

1)

a)

- i) A core object in DBSCAN algorithm with respect to ϵ and MinPts is referred to as the object which has more objects than MinPts ~~is~~ within the ϵ radius or the ϵ neighbourhood.

A border object in DBSCAN algorithm with respect to ϵ and MinPts is referred to as the object which has less objects than MinPts within the ϵ radius ϵ or the ϵ neighbourhood, but is a neighbour of a core object, i.e. is density-reachable from at least one core object.

(A)
1
2

An object that is neither a core object nor a border object is a noise object. A noise object is not density-reachable from any other core object.

- ii) The points p ~~are~~ meet ϵ and MinPts ~~are~~ said to be density-reachable if there exists a sequence of points p_1, p_2, \dots, p_n , where $p_1 = p$ and $p_n = q$, such that p_{i+1} is directly density reachable from p_i , for all of the points i .

Two points p and q are said to be density connected to each other meet ϵ and MinPts if there exists a point o such that both p and q are density-reachable from point o through a sequence.

- b) i) Maximality condition of the DBSCAN algorithm states that if $p \in C$, (a cluster) is from a set of points in D and q is another point $\in D$ which is density reachable from point C , then point $q \in C$.
- ii) i.e. the same cluster.
- 2) Connectivity condition of the DBSCAN algorithm states that if p and q are two points in D , such that $p, q \in C$, then the points p and q are density connected to each other in the same cluster.
- c) There are two parameters Eps and Minpts of the DBSCAN algorithm which can be determined by finding the value of K , or the no. of clusters.
- 3) For each object, the K -distance between the object and its K th neighbour, is also known as the K -distance. Now, for the core objects, the value of K -distance will be small if the cluster is small and the K th neighbour is inside the cluster. Also, for the border objects, the value of the K -distance will be very large. For each object, its K -distance is found out and then their values of the K -distances are plotted in a graph. A point in the graph where there is a sharp change in the value, or an inflection determines the eps , value for a particular K .

Min Pts

The value of k gives us the value of the ~~Min Pts~~ Min Dist for the algorithm. Several values of k are tried and a particular value of k is finally selected, such that it gives the proper number of clusters. If the value of k selected is very high, then a large number of noise objects will be treated as clusters. Also, if the value of k selected is very small, then a small compact cluster will be treated as noise. So, the value of k is chosen very carefully, so that the value of Min Pts is successfully found out.

- 2) as is purity refers to the maximum weight of an object (n_i) in the cluster i divided by the size of the cluster C . Mathematically, purity = $\frac{1}{n_i} \max(n_{ij})$, where $n_j \in C$.
- The value of purity is biased, because it attains the maximum value when the total number of clusters is n . For example, if there is a cluster which contains 5, 2 and 1 objects respectively, then the purity of the cluster $\text{cluster } C$ would be = $\frac{1}{8} \max \{ (5, 2, 1) \} = 5/8 = 0.625$

"> Rand-Index is an external validation criterion which assess how many samples are correctly clustered by the clustering algorithm, with the help of ground truth.

~~For example~~ Mathematically, RI (Rand Index) =

$$\frac{\text{No. of samples clustered correctly}}{\text{Total no. of samples present}}$$

For example, if for an example, the following table is created,

3	No. of samples as A in ground truth	No. of samples as B in the ground truth
No. of samples as A in the cluster	70	60
No. of samples as B in the cluster	50	30

Then the Rand Index for the following sample will be-

$$RI = \frac{70 + 30}{70 + 60 + 50 + 30} = \frac{100}{210}$$

$$= 0.476$$

by intercluster distance between two objects x_i and x_j is defined if the two objects belong to two different clusters. It is represented by $\delta(x_i, x_j)$.

Intracluster distance between two objects x_i and x_j is defined if the two objects belong to the same cluster. It is represented by $\Delta(x_i, x_j)$.

The two intercluster distances are as follows:-

- i) single linkage distance - It refers to the minimum distance between two objects belonging to two different clusters.

$$\text{Mathematically, } \delta_1(x_i, x_j) = \min_{x_i \in S, x_j \in T} \{ \text{dist}(x_i, x_j) \}$$

ii) complete linkage distance - It refers to the distance between the most two most remote objects belonging to two different clusters. Mathematically, $\delta_2(x_i, x_j) = \max_{x_i \in S, x_j \in T} \{ \text{dist}(x_i, x_j) \}$

The two intracluster distances are as follows:-

- i) complete diameter distance - It refers to the maximum distance between two objects belonging to the same cluster.

$$\text{Mathematically, } \Delta_1(x_i, x_j) = \max_{x_i, x_j \in S} \{ \text{dist}(x_i, x_j) \}$$

ii) Average diameter distance - It refers to the average of all the distances between objects that are present in the same cluster.

Mathematically,

$$\Delta_2(X_k) = \frac{1}{|S| |S|-1} \sum_{\substack{x_i, x_j \in S \\ i \neq j}} \text{dist}(x_i, x_j)$$

c) The Dunn's Index or D-Index for cluster validation is defined as

$D\text{Index}(U) = \frac{\min_{i=1}^n \left\{ \min_{\substack{j=1 \\ j \neq i}}^n d(x_i, x_j) \right\}}{\max \{ \Delta(X_k) \}}$

$D\text{Index}(U) = \min_{i=1}^n \left\{ \min_{\substack{j=1 \\ j \neq i}}^n d(x_i, x_j) \right\} / \max \{ \Delta(X_k) \}$

The Dunn's Index is an important index as it compares both the intercluster and intracluster distances. Here $d(x_i, x_j)$ is the intercluster distance and $\Delta(X_k)$ denotes the intracluster distance. The main aim of Dunn's index is to maximise the intercluster distance and minimise the intracluster distance, so that the compactness and connectedness is more between two clusters and the spatial separation is more between the two clusters. The more the value of the Dunn's Index, the more the cluster quality. Thus Dunn's Index helps in the usefulness of cluster evaluation.

3) The support of an association rule of the form $X \rightarrow Y$, is denoted by $(X \cup Y)$ -count where if ~~it~~ $(X \cup Y)$ appears in the transaction divided times by the total no. of transactions. If the no. of transactions is n , then support = $\frac{(X \cup Y)$ -count}{n}

The confidence of an association rule of the form $X \rightarrow Y$, is the no. of transactions where if X occurs, then Y will occur.

Mathematically, confidence $(X \rightarrow Y) = \frac{(X \cup Y)$ -count}{X-count}

An item set is referred to as a frequent item-set if the support of the itemset is greater than the minimum support, i.e. $\text{support}(X) \geq \text{minSup}$.

An association rule is referred to as an important rule if it follows these three conditions-

- i) confidence $(X \rightarrow Y) \geq \text{minConf}$
- ii) $\text{Support}(X \rightarrow Y) = \text{support}(X \cup Y) = \text{support}(Z)$
- iii) $\text{confidence}(X \rightarrow Y) = \frac{\text{support}(X \cup Y)}{\text{support}(X)}$

b) The transactions defined are as follows-

Transactions

T1

Items

Bread, Butter

T2

Bread, MILK, Butter

T3

Bread, Jelly, Butter

T4

Bread, Coke

T5

Milk, Coke

T6

According to a priori algorithm, we have to find the frequent item-sets.
Now let $C_1 = \{\text{Bread}\}, \{\text{Butter}\}, \{\text{Milk}\}, \{\text{Jelly}\}, \{\text{Coke}\}$

$$\text{Now, support}(\text{Bread}) = \frac{5}{6} = 0.83$$

$$\text{support}(\text{Butter}) = \frac{3}{6} = 0.5$$

$$\text{support}(\text{Milk}) = \frac{3}{6} = 0.5$$

$$\text{support}(\text{Jelly}) = \frac{1}{6} = 0.17$$

$$\text{support}(\text{Coke}) = \frac{2}{6} = 0.33$$

$$\text{Now, minSup} = 30\% = 0.3.$$

So, $F_1: \{\text{Bread}\}, \{\text{Butter}\}, \{\text{Milk}\}, \{\text{Coke}\}$

Now, let $C_2: \{\text{Bread, Butter}\}, \{\text{Bread, Milk}\}, \{\text{Bread, Coke}\}, \{\text{Bread, Butter, Milk}\}, \{\text{Butter, Coke}\}, \{\text{Milk, Coke}\}$

Now, support (Bread Butter) = $\frac{3}{6} = 0.5$

support (Bread, Milk) = $\frac{2}{6} = 0.33$

support (Bread, Coke) = $\frac{1}{6} = 0.17$

support (Butter, Milk) = $\frac{1}{6} = 0.17$

support (Butter, Coke) = 0

support (Milk, Coke) = $\frac{1}{6} = 0.17$

support (Bread, Butter, Milk) = ?

So, $F_2: \{ \text{Bread, Butter}, \{ \text{Bread, Milk} \} \}$

Now, $C_3: \{ \text{Bread, Butter, Milk} \} \neq \emptyset$

(No other candidates are generated from F_2).

So, the frequent itemset = $\{ \text{Bread, Butter} \}$, and

$\{ \text{Bread, Milk} \}$.

Now, for the generation of important association rules, we may have the following -

i) Bread \rightarrow Butter, confidence = $\frac{3}{5} = 60\%$

ii) Butter \rightarrow Bread, confidence = $\frac{3}{3} = 100\%$

iii) Bread \rightarrow Milk, confidence = $\frac{2}{5} = 40\%$

iv) Milk \rightarrow Bread, confidence = $\frac{2}{3} = 66.6\%$

Now, minimum confidence = 80%.

So only association rule satisfies this condition.

So, the ~~most~~ important association rule is -

Butter \rightarrow Bread.

c) if there are m ^{items} ~~instances~~ available then the total number of transaction rules could be of the order, $O(2^m)$. Since this ~~is~~ grows exponentially ~~is~~ the cost of computational and memory requirements are very high. The apriori algorithm exploits the sparseness of the data and high support and confidence. Still the apriori algorithm generates thousands, and even millions of association rules. This is a major drawback of apriori algorithm as the computational cost and cost of memory becomes very high.

Q) as the decision tree model trained by the dataset is as follows-

```
graph TD; Outlook([Outlook]) -- Sunny --> SunnyNode(( )); Outlook -- Overcast --> OvercastNode(( ));
```

4) a) The dataset is given as follows

Outlook	Humidity	Wind	Play Tennis (Class Variable)
Sunny	Normal	Strong	Yes
Sunny	High	Weak	No
Overcast	High	Strong	Yes
Rain	Normal	Strong	No
Rain	High	Weak	Yes
Sunny	Normal	Weak	No
Sunny	High	Strong	No
Rain	High	Strong	Yes
Rain	Normal	Weak	Yes
Overcast	Normal	Weak	Yes

Total no of samples, $N = 10$

Now, the information required would -

$$\text{Info (I)} = \frac{-p}{p+n} \log \left(\frac{p}{p+n} \right) - \frac{n}{p+n} \log \left(\frac{n}{p+n} \right)$$

$$\begin{aligned}
 &= -\frac{6}{10} \log_2 \left(\frac{6}{10} \right) - \frac{4}{10} \log_2 \left(\frac{4}{10} \right) \\
 &= -\frac{3}{2} \log_2 \left(\frac{3}{2} \right) - \log_2 \left(\frac{1}{2} \right) \\
 &= \frac{3}{2} (\log_2(3) - \log_2(2)) - 1 \times 0 \\
 &= -\frac{3}{2}
 \end{aligned}$$

$$\begin{aligned}
 &= -\frac{6}{10} \log_2 \left(\frac{6}{10}\right) - \frac{4}{10} \log_2 \left(\frac{4}{10}\right) \\
 &= -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) \\
 &= -\frac{3}{5} (\log_2(3) - \log_2(5)) - \frac{2}{5} (\log_2(2) - \log_2(5)) \\
 &= -\frac{3}{5} (1.5 - 2.3) - \frac{2}{5} (1 - 2.3) \\
 &= \frac{3}{5} \times 0.7 + \frac{2}{5} \times 1.3 = 0.42 + 0.52 = 0.94
 \end{aligned}$$

Now, for each of the attributes, we consider the information and information gain.

i) Outlook

$$\text{Info}_{\text{Outlook}}(D) = \sum_{v=1}^{\infty} \frac{|D_v|}{|D|} \text{Info}(D_v)$$

$$\begin{aligned}
 \text{Info}(\text{Sunny}) &= -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) \\
 &= \frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) \\
 &= \frac{1}{2} \log_2 2 + \frac{1}{2} \log_2 2 = \log_2 2 = 1.
 \end{aligned}$$

$$\text{Info}(\text{Rain}) = -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right) = 1.$$

$$\begin{aligned}
 \text{Info}(\text{Overcast}) &= -\frac{2}{2} \log_2 \left(\frac{2}{2}\right) - 0 \\
 &= 0.
 \end{aligned}$$

$$\begin{aligned}
 \text{So, Info}_{\text{Outlook}}(D) &= \frac{4}{10} \times 1 + \frac{4}{10} \times 1 + \frac{2}{10} \times 0 \\
 &= \frac{8}{10} = 0.8
 \end{aligned}$$

~~∴~~ Gain(Outlook) = 0.14

ii) Humidity -

$$\text{Info (Normal)} = -\frac{4}{5} \log_2(4/5) - \frac{1}{5} \log_2(1/5)$$

$$= -\frac{4}{5} (\log_2(4) - \log_2(5)) - \frac{1}{5} (\log_2(1) - \log_2(5))$$

$$= -\frac{4}{5} \times 2 (2 - 2.3) - \frac{1}{5} (0 - 2.3)$$

$$= \frac{4}{5} \times 0.3 + \frac{1}{5} \times 2.3 = 0.24 + 0.46 = 0.7$$

$$\text{Info (High)} = -\frac{2}{5} \log_2(2/5) - \frac{3}{5} \log_2(3/5) = 0.94$$

$\therefore \text{Info Humidity} \quad (1) = \frac{5}{10} \times 0.7 + \frac{5}{10} \times 0.94$

$$= 0.35 + 0.47 = 0.82$$

$$\therefore \text{Gain (Humidity)} = 0.94 - 0.82 = 0.12$$

iii) Wind -

$$\text{Info (Strong)} = -\frac{2}{5} \log_2(2/5) - \frac{3}{5} \log_2(3/5) = 0.94$$

$$\text{Info (Weak)} = -\frac{4}{5} \log_2(4/5) - \frac{1}{5} \log_2(1/5)$$

$$= 0.7$$

$$\therefore \text{Info Wind} \quad (1) = \frac{5}{10} \times 0.94 + \frac{5}{10} \times 0.7 = 0.35 + 0.47$$

$$= 0.82$$

$$\therefore \text{Gain (Wind)} = 0.94 - 0.82 = 0.12$$

So, among the 3 attributes, gain of attribute ^{Outlook} is the highest, so the root of the decision tree will be the attribute ^{Outlook}.

b) The possible terminating criteria of a decision tree algorithm are -

- (3) i) When all of the nodes of the decision tree have the same class label. The decision tree stops growing and the leaf node is assigned the value of the class node.
- ii) When none of the attributes are left for splitting a node into further nodes. In such cases, a majority decision is taken for assigning the class label to the leaf node.
- iii) When all of the data samples are exhausted. In such cases each leaf node is assigned a class label according to the value of the node.
- c) Over-fitting happens when a model tries to fit all of the data including noise and outliers into itself, resulting in very low error rates in the training data set but high error rates in the unseen data.
- (3) Over-fitting also occurs when all of the training data is used for model construction resulting in very low accuracy for the test data.

In decision trees, the problem of overfitting can be solved in two ways-

i) Pre-pruning- In pre-pruning, the growth of the tree is stopped, if a split to the node does not increase the goodness measure beyond a certain threshold. This threshold is generally determined by the user.

ii) Post pruning- In post pruning, the decision tree is let to fully grow and then those branches are removed which reflects the noise and outliers in data. In such cases, the decision tree removes the more complex decisions which may be influenced by noisy data.