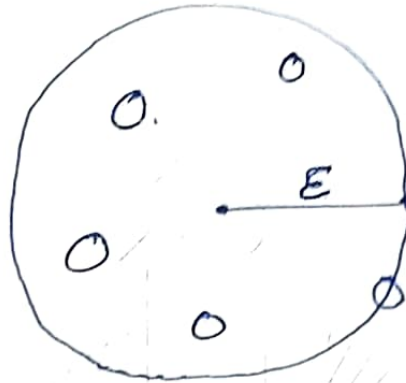


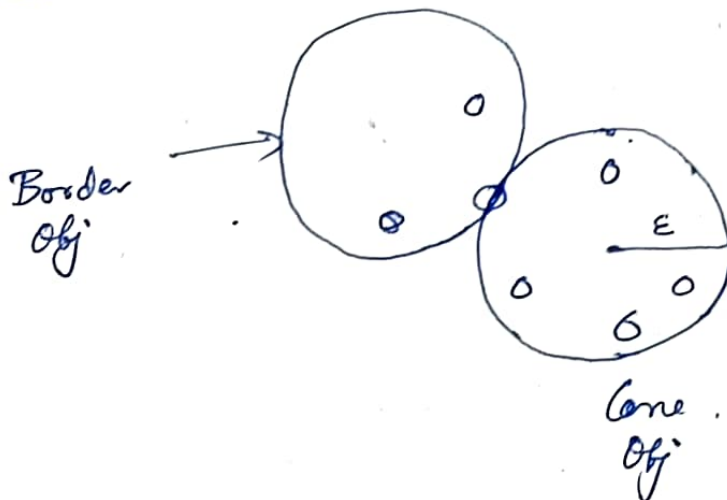
1) a) Core Object: Objects that are within the range of specified Radius ( $\epsilon$ ) & are greater than or equal to a certain number ( $\text{min\_pts}$ ) are known as Core objects.

E.g.  $\text{min\_pts} = 4$

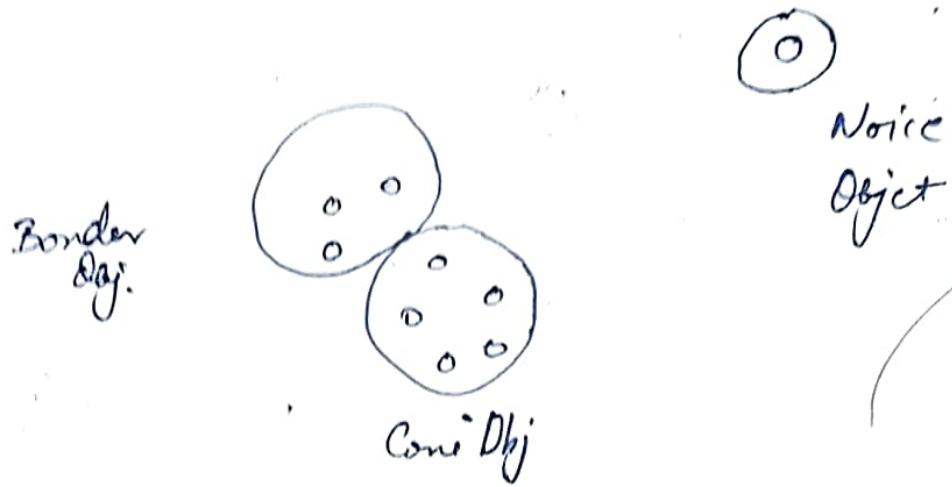


b) Border Object: Objects that are fewer than  $\text{min\_pts}$  & are at  $\epsilon$ -neighbourhood of core object are known as border objects.

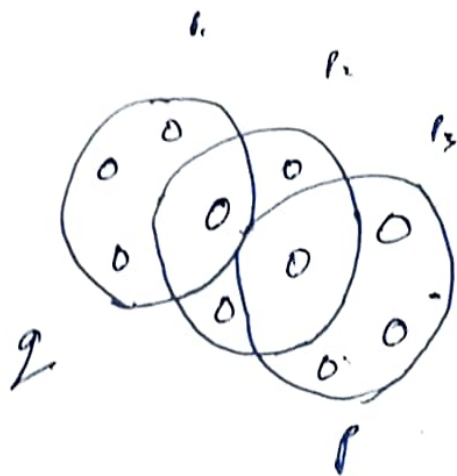
E.g.  $\text{min\_pts} = 4$



a) Noice Object: Object that are not Core objects or Border Objects are Noice Objects. They are far away from Clusters.



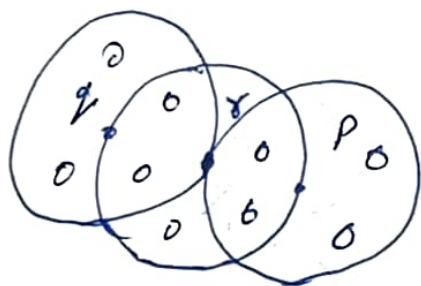
1) (ii) Density Reachability: An object  $q$  is said to be Density reachable to  $p$  if there exist  $p_1, p_2, \dots, p_n$ , where  $p_1$  is  $q$  &  $p_n$  is  $p$  that are Direct Density Reachable. In a series  $p_k$  is directly Density Reachable to  $p_{k-1}$ . Then  $p$  &  $q$  are called Density reachable.



hence  $p$  &  $q$  are density reachable.

1) a) ii) Density Connectivity: Object  $p$  &  $q$  are said to be Density connected if they are directly density connected to a common object  $o$ .

4 1/2



Hence  $x$  is direct density connected to both  $p$  &  $q$ .

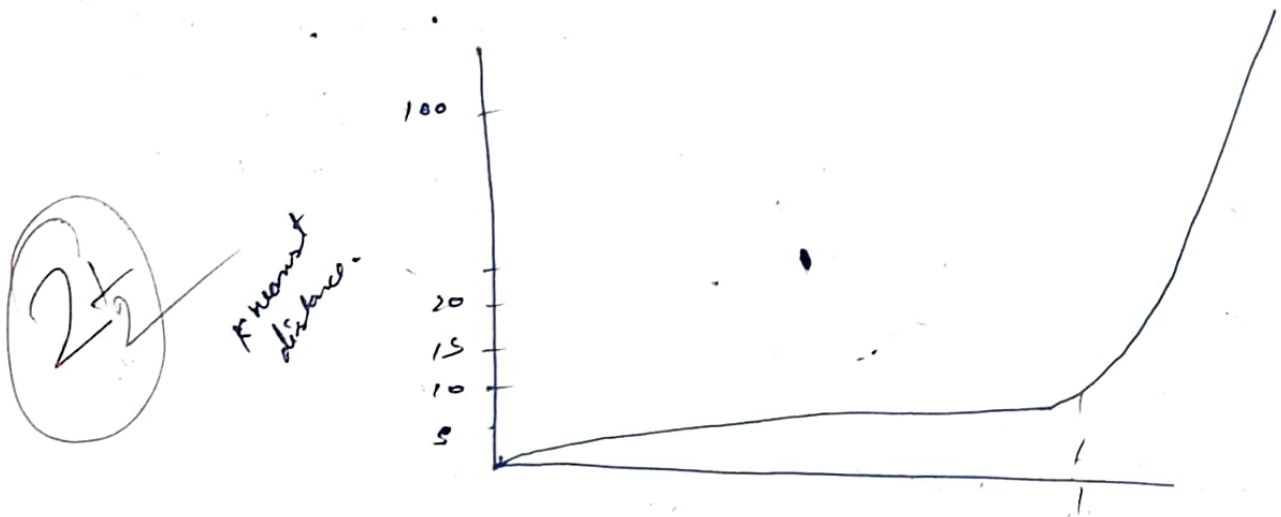
f) b) i) Maximality Condition: If  $q \in C_1$  & it is density reachable to  $p \in C_2$ , then both  $p, q \in C$ .

ie ~~the~~ Density Reachable clusters are combined to

2 1/2 from a single cluster, ~~as~~ until there is no  $q$  left that can be added to the Cluster.

ii) Connectivity Condition: If  $q$  is border object to the core object  $p \in C$ , then,  $q$  is combined to the cluster  $C$  ie  $q, p \in C$ : & the Noices that are far away are discarded.

1) c)  $Eps$  is specified radius  $E$ , &  $Minpts$  is the minimum number of points within that radius. The  $Eps$  &  $Minpts$  are determined using distance of  $k$ -nearest neighbours.



The sharp turn in graph decides the optimal distance radius  $E$ .  $Minpts$  value is decided using above. Higher value of  $Minpts$  will isolate many cluster & lower values of  $Minpts$  will make Noise a cluster. So, deciding  $E$  &  $Minpts$  is crucial part of Density-Based spatial Clustering of Application with Noise (DBSCAN).



2) Let  $I$  be the item set  $\{i_1, i_2, i_3, \dots, i_n\}$  &  $T$  be the transaction  $(t_1, t_2, t_3, \dots, t_n)$ .

Support: It is the <sup>Ratio</sup> number of <sup>transactions</sup>  $X$  occurred in to the total Transactions number.

$$\text{Sup}(X) = \frac{(X \cup Y) \cdot \text{count}}{n}$$

$\{n$  is total no. of transactions.

3

Confidence: Confidence is ratio of number of times  $X$  occur in different transaction to the Support of  $X$ .

$$\text{Conf}(X) = \frac{(X \cup Y) \cdot \text{count}}{X \cdot \text{count}}$$

3a) A  $k$ -item set is referred to a frequently <sup>itemset</sup> if it is found in at least  $k-1$  transactions.

3a) An association rule is referred as important if its confidence is 100%. A higher value of Confidence indicates high chance of association rule to be valid.

3)b) Min Support 30%,  
Min Confidence 80%,

Transactions

$T_1$

$T_2$

$T_3$

$T_4$

$T_5$

$T_6$

Items

(I<sub>1</sub>)

(I<sub>2</sub>)

Bread, Butter

(I<sub>3</sub>)

Bread, Milk, Butter

(I<sub>4</sub>)

Bread, Jelly, Butter

(I<sub>5</sub>)

Bread, Coke

Bread, Milk

Milk, Coke

Support  $I_1 = \frac{5}{6}$  ~~2/3~~

$$I_2 = \frac{3}{6}$$

$$I_3 = \frac{3}{6}$$

$$I_4 = \frac{1}{6}$$

$$I_5 = \frac{2}{6}$$

Confidence  $\{I_1, I_2\} : 3, \{I_1, I_3\} : 2, \{I_1, I_4\} : 1$

$\{I_1, I_5\} : 1, \{I_2, I_3\} : 1, \{I_2, I_4\} : 1$

$\{I_3, I_5\} : 1,$

Frequent one  $\{I_1, I_2\}, \{I_1, I_3\},$

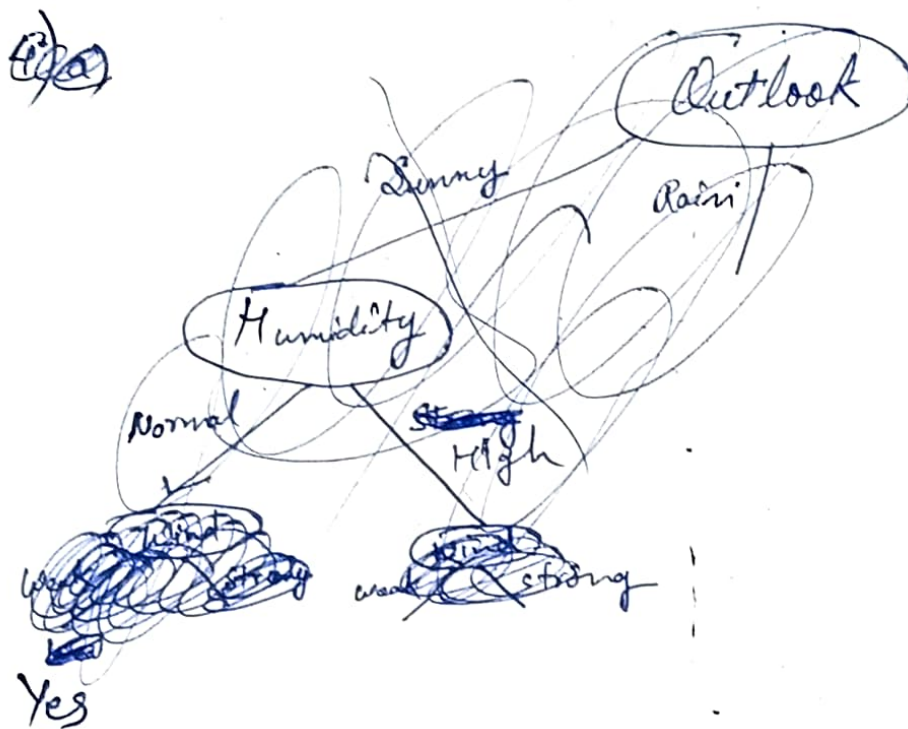
3-Items

$\{I_1, I_3, I_2\} : 1, \{I_1, I_4, I_2\} : 1$   $\{I_2, I_4\}$  are not frequent so, Not considered?

22

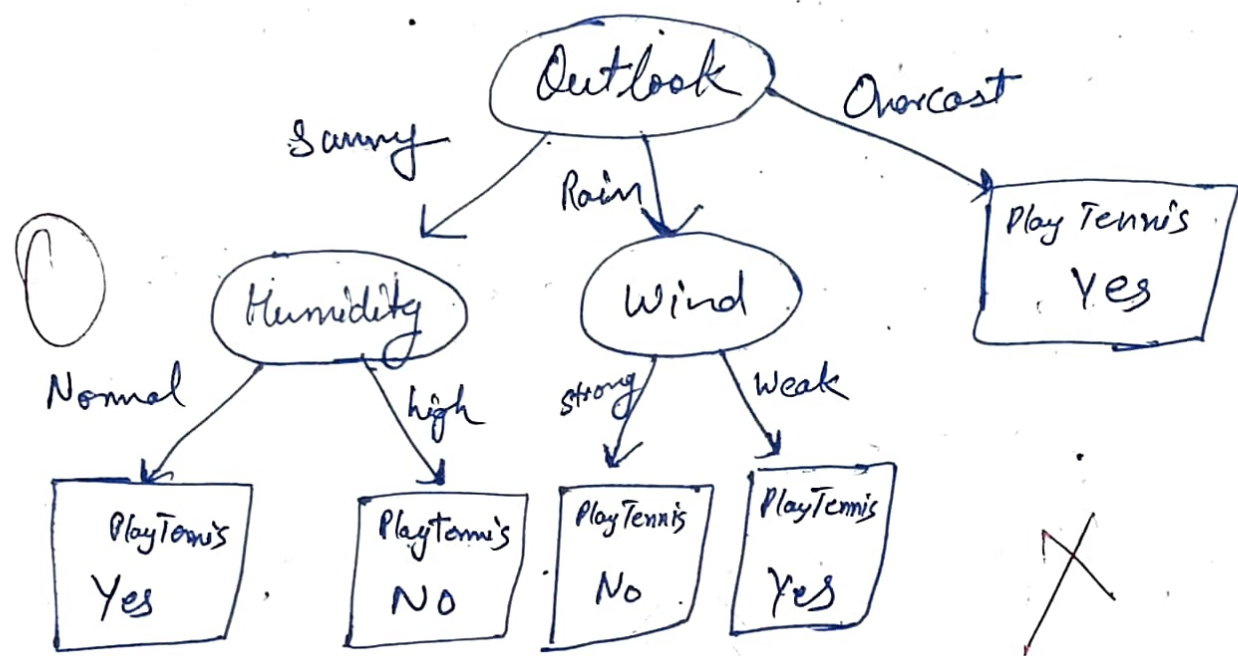
X	Y	Confidence
$I_1$	$I_2, I_3$	$\frac{1}{5}$
$I_2$	$I_1, I_3$	$\frac{2}{3}$
$I_3$	$I_1, I_2$	$\frac{3}{3}$
$I_2, I_3$	$I_1$	$\frac{1}{1}$
$I_1, I_3$	$I_2$	$\frac{1}{2}$
$I_1, I_2$	$I_3$	$\frac{3}{3}$

3)c) Drawback of a-priori algorithm is there  
can be different <sup>minimum</sup> support for different items  
In the item set, a-priori does not consider  
the multiple minimum support. To overcome  
this problem we use MS-priori algorithm.





4)a)



Root of a decision Tree is Outlook.

4)b) Possible Terminating criteria for Decision Tree

algorithm:

- If a node can decide the <sup>output</sup> class ~~off~~, then there is no need for further branching
- If all the decision nodes are used no further nodes can be made & only leaves will be the output class.
- If a decision node is ~~too~~ branched but their output is same, then the succer decision node can be removed.

4) c) Causes of Overfitting occurs when the cases for which the model was not trained gives wrong output or inconsistent decision nodes are formed.

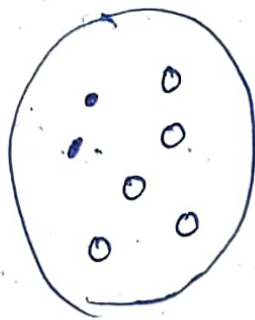
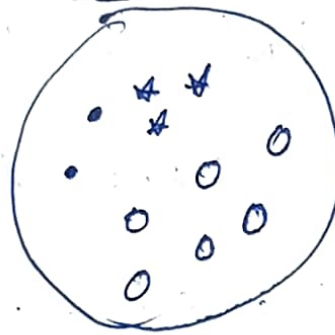
2) The overfitting can be solved in 2 ways

Pre Prunning: Removing the branch of tree that is not required & is not able to perform well classification.

Post Prunning: Adding a decision node in the tree that will ~~be~~ correctly ~~predict the output~~ of the classify the output.

2) a) i) Purity: is a measure of validation of the cluster formed.

$$\text{Purity} = \frac{\max(w_1, w_2, \dots, w_n)}{\text{Number of points in Cluster}}$$

C<sub>1</sub>C<sub>2</sub>

$$P(C_1) = \frac{\max(2, 5)}{7}$$

$$= \frac{5}{7}$$

$$P(C_2) = \frac{\max(2, 3, 6)}{11}$$

$$= \frac{6}{11}$$

3

2) ii) Rand Index :-

Predicted Output

		Predicted Output	
		C	~C
Ground Truth	C	TP	FN
	~C	FP	TN

Rand Index

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

$$TP + FN + FP + TN$$

TP & TN represent correctly classified clusters.

2) b) Inter Cluster Distance: Inter Cluster distance is the distance between the 2 clusters. A good classifying algorithm will try to increase Inter Cluster Distance.

Intra Cluster Distance: Intra Cluster distance is the distance between 2 points in the same cluster. A good classification algorithm will try to reduce the intra cluster distance.

These distances are calculated using different methodologies.

Inter Cluster Distances:

$$\text{Mean Cluster Distance} = \frac{\text{Summation of All possible inter cluster distances}}{\text{No. of Possible Cluster distances}}$$

$$\text{root mean Square} = \sqrt{\frac{\sum_{i \neq j}^k (\delta(x_i, x_j))^2}{N}}$$

Intra Cluster Distances:

$$\text{Mean Intra Cluster Distance} = \frac{\sum_{k=1}^n \Delta x_k}{N}$$

$$\text{root mean square Intra Cluster Distance} = \sqrt{\frac{\sum_{k=1}^n (\Delta x_k)^2}{N}}$$



2)c) Dunn's Cluster validation Index: let  $\delta(x_i, x_j)$  be the inter cluster distance &  $\Delta x$  be the intra cluster distance then Dunn's Index is defined as

$$\text{Dunn's Index} = \min_{1 \leq i \leq c} \left\{ \min_{1 \leq j \leq c} \left\{ \frac{\delta(x_i, x_j)}{\max_{1 \leq k \leq c} (\Delta x_k)} \right\} \right\}$$

A higher value of Dunn's index indicate ~~more~~ high inter cluster distance & low intra cluster distance. This is used in cluster evaluation to decide how well a cluster classification performed. A lower Dunn's Index will indicate low inter cluster distance & high intracuster distance which indicates the clusters formed are not well classified as the aim of a good reclassification Algorithm is to maximize intercluster distance & minimize intra-cluster distance.