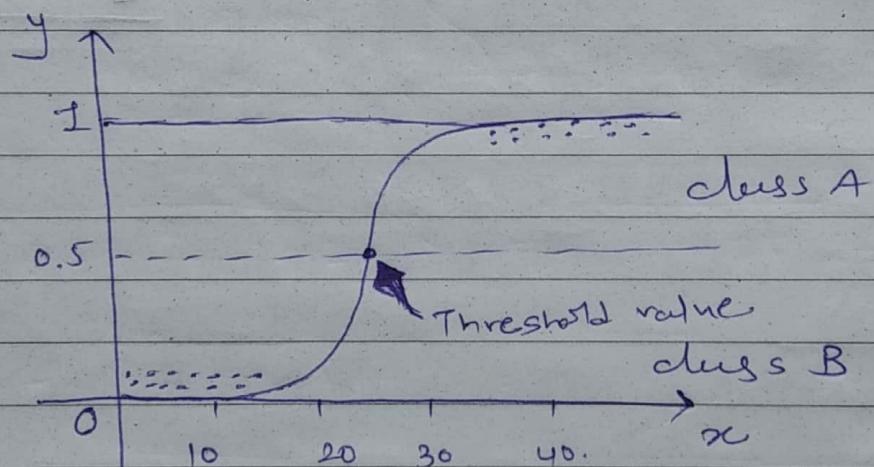
$$= 20 - (5)(2)$$
$$\therefore \boxed{b = 10}$$

$$\therefore \text{Estimated regression eqn} \Rightarrow \hat{y} = aX + b$$
$$\hat{y} = 5X + 10$$

* Logistic Regression :-

- It gives probabilistic values which lie between 0 and 1.
- Used to solve classification problems.
- sigmoid function $P = \frac{e}{1 + e^{-x}}$; $e = 2.718$

Sigmoid function is simply trying to convert the independent variable (x) into expression of probability that ranges b/w 0 and 1 w.r.t. the dependent variable.



- No Dataset should be free of missing values.
- Data points may vary between 30 - 50 for one class. ∵ 60 - 100 datapoints. Should be there
- Eqⁿ: $\log \left[\frac{P}{1-P} \right] = \alpha_0 + \alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_m x_m$

* Types of logistic regression :-

Binomial : two possible types of dependent variable.

→ Yes or No, 0 or 1, Pass or Fail

Multinomial : 3 or more possible unordered types of dependent variable.

→ cat, dog, sheep.

Ordinal : 3 or more possible ordered types of dependent variable.

→ low, medium, high.

* K nearest neighbour classification :-

* Ex:- $x = (\text{Maths} = 6, \text{CS} = 8)$, $K=3$

	Maths	CS	Result
1)	4	3	F
2)	6	7	P
3)	7	8	P
4)	5	5	F
5).	8	8	P.

$$d = \sqrt{(x_{01} - x_{A1})^2 + (x_{02} - x_{A2})^2}$$

↑ ↑
 observed Actual value
 (maths: 6, CS: 8)

$$\textcircled{1} \quad \sqrt{|6-4|^2 + |8-3|^2} = \sqrt{4+25} = 5.38$$

$$\textcircled{2} \quad \sqrt{|6-6|^2 + |8-7|^2} = \sqrt{1} = \textcircled{1}$$

$$\textcircled{3} \quad \sqrt{|6-7|^2 + |8-8|^2} = \sqrt{1+0} = \textcircled{0.0001}$$

$$\textcircled{4} \quad \sqrt{|6-5|^2 + |8-5|^2} = \sqrt{1+9} = 3.16$$

$$\textcircled{5} \quad \sqrt{|6-8|^2 + |8-8|^2} = \sqrt{4} = \textcircled{2}$$

→ Here, 2nd, 3rd and 5th data points/neighbors are nearest.

- 2nd → Pass
 → 3rd → Pass
 → 5th → Pass.

∴ 3 of them are pass. ~~and hence fail~~.

∴ Probability of pass is greater than fail.

→ Therefore, we can declare he is pass.

* Conditional Probability :-

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

$$\text{ex:- } \rightarrow P(B) = \frac{30}{100} = 0.3$$

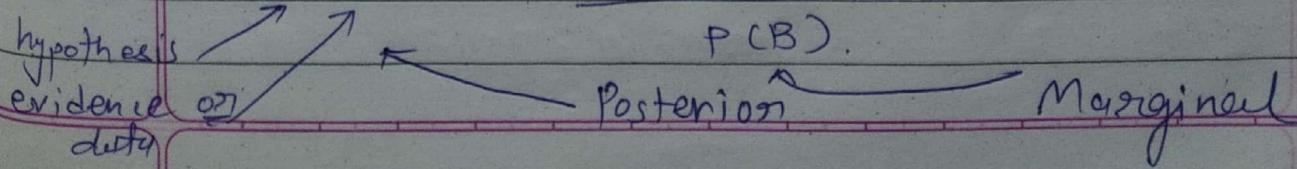
$$P(A \cap B) = \frac{20}{100} = 0.2$$

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

$$= \frac{0.2}{0.3} = 0.67$$

* Baye's Theorem :-

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



* Naive Bayes Classifier :-

Given: Fruits = { yellow, Sweet, Long }

Attributes.

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

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

→ Probability of yellow given that my fruit is orange,

$$P(\text{Yellow} | \text{Orange}) = \frac{P(\text{Orange} | \text{Yellow}) \cdot P(\text{Yellow})}{P(\text{Orange})}$$

$$\begin{aligned}
 & \text{orange + yellow} \\
 & = \frac{350}{800} \times \frac{800}{1200} \\
 & \quad \text{total yellow fruits} \\
 & \quad \text{total fruits} \\
 & = \frac{650}{1200} \\
 & = 0.538
 \end{aligned}$$

→ Probability of sweet given that my fruit is orange,

$$\begin{aligned} P(\text{Sweet} | \text{Orange}) &= \frac{P(\text{Orange} | \text{Sweet}) \cdot P(\text{Sweet})}{P(\text{Orange})} \\ &= \frac{450}{850} \times \frac{850}{1200} \\ &\quad \frac{650}{1200} \\ &= 0.692 \end{aligned}$$

→ Probability of long given that my fruit is Orange,

$$\begin{aligned} P(\text{long} | \text{Orange}) &= \frac{P(\text{Orange} | \text{long}) \cdot P(\text{long})}{P(\text{Orange})} \\ &= \frac{0}{400} \times \frac{0}{400} \\ &\quad \frac{650}{1200} \\ &= 0 \end{aligned}$$

Now, Probability of Fruit given that orange is my fruit,

$$\begin{aligned} P(\text{Fruit} | \text{Orange}) &= P(\text{Yellow} | \text{Orange}) \times \\ &\quad P(\text{Sweet} | \text{Orange}) \times \\ &\quad P(\text{long} | \text{Orange}) \\ &= 0.538 \times 0.692 \times 0 \\ &= 0 \end{aligned}$$

Now, probability of Fruit given that banana is my fruit,

$$P(\text{Fruit} | \text{Banana}) = P(\text{Yellow} | \text{Banana}) \times P(\text{Sweet} | \text{Banana}) \times P(\text{long} | \text{Banana}).$$

$$\therefore P(\text{Yellow} | \text{Banana}) = \frac{P(B|Y) \cdot P(Y)}{P(B)}$$

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

$$= \frac{400}{1200}$$

$$= 1$$

$$\& P(\text{Sweet} | \text{Banana}) = \frac{P(B|S) \cdot P(S)}{P(B)}$$

$$= \frac{300}{880} \times \frac{850}{1200}$$

$$= 0.75$$

$$\& P(\text{long} | \text{Banana}) = \frac{P(B|L) \cdot P(L)}{P(B)}$$

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

$$= \frac{400}{1200}$$

$$= 0.875$$

$$\therefore P(\text{Fruit} / \text{Banana}) = 1 \times 0.75 \times 0.875 \\ = 0.656$$

→ Now, probability of Fruit given that the fruit is from others category,

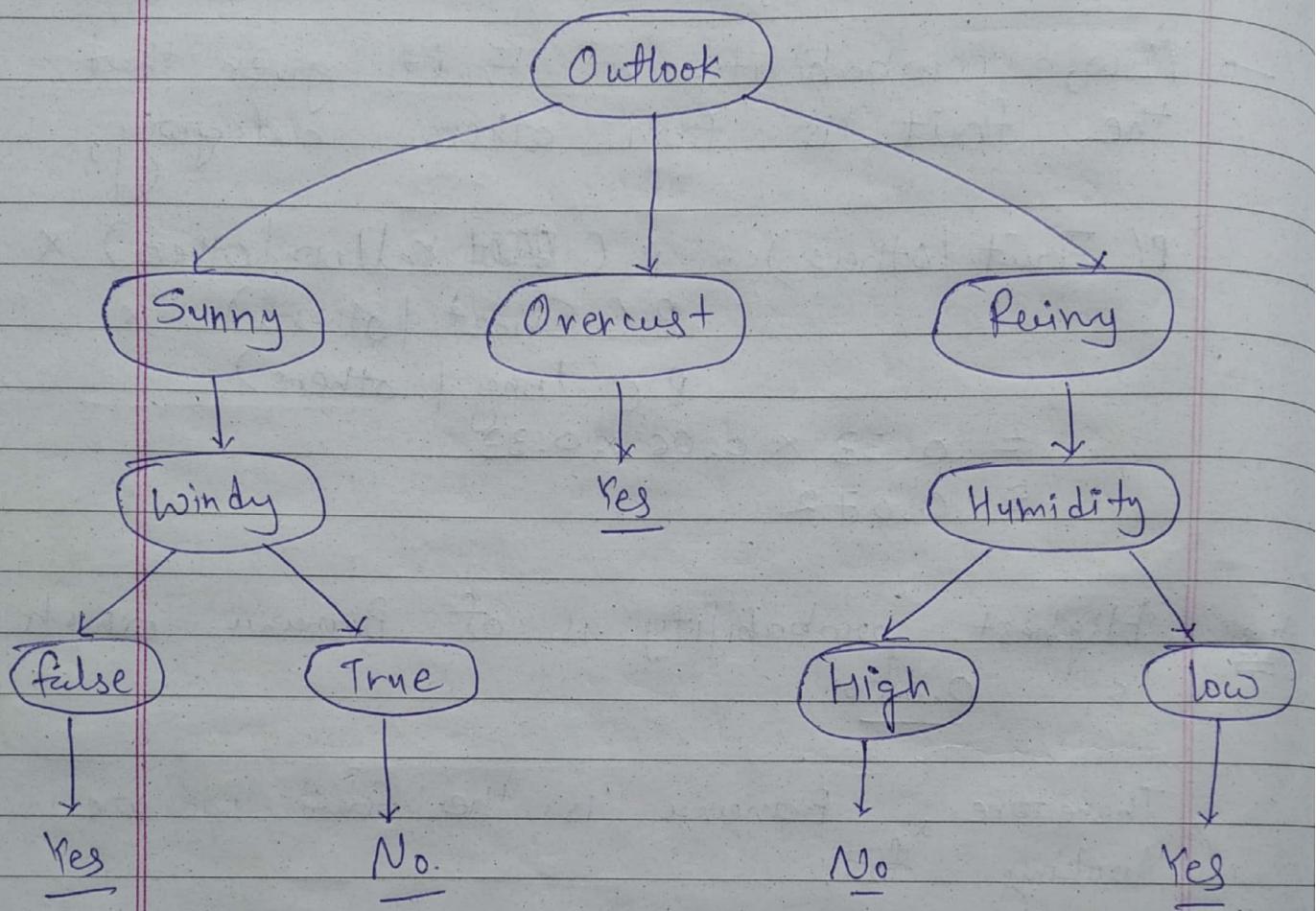
$$P(\text{Fruit} / \text{others}) = P(\text{Red yellow} / \text{others}) \times \\ P(\text{Sweet} / \text{others}) \times \\ P(\text{long} / \text{others}), \\ = 0.33 \times 0.66 \times 0.33 \\ = 0.072$$

Ans: Highest probability is of Banana which is 0.656.

Therefore, Banana is the fruit we are looking for.

* Decision Tree :-

Play Golf



* Entropy :-

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

→ One attribute :-

Play Golf	
Yes	No
9	5

$$E(5, 9) = \text{Entropy}\left(\frac{5}{14}, \frac{9}{14}\right)$$

$$= \text{Entropy}(0.36, 0.64)$$

$$= - (0.36 \cdot \log_2 0.36)$$

$$- (0.64 \cdot \log_2 0.64) = 0.94$$

→ Two attributes :-

		Play Golf		
		Yes	No	
Outlook	Sunny	3	2	5
	Overcast	4	0	4
	Rainy	2	3	5
		14		

$$\begin{aligned}
 E(\text{Play Golf}, \text{Outlook}) &= P(\text{Sunny}) \cdot E(3, 2) \\
 &\quad + P(\text{Overcast}) \cdot E(4, 0) \\
 &\quad + P(\text{Rainy}) \cdot E(2, 3).
 \end{aligned}$$

$$\begin{aligned}
 &= \left(\frac{5}{14} \right) \left(-\frac{3}{5} \cdot \log_2 \frac{3}{5} - \frac{2}{5} \cdot \log_2 \frac{2}{5} \right) \\
 &\quad + \left(\frac{4}{14} \right) \left(-\frac{4}{4} \cdot \log_2 \frac{4}{4} - \frac{0}{0} \cdot \log_2 0 \right) \\
 &\quad + \left(\frac{5}{14} \right) \left(-\frac{2}{5} \cdot \log_2 \frac{2}{5} - \frac{3}{5} \cdot \log_2 \frac{3}{5} \right) \\
 &= (0.357)(-0.6 \cdot (-0.73) - 0.4 \cdot 0) \\
 &\quad + (0.285)(0) \\
 &\quad + (0.357)(-0.4 \cdot (-1.32) - 0.6 \cdot 0) \\
 &= 0.693
 \end{aligned}$$

* Example :-

Day	Outlook	Temp	Humidity	Wind	<u>Play Tennis</u>
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No.

→ 4 Attributes : Outlook, Temp, Humidity, Wind,

→ Max IG will be considered as root node.

→ 5 No's and 9 Yes's

(1) → Entropy of entire dataset

$$\rightarrow E(S) = E(+9, -5) = \frac{9}{14} \log_2 \frac{9}{14} + \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.94$$

(2) →

Outlook :-

sunny, Overcast, Rain
 Yes \downarrow No's in sunny's rows.

$$\rightarrow E(\text{Sunny}) = E(2+, 3-)$$

$$= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$\rightarrow E(\text{Overcast}) = E(4+, 0-) = 0$$

$$\rightarrow E(\text{Rain}) = E(3+, 2-)$$

$$= -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$$

→ Grain (S, outlook)
 no. of times sunny is appearing.

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

Total \rightarrow 14 Rain \rightarrow 5

$$= \frac{5}{14} E(\text{Rain})$$

$$= 0.94 - \frac{5}{14} (0.971) - \frac{4}{14} (0) - \frac{5}{14} (0.971)$$

$$= 0.2464$$

(3) →

Temp :- Hot, Mild, Cool.

$$\rightarrow E(\text{Hot}) = E(2+, 2-) \xrightarrow{\text{same}} 1$$

$$= -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} = 1$$

$$\rightarrow E(\text{Mild}) = E(4+, 2-)$$

$$= -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.9183$$

$$\rightarrow E(\text{Cool}) = E(3+, 1-)$$

$$= -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.8113$$

→ Grain = ~~Entropy~~ (S, Temp)

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

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

$$= 0.94 - \frac{4}{14} (1) - \frac{6}{14} (0.9183)$$

$$- \frac{4}{14} (0.8113)$$

$$= 0.0289$$

(4) →

Humidity :- High, Normal

$$\rightarrow E(\text{High}) = E(3+, 4-)$$

$$= -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.9852$$

$$\rightarrow E(\text{Normal}) = E(6+, 1-)$$

$$= -\frac{6}{7} \log_2 \frac{6}{7} - \frac{1}{7} \log_2 \frac{1}{7} = 0.5916$$

→ Grain(S, Humidity)

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

$$= 0.94 - \frac{7}{14} (0.9852) - \frac{7}{14} (0.5916)$$

$$= 0.1516$$

(5) →

Wind :- Strong, Weak.

$$\rightarrow E(\text{Strong}) = E(3+, 3-) = 1$$

$$\rightarrow E(\text{Weak}) = E(6+, 2-) = 0.8113$$

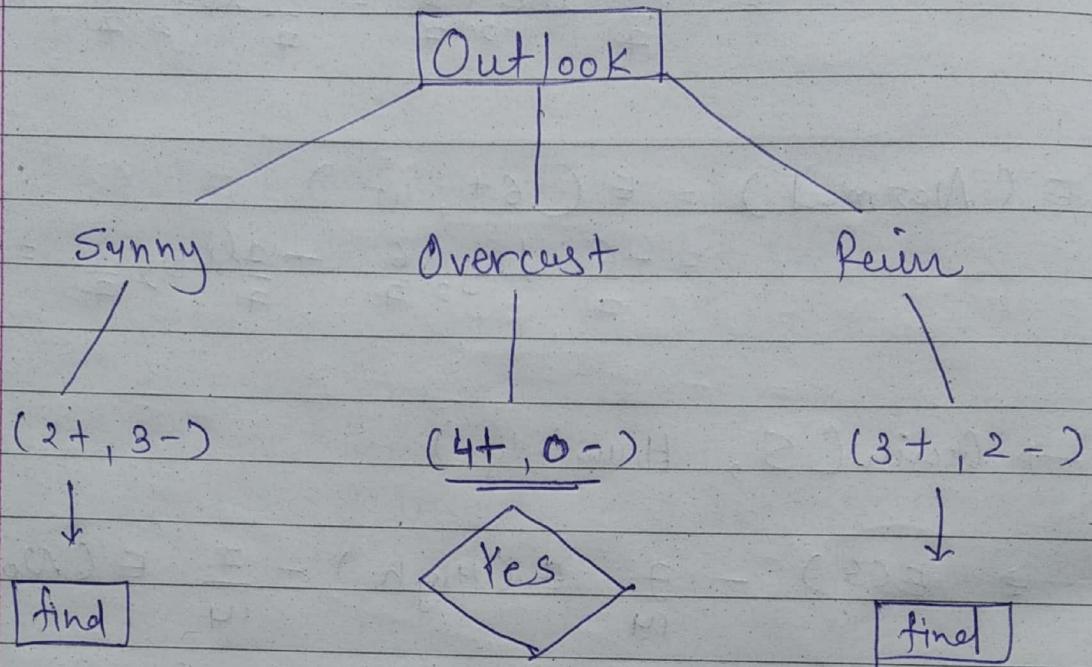
→ Grain(S, Wind)

$$= E(S) - \frac{6}{14} E(\text{Strong}) - \frac{8}{14} E(\text{Weak})$$

$$= 0.94 - \frac{6}{14} (1) - \frac{8}{14} (0.8113)$$

$$= 0.0478.$$

→ From the above calculations, Outlook is having maximum entropy.



→ Now, outlook is considered.

Therefore for sunny, ~~not~~

Day	Temp	Humidity	Wind.	Play Tennis
1	Hot	High	Weak	No
2	Hot	High	Strong	No
8	Mild	High	Weak	No
9	Cool	Normal	Weak	Yes
11	Mild	Normal	Strong	Yes.

→ We don't need to consider outlook further.

* Confusion Matrix :-

- * Consider the confusion matrix given below for a binary classifier predicting the presence of a disease.

The classifier made a total of 150 predictions. Out of those 150 cases, the classifier predicted yes 100 times and no 50 times.

In reality, 100 patients in the sample have the disease and 50 patients do not.

		Predicted No	Predicted Yes
Actual No	45 (TN)	5 (FP)	
	5 (FN)	95 (TP)	

* Accuracy :
$$\frac{TN + TP}{\text{total predictions}} \xrightarrow{\text{total correctness}}$$

$$= \frac{45 + 95}{150}$$

$$= 93.33 \%$$

* Precision : When it predicts yes, how often is it correct?

* Precision = $\frac{TP}{\text{total predicted yes}}$ = $\frac{95}{100}$ = 95%

* Recall = $\frac{TP}{\text{total Actual Yes}}$ = $\frac{95}{100}$ = 95%

Ans Example :-

	Predicted No	Predicted Yes
Actual No	2 TN	0 FP
Actual Yes	3 FN	5 TP

→ Accuracy = $\frac{TN + TP}{\text{total predictions}}$
 $= \frac{2+5}{10} = \frac{7}{10} \times 100 = 70\%$

→ Precision = $\frac{TP}{\text{total predicted yes}} = \frac{5}{5} \times 100 = 100\%$
 $(TP + FP)$

→ Recall = $\frac{TP}{\text{total actual yes}} = \frac{5}{8} \times 100 = 62.5\%$
 $(TP + FN)$

→ F1 Score = $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

* 3×3 Confusion matrix :-

	A	B	C	Total	total of actual
	Predicted				
A	Act true	15	2	3	20.
B	Act true	7	15	8	30.
C		2	3	45	50
	Total	24	20	56	100

$$\rightarrow \text{Accuracy} = \frac{TP + TN}{\text{total predicted}} \quad (\text{Cross axis})$$

$$= \frac{15 + 15 + 45}{100} * 100$$

$$= 75 \%$$

$$\rightarrow \text{Precision}_A = \frac{\text{Correctly predicted}}{\text{Total predicted}}$$

$$= \frac{15}{24} = 0.625.$$

$$\text{Precision}_B = \frac{15}{20} = 0.75$$

$$\text{Precision}_C = \frac{45}{56} = 0.80$$

$$\star \rightarrow F-1 \text{ score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{for}(A) = 2 \times \frac{0.625 \times 0.75}{0.625 + 0.75} = \frac{0.9375}{1.375} = 0.681$$



Recall = Correctly classified
Actual

$$\text{Recall}_A = \frac{15}{20} = 0.75$$

$$\text{Recall}_B = \frac{15}{30} = 0.5$$

$$\text{Recall}_C = \frac{45}{50} = 0.9.$$

* Cross validation :-

→ It is a statistical method used to estimate the performance (or accuracy) of ML models.

Types of cross validation

Exhaustive

- 1) Hold out method
- 2) k-fold cross validation
- 3) stratified k-fold cross validation

Non-exhaustive

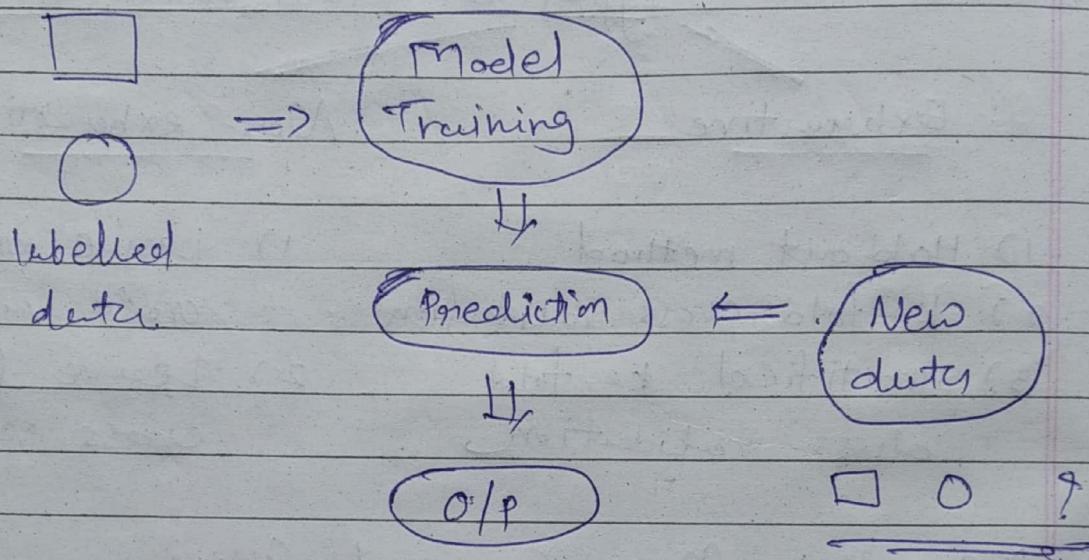
- 1) Leave one out cross validation
- 2) Leave P out cross validation

* MSE : Mean Squared Error $\therefore \frac{\sum (Actual - Forecast)^2}{n}$

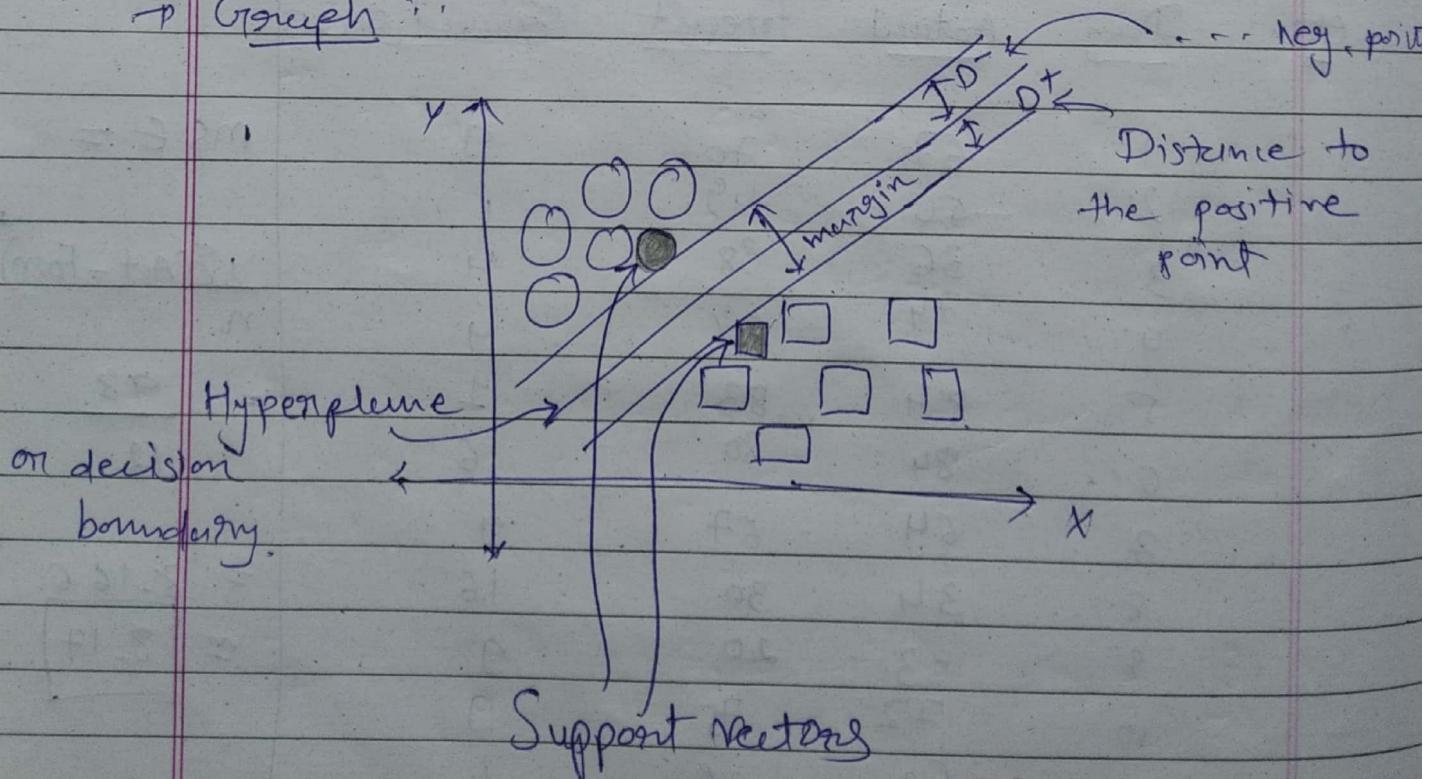
Day	Actual	Forecast	Squared error	MSE =
1	67	70	9	
2	50	49	1	
3	36	38	4	
4	74	76	4	
5	84	83	1	$= \frac{\sum (Act - Fore)^2}{n}$
6	84	80	16	
7	64	67	9	
8	34	30	16	
9	23	20	9	
10	72	75	9	
11	62	68	4	
12	42	38	16	
			98	

* Support Vector Machine :-

- Follows Supervised Learning.
- Diagram:



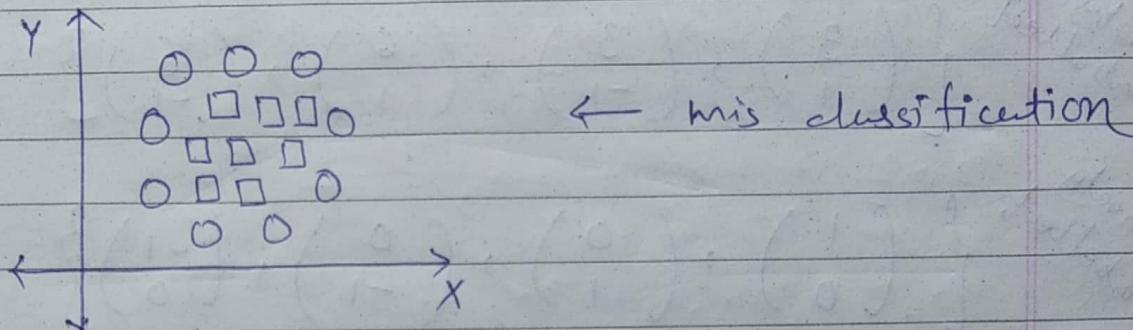
→ Graph :-



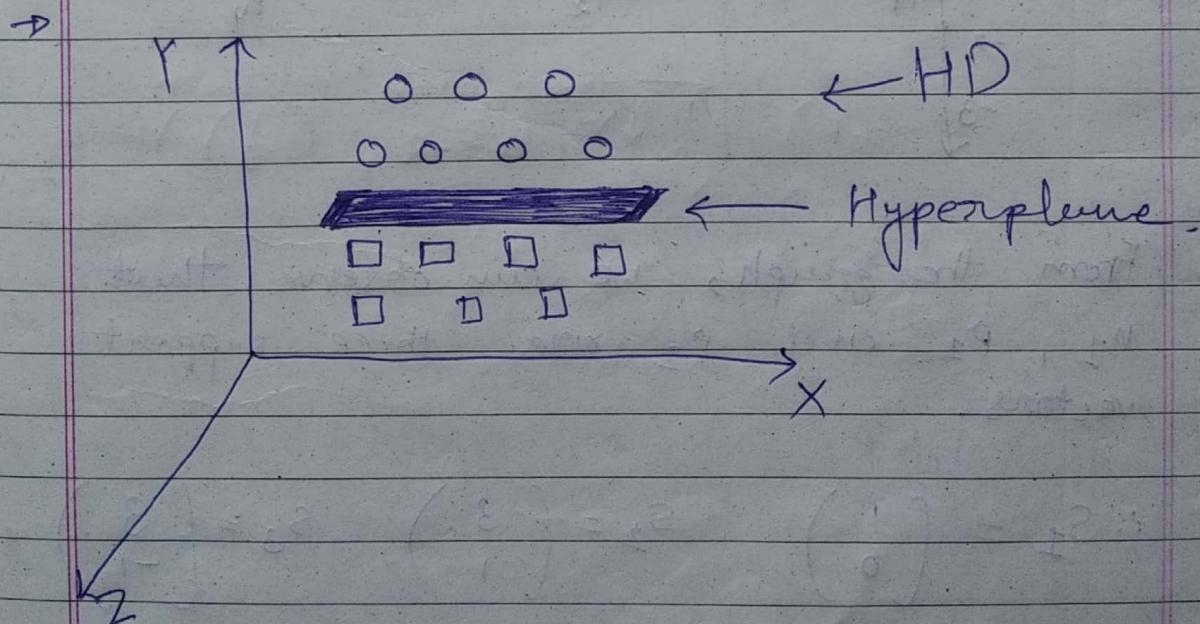
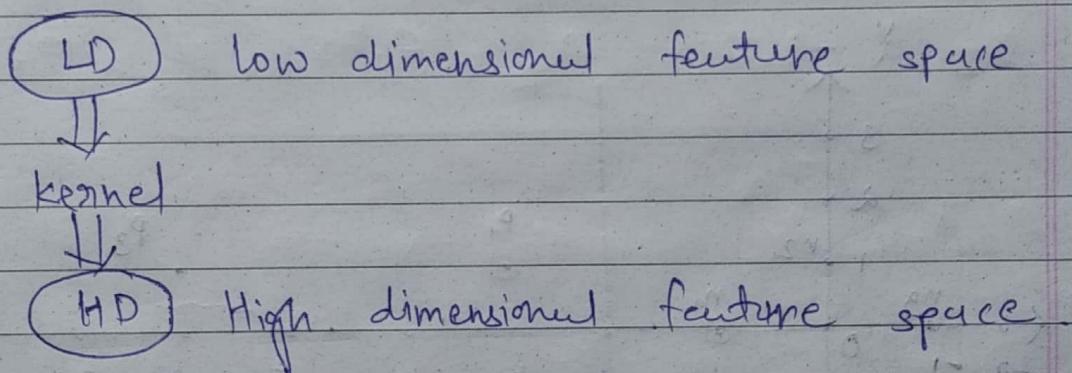
→ Maximum width of hyperplane is chosen to get better result.

→ Maximal Margin Hyperplane is selected

* Non-linear SVM :-



→ To classify this, kernel function is used



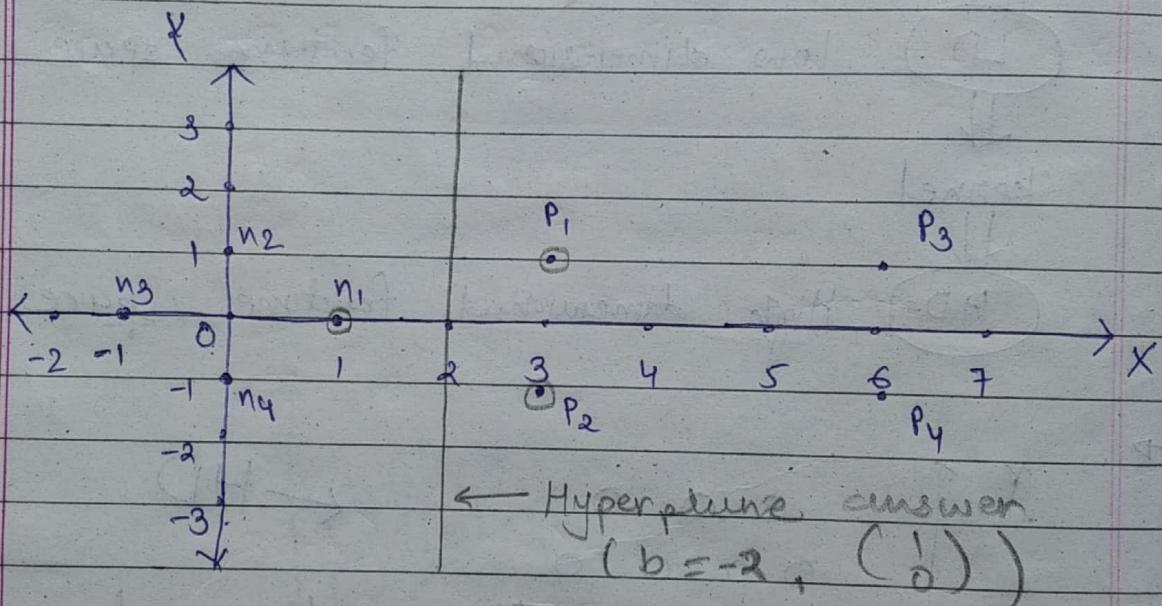
* Linear SVM Examples :-

ex:- Dataset :-

Positively labelled $\{ (3, 1), (3, -1), (6, 1), (6, -1) \}$

Negatively labelled $\{ (1, 0), (0, 1), (0, -1), (-1, 0) \}$

Soln :- Plot a graph.



- From the graph, we can observe that n_1, p_1 and p_2 are three support vectors.

$$\therefore s_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, s_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, s_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix}$$

→ Now, augment each vector with a '1' as a bias input.

$$\underline{s}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \Rightarrow \bar{\underline{s}}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \leftarrow 1 \text{ is added}$$

$$\underline{s}_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix} \Rightarrow \bar{\underline{s}}_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix}$$

$$\underline{s}_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix} \Rightarrow \bar{\underline{s}}_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

→ Calculate α ,

$$\cancel{\alpha_1 \bar{\underline{s}}_1 \cdot \bar{\underline{s}}_1} + \alpha_2 \bar{\underline{s}}_2 \cdot \bar{\underline{s}}_1 + \alpha_3 \bar{\underline{s}}_3 \cdot \bar{\underline{s}}_1 = -1$$

($\because \underline{s}_1$ is present in the other side)

$$\& \alpha_1 \bar{\underline{s}}_1 \cdot \bar{\underline{s}}_2 + \alpha_2 \bar{\underline{s}}_2 \cdot \bar{\underline{s}}_2 + \alpha_3 \bar{\underline{s}}_3 \cdot \bar{\underline{s}}_2 = +1$$

$$\& \alpha_1 \bar{\underline{s}}_1 \cdot \bar{\underline{s}}_3 + \alpha_2 \bar{\underline{s}}_2 \cdot \bar{\underline{s}}_3 + \alpha_3 \bar{\underline{s}}_3 \cdot \bar{\underline{s}}_3 = +1$$

$$\therefore \alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\& \alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} = +1$$

$$\& \alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = +1$$

$$\therefore \alpha_1(1+0+1) + \alpha_2(3+0+1) + \alpha_3(3+0+1) = -1$$

$$\& \alpha_1(3+0+1) + \alpha_2(9+1+1) + \alpha_3(9-1+1) = 1$$

$$\& \alpha_1(3+0+1) + \alpha_2(9-1+1) + \alpha_3(9+1+1) = 1$$

$$\therefore 2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$

$$\& 4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1$$

$$\& 4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1$$

$$\therefore \alpha_1 = -3.5, \alpha_2 = 0.75, \alpha_3 = 0.75$$

→ Calculate weight vector:-

$$\bar{w} = \sum_{i=1}^n \alpha_i \bar{s}_i$$

$$= (-3.5) \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ 1 \\ -1 \end{pmatrix}$$

$$= \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix}$$

its a bias

→ We can equate the last entry in \bar{w} as the hyperplane offset b and write the separating hyperplane equation;

$$y = w_1 x + b$$

where, $w = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

$$b = -2$$

if $\begin{pmatrix} 1 \\ 0 \end{pmatrix} \Rightarrow$ parallel to y axis

$\begin{pmatrix} 0 \\ 1 \end{pmatrix} \Rightarrow$ parallel to x axis

$\begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow$ us with respect to x and y axis.

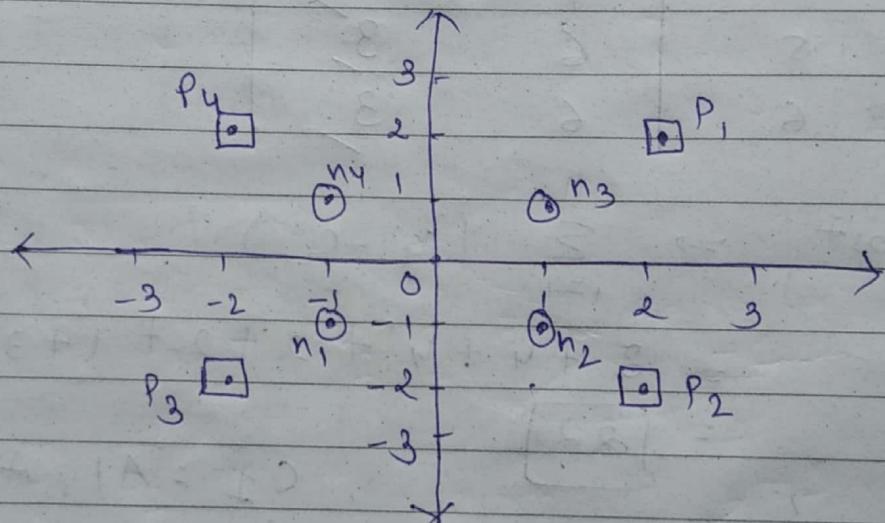
* Support Vector Machine :-

Non - Linear

Ex :- positively labelled : $\left\{ \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 \\ -2 \end{pmatrix}, \begin{pmatrix} -2 \\ -2 \end{pmatrix}, \begin{pmatrix} -2 \\ 2 \end{pmatrix} \right\}$

Negatively labelled : $\left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right\}$

Soln.
=



→ To convert these points into non-linear plane, this formula is used.

$$\phi_1 = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{cases} \begin{pmatrix} 4 - x_2 + |x_1 - x_2| \\ 4 - x_1 + |x_1 - x_2| \end{pmatrix}; & \text{if } \sqrt{x_1^2 + x_2^2} > 2 \\ \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} & \text{otherwise} \end{cases}$$

→ For positively labelled,

$$\begin{pmatrix} 2 \\ 2 \end{pmatrix} \Rightarrow \begin{pmatrix} 4 - 2 + 0 \\ 4 - 2 + 0 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

$$\begin{pmatrix} 2 \\ -2 \end{pmatrix} \Rightarrow \left(\frac{4 - (-2) + |2 - (-2)|}{4 - 2 + |2 - (-2)|} \right) = \begin{pmatrix} 10 \\ 6 \end{pmatrix}$$

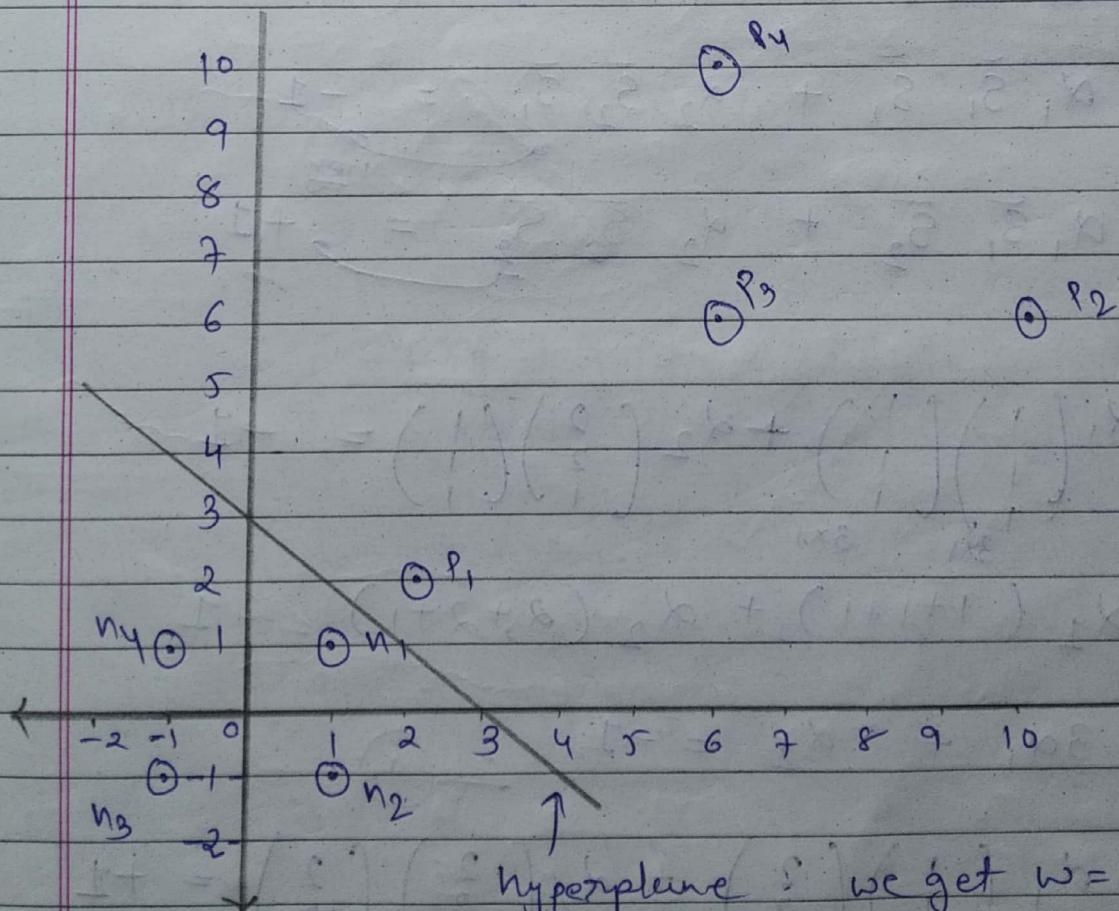
$$\begin{pmatrix} -2 \\ -2 \end{pmatrix} \Rightarrow \left(\frac{4 - (-2) + |-2 - (-2)|}{4 - (-2) + |-2 - (-2)|} \right) = \begin{pmatrix} 6 \\ 6 \end{pmatrix}$$

$$\begin{pmatrix} -2 \\ 2 \end{pmatrix} \Rightarrow \left(\frac{4 - 2 + |-2 - 2|}{4 + 2 + |-2 - 2|} \right) = \begin{pmatrix} 6 \\ 10 \end{pmatrix}$$

→ For negatively labelled,

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ -1 \end{pmatrix} \Rightarrow \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$\begin{pmatrix} -1 \\ -1 \end{pmatrix} \Rightarrow \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \Rightarrow \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$



Hyperplane : we get $w = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$
which means that it will make
45° angle with x and y axis.

→ Here, n_1 and p_1 are nearest to each other belongs from two different labelled.

$$\therefore \{ S_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, S_2 = \begin{pmatrix} 2 \\ 2 \end{pmatrix} \}$$

→ Add 1 as a bias input.

$$\tilde{S}_1 = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \quad \tilde{S}_2 = \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix}$$

→ find α ,

$$\alpha_1 \tilde{S}_1 \tilde{S}_1 + \alpha_2 \tilde{S}_2 \tilde{S}_1 = -1$$

$$\alpha_1 \tilde{S}_1 \tilde{S}_2 + \alpha_2 \tilde{S}_2 \tilde{S}_2 = +1$$

$$\therefore \alpha_1 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = -1$$

$$\therefore \alpha_1 (1+1+1) + \alpha_2 (2+2+1) = -1$$

$$\therefore 3\alpha_1 + 5\alpha_2 = -1 \quad \text{--- (1)}$$

$$\text{and } \alpha_1 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix} \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix} = +1$$

$$\therefore \alpha_1(2+2+1) + \alpha_2(4+4+1) = +1$$

$$\therefore 5\alpha_1 + 9\alpha_2 = 1$$

(2)

$$\therefore 5\alpha_1 + 9\alpha_2 = 1$$

$$\therefore ① \times 5 - ② \times 3.$$

$$\therefore 15\alpha_1 + 25\alpha_2 = -5$$

$$= 15\alpha_1 + 27\alpha_2 = -3$$

$$-2\alpha_2 = -8$$

$$\therefore \boxed{\alpha_2 = 4}$$

$$\therefore 5\alpha_1 + 9(4) = 1$$

$$\therefore 5\alpha_1 = -35 \Rightarrow \boxed{\alpha_1 = -7}$$

→ Calculate weight vector :-

$$\begin{aligned} \vec{w}^N &= \sum_{i=1}^n \alpha_i \vec{s}_i \\ &= (-7) \left(\begin{array}{c} 1 \\ 1 \end{array} \right) + (4) \left(\begin{array}{c} 2 \\ 1 \end{array} \right) \\ &= \left(\begin{array}{c} -7 \\ -7 \end{array} \right) + \left(\begin{array}{c} 8 \\ 4 \end{array} \right) \\ &= \left(\begin{array}{c} 1 \\ 1 \\ -3 \end{array} \right) \end{aligned}$$

bias input = b.

$$\therefore y = w \cdot x + b$$

$$w = \left(\begin{array}{c} 1 \\ 1 \end{array} \right)$$

* K-means clustering :-

→ Suppose that the data mining task is to cluster points into three clusters,

where points are,

$A_1(2, 10)$, $A_2(2, 5)$, $A_3(8, 4)$.

$B_1(5, 8)$, $B_2(7, 5)$, $B_3(6, 4)$.

$C_1(1, 2)$, $C_2(4, 9)$.

→ The distance function is Euclidean distance.

Soln:- A_1 , B_1 , C_1 are the centre of each cluster. (if not given, we can select).

∴ Initial Centroids : A_1 , B_1 , C_1 .

Dataset	Distance to centroid	Cluster	New Cluster
$x_1 \ y_1$	$x_2 \ y_2$	$x_2 \ y_2$	$x_2 \ y_2$
A1 2 10	0	3.61	8.06
A2 2 5	5	4.24	3.16
A3 8 4	8.49	5	7.28
B1 5 8	3.61	0	7.21
B2 7 5	7.07	3.61	6.71
B3 6 4	7.21	4.12	5.39
C1 1 2	8.06	7.21	0
C2 4 9	2.24	1.41	7.62

cluster no. ① ② ③

→ Now, calculate Euclidean distance using formulae,

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

→ Now, smallest cluster will be assigned in each row.

ex: for A1 : 0, 3.61, 8.06

0 is the smallest \Rightarrow cluster ① will be assigned.

for A2 : 5, 4.24, 3.16 \Rightarrow 3.16 is smallest

\therefore cluster ③ will be assigned

like wise for all Datapoints.

→ After assigning all the clusters, we need to calculate new centroids for each one.

→ for cluster no. ③ :

$$\frac{2+1}{2}, \frac{5+2}{2} \Rightarrow (1.5, 3.5)$$

→ for cluster no. ② :

$$\left(\frac{8+5+7+6+4}{5}, \frac{4+8+5+4+9}{5} \right) \\ \Rightarrow (6, 6)$$

→ for cluster no. ① :

$$\left(\frac{2}{1}, \frac{10}{1} \right) \Rightarrow (2, 10).$$

∴ Current Centroids : A1 (2, 10)
B1 (6, 6)
(1 (1.5, 3.5))

→ Now, using current centroids, we need to perform the same operations.

Datapoints		Distance to Centroid				Cluster	New cluster
	x _i y _i	2 10	6 6	1.5 3.5			
A1	2 10	0	5.66	6.52	1		(1)
A2	2 5	5	4.12	1.58	3		(3)
A3	8 4	8.49	2.83	6.52	2		2
B1	5 8	3.61	2.24	5.7	2		2
B2	7 5	7.07	1.41	5.7	2		2
B3	6 4	7.21	2	4.53	2		2
C1	1 2	8.06	6.4	1.58	3		(3)
C2	4 9	2.24	3.61	6.04	2	X 2	(1)

→ Here, cluster is what we found in our first round and new cluster is the one having smallest distance.

∴ Assign numbers to new clusters

→ C2 is moving from cluster 2 to 1
 ∴ We need to again calculate the new centroids.

→ for cluster no (3),

$$(1.5, 3.5)$$

for cluster no (2),

$$\left(\frac{8+5+7+6}{4}, \frac{4+8+5+4}{4} \right) \Rightarrow (6.5, 5.25)$$

for cluster ①,

$$\left(\frac{2+4}{2}, \frac{10+9}{2} \right) \Rightarrow (3, 9.5)$$

→ Current centroids, A1 (3, 9.5)
 B1 (6.5, 5.25)
 C1 (8.6, 1.5, 3.5)

Data points	Distance to centroid					Cluster	New Cluster
	3 9.5	6.5 5.25	1.5 3.5				
A1	2 10	1.12	6.54	6.52	1	1	
A2	2 5	4.61	4.51	1.58	3	3	
A3	8 4	7.43	1.95	6.52	2	2	
B1	5 8	2.5	3.13	5.7	2	1	
B2	7 5	6.62	0.56	5.7	2	2	
B3	6 4	6.26	1.35	4.53	2	2	
C1	1 2	7.76	6.39	1.58	3	3	
C2	4 9	1.12	4.51	6.04	1	1	

↑ New cluster of last round.

→ Again, cluster 2 is moved to 1.

∴ Calculate new centroids.

$$A1 (3.67, 9)$$

$$B1 (7, 4.33)$$

$$C1 (1.5, 3.5)$$

Data points		Distance to Centroid			Cluster	New Cluster
		3.67	9	7.43	1.5	3.5
A1	2	10	1.94	7.56	6.52	1
A2	2	5	4.33	5.04	1.58	3
A3	8	4	6.62	1.05	6.52	2
B1	5	8	1.67	4.18	5.7	1
B2	7	5	5.21	0.67	5.7	2
B3	6	4	5.52	1.05	4.53	2
C1	1	2	7.49	6.44	1.58	3
C2	4	9	0.33	5.55	6.04	1

→ Here, All the clusters and new clusters are same.

∴ ①st cluster = A1, B1, C2

②nd cluster : A3, B2, B3

③rd cluster : A2, C1.

* Newest Neighbour Clustering

Ex: Dataset is of points (A, B, C, D, E)

	A	B	C	D	E	
A	0	E	2	2	3	Column wise
B		0	2	4	3	
C			0	1	5	
D				0	3	
E					0.	distance given

threshold = 2

Step \rightarrow

$$K_1 = \{ A, B, C, D \} , K_2 = \{ E \}$$

initially added

$\leq t$

\rightarrow For threshold = 1 :

$$K_1 = \{ A, B \} (\because \text{dist}(A, B) \leq t)$$

initially added

$$K_2 = \{ C, D \} (\because \text{dist}(C, D) \leq t)$$

$$K_3 = \{ E \}$$

<u>ex:</u>	$A_1 = (2, 10)$	$A_4 = (5, 8)$	$A_7 = (1, 2)$
	$A_2 = (2, 5)$	$A_5 = (7, 5)$	$A_8 = (4, 9)$
	$A_3 = (8, 4)$	$A_6 = (6, 4)$	$t = 4$

Soln: A_1 is placed in a cluster by it self.

$$\therefore K_1 = \{A_1\}$$

$$A_2 \rightarrow d(A_1, A_2) = \sqrt{2+5} = 5 > t$$

$$\therefore K_2 = \{A_2\}$$

$$A_3 \rightarrow d(A_1, A_3) = \sqrt{36+36} = 8.40 > t$$

$$d(A_2, A_3) = \sqrt{37} > t \quad \text{form a new cluster}$$

$$\therefore K_3 = \{A_3\}$$

$$A_4 \rightarrow d(A_1, A_4) = \sqrt{13} = < t$$

$$d(A_2, A_4) = \sqrt{18}$$

$$d(A_3, A_4) = \sqrt{25}$$

Add into already
initiated cluster having
nearest point or
smallest distance

$$\therefore K_1 = \{A_1, A_4\}$$

$$A_5 \rightarrow d(A_1, A_5) = \sqrt{50}$$

$$d(A_2, A_5) = \sqrt{25}$$

$$d(A_3, A_5) = \sqrt{2} < t$$

$$d(A_4, A_5) = \sqrt{13}$$

$$\therefore K_3 = \{A_3, A_5\}$$

$$A_6 \rightarrow d(A_5, A_6) = \sqrt{2} < t$$

$$d(A_3, A_6) = \sqrt{2} < t$$

$$\therefore K_3 = \{A_3, A_5, A_6\}$$

$$A_7 \rightarrow d(A_2, A_7) = \sqrt{10} < t$$

$$\therefore K_2 = \{A_2, A_7\}$$

$$A_8 \rightarrow d(A_8, A_4) = \sqrt{2} < t$$

$$\therefore K_1 = \{A_1, A_4, A_8\}$$

Clusters:

$$K_1 = \{A_1, A_4, A_8\}$$

$$K_2 = \{A_2, A_7\}$$

$$K_3 = \{A_3, A_5, A_6\}$$

if $> t \Rightarrow$ form a new cluster

else if $\leq t$

\Rightarrow Add into an already assigned cluster

(nearest distance)

* Clusters Using Single and Complete Link Clustering :-

Agglomerative Clustering :-

- Consider the following set of one dimensional data (points).

18, 22, 25, 42, 27, 43.

- Apply Agglomerative hierarchical clustering algorithm to build the hierarchical clustering dendrogram.

- Merge the clusters using Min distance and update the proximity matrix accordingly.

- distance matrix :

	18	22	25	27	42	43
18	0	4	7	9	24	25
22	4	0	3	5	20	21
25	7	3	0	2	17	18
27	9	5	2	0	15	16
42	24	20	17	15	0	1
43	25	21	18	16	1	0

(1) ← min distance.

Single Link

Date _____
Page _____

Solⁿ

: - Min distance among all is 1
which is b/w 42 and 43.

: Merge 42 and 43.

	18	22	25	27	42, 43
18	0	4	7	9	24
22	4	0	3	5	20
25	7	3	0	②	17
27	9	5	2	0	15 min of (17, 15)
42, 43	24	20	17	15	0

→ Now, 2 is min dist. b/w 25 and 27.
∴ merge both of them.

	18	22	25, 27	42, 43
18	0	4	7	24
22	4	0	③	20
25, 27	7	3	0	17, 15
42, 43	24	20	17, 15	0

→ 3 is min dist. ∴ merge 22, ~~17, 15~~ 25, 27
b/w 22 and 25, 27

	18	22, 25, 27	42, 43
18	0	④	24
22, 25, 27	4	0	20, 15
42, 43	24	20, 15	0

→ 4 is min dist b/w 18, and 22, 25, 27.

18, 22, 25, 27

42, 43.

18, 22, 25, 27

0

25 (15)

42, 43

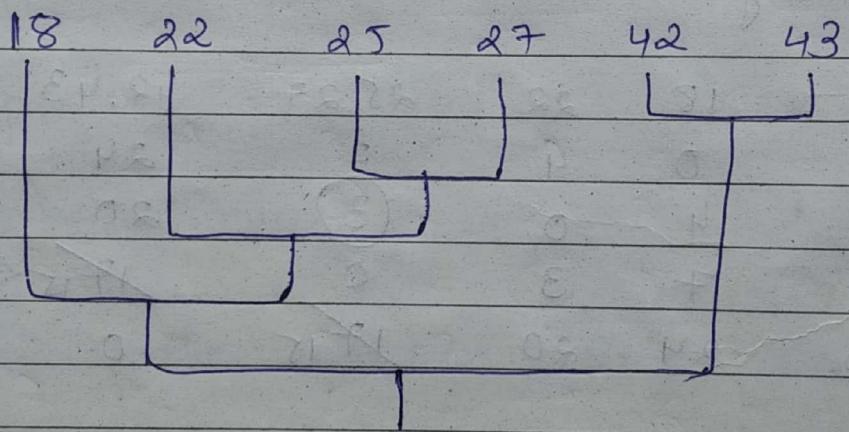
24 15

0

$\rightarrow \min = 15 \therefore 18, 22, 25, 27 + 42, 43.$

		18, 22, 25, 27, 42, 43					
18, 22, 25, 27, 42, 43		0					

\therefore dendrogram :-



* Complete Link :-

distance mat:

	1	2	3	4	5
1	0	4	7	9	(1)
2	4	0	3	5	3
3	7	3	0	2	6
4	9	5	2	0	8
5	1	3	6	8	0

Complete Link

Date _____
Page _____

Sol^y.
=

min dist = 1. \therefore merge 1 and 5
b/w 1 & 5.
and TAKE MAX
of two distance except 0.

	1, 5	2	3	4
1, 5	0.	4	7	9
2	4	0	3	5
3	7	3	0	2
4	9	5	2.	0

\rightarrow min is 2 b/w 3 and 4
 \therefore merge 3 and 4.

	1, 5	2	3, 4
1, 5	0	9	7
2	9	0	3
3, 4	7	3	0

	1, 5	2	3, 4
1, 5	0	4	9
2	4	0	5
3, 4	9	5	0

\rightarrow min = 4 b/w 1, 5 and 2 \therefore merge them.

	1, 5, 2	3, 4
1, 5, 2	9	9
3, 4	9	0

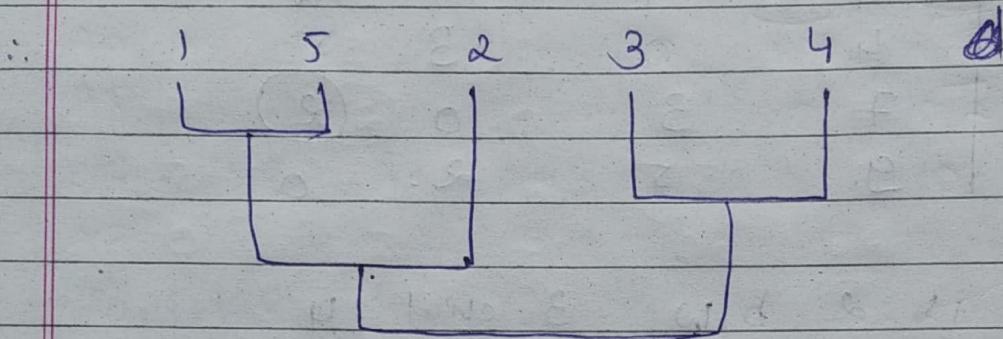
\therefore min = 9 b/w 1, 5, 2 & 3, 4.

1, 5, 2, 3, 4

1, 5, 2, 3, 4

O

∴ 1 2 3 4 5 — development

* Average Link :-

	A	B	C	D	E	
A	0	(1)	2	2	3	<u>Avg</u>
B	1	0	2	5	4	
C	2	2	0	5	4	
D	2	5	3	0	6	
E	3	4	6	6	0	

Sel → ①

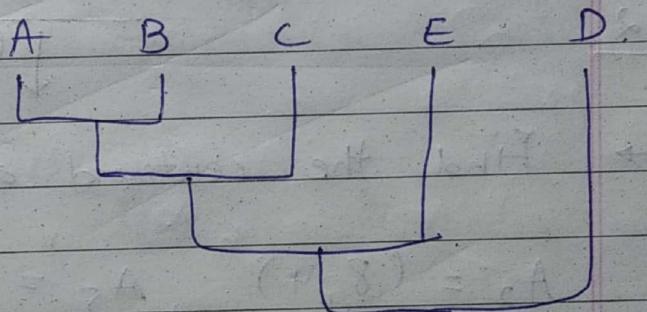
	A, B	C	D	E	
A, B	0	(2)	3.5	3.5	
C	2	0	5	4	
D	3.5	5	0	6	
E	3.5	4	6	0	

2	A, B, C	D	E
A, B, C	0	4.25	3.75 min
D	4.25	0	6
E	3.75	6	0
			Ans

3	A, B, C, E	D.
A, B, C, E	0	5.125 min.
D	5.125	0.

4	A, B, C, E, D
	0.

denoegerm



* Centroid Link :-

$$A_1 = (2, 10)$$

$$A_2 = (2, 5)$$

$$A_3 = (8, 4)$$

$$A_4 = (5, 8)$$

$$A_5 = (7, 5)$$

$$A_6 = (6, 4)$$

$$A_7 = (8, 2)$$

$$A_8 = (4, 9)$$

→ dist. matrix:

	A1	A2	A3	A4	A5	A6	A7	A8
A1	0	$\sqrt{25}$	$\sqrt{36}$	$\sqrt{13}$	$\sqrt{50}$	$\sqrt{52}$	$\sqrt{65}$	$\sqrt{5}$
A2		0	$\sqrt{37}$	$\sqrt{18}$	$\sqrt{25}$	$\sqrt{17}$	$\sqrt{10}$	$\sqrt{20}$
A3			0	$\sqrt{25}$	$\sqrt{2}$	$\sqrt{2}$	$\sqrt{53}$	$\sqrt{41}$
A4				0	$\sqrt{13}$	$\sqrt{17}$	$\sqrt{52}$	$\sqrt{2}$
A5					0	$\sqrt{2}$	$\sqrt{45}$	$\sqrt{25}$
A6						0	$\sqrt{29}$	$\sqrt{29}$
A7							0	$\sqrt{58}$
A8								0

→ $\overbrace{\{A_1\} \{A_2\}}^{\text{Dk52}} \overbrace{\{A_3, A_5, A_6\}}^{\text{B}} \overbrace{\{A_4, A_8\}}^{\text{C}} \overbrace{\{A_7\}}$

→ Find the centroid of these two clusters

$$A_3 = (8, 4), A_5 = (7, 5), A_6 = (6, 4)$$

$$\therefore \text{centroid } B = \left(\frac{8+7+6}{3}, \frac{4+5+4}{3} \right)$$

$$\therefore B = (7, 4.33)$$

$$A_4 = (5, 8) \quad A_8 = (4, 9)$$

$$\therefore \text{centroid } C = \left(\frac{5+4}{2}, \frac{8+9}{2} \right) = (4.5, 8.5)$$

$$\therefore C = (4.5, 8.5)$$

$$\rightarrow d(A_1, B) = \sqrt{(7-2)^2 + (4-3-10)^2}$$

D ≤ 3

$$A_1: (2, 10) = 7.55$$

$$B: (7, 4, 3)$$

$$d(A_1, C) = \sqrt{(4.5-2)^2 + (8-5-10)^2}$$

$$A_1: (2, 10) = 2.91 \leq 3$$

$$C: (4.5, 8.5)$$

$$\therefore \{A_1, A_4, A_8\} \quad \{A_2\} \quad \{A_3, A_5, A_6\} \quad \{A_7\}$$

$$\boxed{\text{dist}} = \sqrt{10}$$

$$= 3.16 > 3$$

\rightarrow Centroid of A_1, A_4, A_8

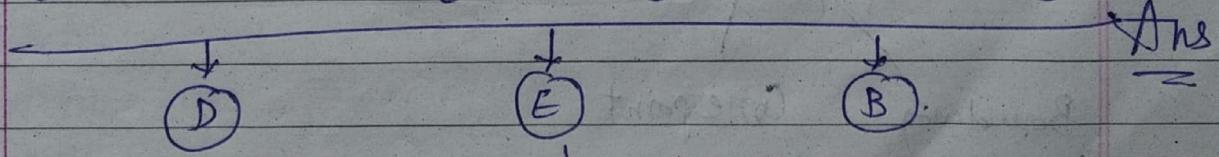
$$D_b \leq 4 \quad = \left(\frac{2+5+4}{3}, \frac{10+8+9}{3} \right)$$

$$D = (3.66, 9)$$

~~and E (Centroid of A_2, A_7)~~

$$d(A_2, D) > D_b \leq 4$$

$$\therefore \rightarrow \{A_1, A_4, A_8\} \quad \{A_2, A_7\} \quad \{A_3, A_5, A_6\}$$



$$d(D, B) = 5.74$$

$$d(D, E) = 5.9$$

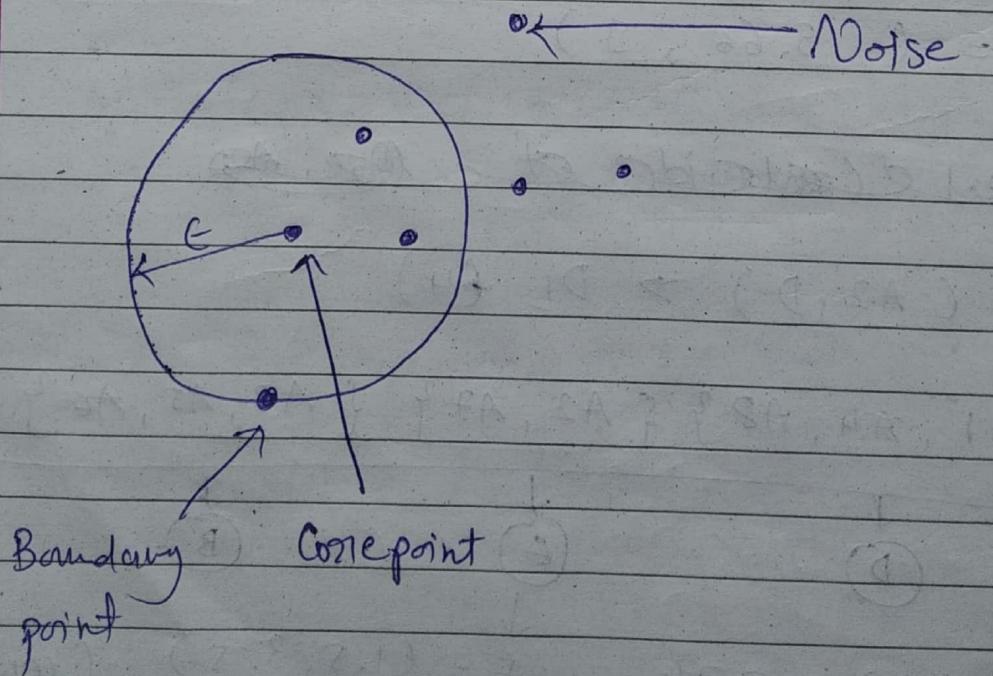
$$E = (1.5, 3.5) \quad (\text{centroid})$$

\rightarrow Ans

* DBScan :-

(Density based Clustering method)

- ϵ , Minpts ($\epsilon > 0$)
 - ↑ ↑
 - Radius Neighbourhood
points within ϵ
- Core object : got at least minpts objects or points within ϵ . (including itself)
- Boundary points : An object or a point which is not a core point but it is in neighbourhood of the core point within radius ϵ .
- Noise / outlier : Neither a core point nor a boundary point



Ex:- $\epsilon = 3.5$, MinPts = 3

$$A_1 = (5, 7)$$

$$A_2 = (8, 4)$$

$$A_3 = (3, 3)$$

$$A_4 = (4, 4)$$

$$A_5 = (3, 7)$$

$$A_6 = (6, 7)$$

$$A_7 = (6, 1)$$

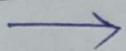
$$A_8 = (5, 5)$$

→ dist. matrix:



A ₁	A ₂	A ₃	A ₄	A ₅	A ₆	A ₇	A ₈
----------------	----------------	----------------	----------------	----------------	----------------	----------------	----------------

A ₁	0
----------------	---



A ₂	4.24	0
----------------	------	---

A ₃	4.47	5.1	0
----------------	------	-----	---

A ₄	3.16	4	1.41	0
----------------	------	---	------	---

A ₅	2	5.83	4	3.16	0
----------------	---	------	---	------	---

A ₆	1	3.61	5	3.61	3	0
----------------	---	------	---	------	---	---

A ₇	6.08	3.61	3.61	3.61	6.71	6	0
----------------	------	------	------	------	------	---	---

A ₈	2	3.16	2.83	1.41	2.83	2.24	4.12	0
----------------	---	------	------	------	------	------	------	---

Sol? ① Points : Their Neighbours within 3.5 Check

A₁ A₄, A₅, A₆, A₈

A₂ A₈

A₃ A₄, A₈

A₄ A₅, A₈, A₁, A₃

A₅ A₆, A₈, A₁, A₄

A₆ A₈, A₁, A₅

A₇ None

A₈ A₁, A₃, A₄, A₅, A₆, A₂

Rowwise &
Columnwise

② → Checking whether each point is core point or not:

If a point has minimum / at least MinPts points (including itself) in the radius ϵ , it is declared as core point.

$S_1, S_3, S_4, S_5,$
∴ $A_1, \cancel{A_2}, A_4, A_5, A_6, A_8$ are core points.
 A_3 ← (∴ including itself)
and B, A_2, A_7 are Noise.

③ → Now, checking whether A_2 and A_7 can be qualified as a boundary point?

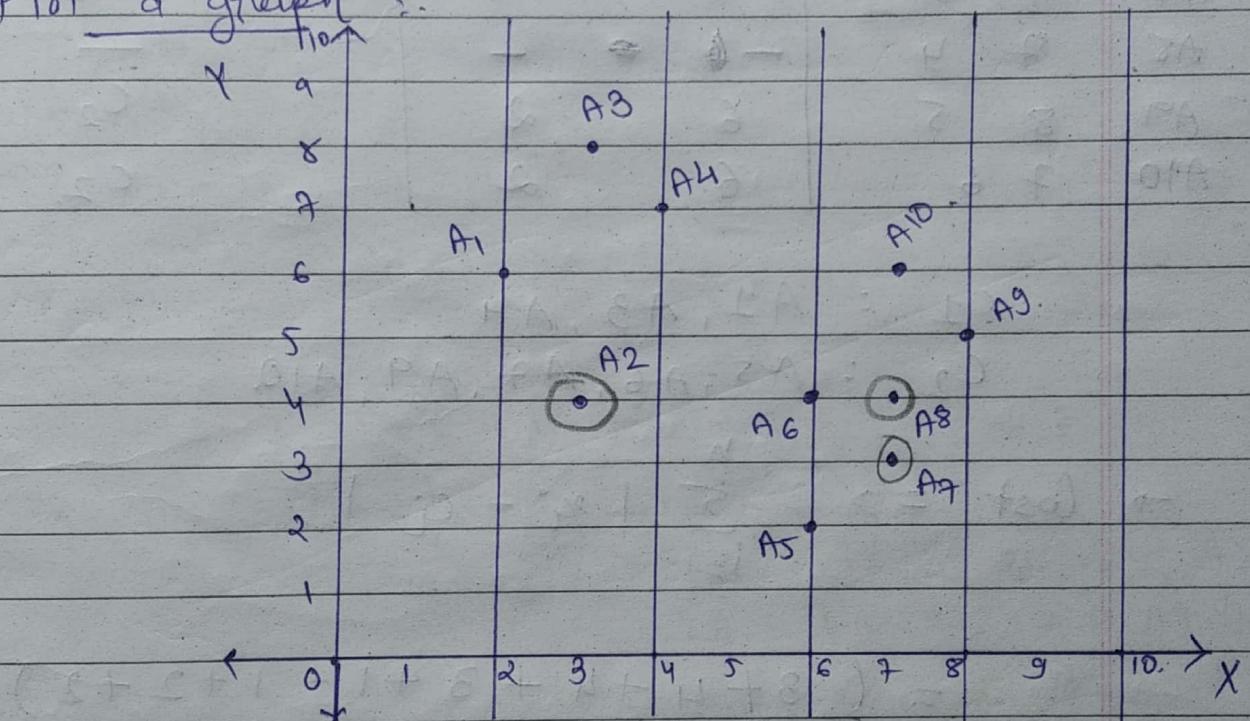
→ Here, A_2 is in the vicinity of A_8
∴ It can be considered as a boundary point.
while A_7 is not the neighbour of any point (core point).

∴ $A_2 \rightarrow$ Boundary Point
 $A_7 \rightarrow$ Noise

* Medoids Link :-

$$\begin{array}{ll}
 A_1 = (2, 6) & A_6 = (6, 4) \\
 A_2 = (3, 4) & A_7 = (7, 0) \\
 A_3 = (3, 8) & A_8 = (7, 4) \\
 A_4 = (4, 7) & A_9 = (8, 5) \\
 A_5 = (6, 2) & A_{10} = (7, 6)
 \end{array}$$

① → Plot a graph :



→ Select two random (within two diff. clusters observing from the graph) objects.

$$c_1 = (3, 4)$$

$$c_2 = (7, 4)$$

② → Now, create a table of distance from all objects to c_1 and c_2 .

(Manhattan distance
 $|x_1 - x_2| + |y_1 - y_2|$
 is used) $C_1 = (3, 4)$
 $C_2 = (7, 4)$

Date _____
 Page _____

Data points	distance from		Cluster
	C_1	C_2	
A1 2 6	3	7	$\text{No. min } C_1$
A2 3 4	-	-	-
A3 3 8	4	8	C_1
A4 4 7	4	6	C_1
A5 6 2	5	3	C_2
A6 6 4	3	1	C_2
A7 7 3	5	1	C_2
A8 7 4	-	-	-
A9 8 5	6	2	C_2
A10 7 6	6	2	C_2

$C_1 : A_1, A_3, A_4$

$C_2 : A_5, A_6, A_7, A_9, A_{10}$

$$\rightarrow \text{Cost} \Rightarrow \sum_{i=1}^n |x_i - c_i|$$

$$= (3+4+4+3+1+1+2+2) \\ = 20.$$

\rightarrow Again we need to select one more point as C_3 .

(Choose such a point which is near to either C_1 or C_2 .)

Here, let $C_3 = (7, 3)$,

$$C_1 = (3, 4)$$

$$C_3 = (7, 3)$$

Date _____
Page _____

Data points

distance from

cluster.

A1 2 6

$C_1 \quad C_3$

3 8

min

C_1

A2 3 4

—

—

A3 3 8

4 9

C_1

A4 4 7

4 7

C_1

A5 6 2

5 2

C_2

A6 6 4

3 2

C_2

A7 7 3

8 —

—

A8 7 4

4 ~~1~~ 1

C_2

A9 8 5

6 3

C_2

A10. 7 6.

6 3

C_2

$$\text{Cost} \Rightarrow \sum_{i=1}^n |x_i - c_i|$$

$$= 3 + 4 + 4 + 2 + 2 + 1 + 3 + 3$$

$$= 22$$

$C_1 : A1, A3, A4, A2$

$C_3 : A5, A6, A8, A9, A10, A7$

→ Here, cost increases.

∴ Previously selected C_1 and C_2 were best.

∴ C_1 and C_2 are center of clusters.

$C_1 : \{(2, 6), (3, 8), (4, 7), (3, 4)\}$

C_1 itself

$C_2 : \{(6, 2), (6, 4), (7, 3), (7, 4), (8, 5), (7, 6)\}$

* Association Rules :-

IMP

Apriori Algorithm :-

- Here are a dozen sales transactions.
- The objective is to use this transaction data to find affinities b/w products, that is, which products sell together often.
- The support level will be set at 33 percent.
The confidence level will be set at 50 percent.
- Rule $X \rightarrow Y$:

$$\text{Support} = \frac{\text{freq}(x, y)}{N}$$

$$\text{Confidence} = \frac{\text{freq}(x, y)}{\text{freq}(x)}$$

Example : Min support = 50 %.

Threshold confidence = 70 %.

Transaction ID

	Items			
T ₁	1	3	4	
T ₂	2	3	5	
T ₃	1	2	3	5
T ₄	2	5		

$$\text{Confidence} = \frac{S(A \cup B)}{S(A)}$$

$$\text{Support} = \frac{\text{freq}}{\text{total transaction}}$$

Date _____
Page _____

Solⁿ

(1) itemset

support

min support = 50%
confidence = 70%

$$1 \quad 2/4 = 50\% \quad \checkmark$$

$$2 \quad 3/4 = 75\% \quad \checkmark$$

$$3 \quad 3/4 = 75\% \quad \checkmark$$

$$4 \quad 1/4 = 25\% \quad \checkmark$$

$$5 \quad 3/4 = 75\% \quad \checkmark$$

$$\therefore \{1, 2, 3, 5\}$$

(2)

itemset

support

$$\{1, 2\} \quad 2/4 = 50\% \quad \checkmark$$

$$\{2, 3\} \quad 2/4 = 50\% \quad \checkmark$$

$$\{3, 5\} \quad 2/4 = 50\% \quad \checkmark$$

$$\{5, 1\} \quad 2/4 = 50\% \quad \checkmark$$

$$\{1, 3\} \quad 2/4 = 50\% \quad \checkmark$$

$$\{2, 5\} \quad 3/4 = 75\% \quad \checkmark$$

$$\therefore \{1, 2, 3, 5\}$$

(3)

itemset

support

$$\{1, 2, 3\} \quad 1/4 = 25\% \quad \checkmark$$

$$\{1, 2, 5\} \quad 1/4 = 25\% \quad \checkmark$$

$$\{2, 3, 5\} \quad 2/4 = 50\% \quad \checkmark$$

$$\{1, 3, 5\} \quad 1/4 = 25\% \quad \checkmark$$

$$\therefore \{2, 3, 5\}$$

(4) Find Association Rule :-

$$\therefore \text{Confidence} = \frac{S(A \cup B)}{S(A)}$$

$$\text{For, } (2 \wedge 3) \rightarrow 5 \Rightarrow \frac{S((2 \wedge 3) \cup 5)}{S(2 \wedge 3)}$$

$$= \frac{2}{2} = 100\%$$

Support of
2, 3, 5.

Support of
2, 3 only.

{2, 3, 5}

Date _____

Page _____

Rules

support

Confidence

$$(2^1 3) \rightarrow 5$$

2

$$2/2 = 100\%$$

$$(2^1 5) \rightarrow 3$$

2

$$2/3 = 66\%$$

$$(3^1 5) \rightarrow 2$$

2

$$2/2 = 100\%$$

$$5 \rightarrow (2^1 3)$$

2

$$2/3 = 66\%$$

$$3 \rightarrow (2^1 5)$$

2

$$2/3 = 66\%$$

$$2 \rightarrow (3^1 5)$$

2

$$2/3 = 66\%$$

divide by
support of

→ Confidence threshold is 70%.

∴ $(2^1 3) \rightarrow 5$ and $(3^1 5) \rightarrow 2$
are applicable rules.

— By Rushik. Rathod