

(a)

(i) Core Object

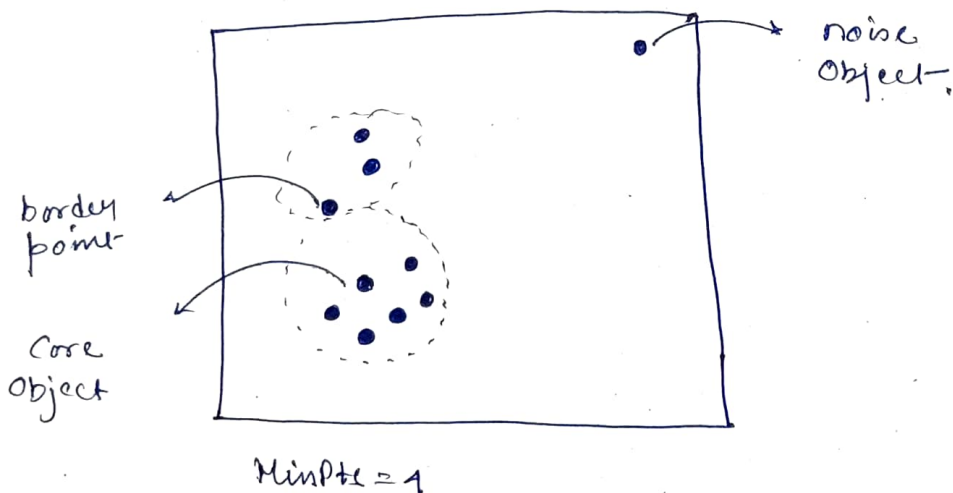
For a certain value of $MinPts$ and ϵ , a point 'p' will be a core object if there are more or equal number of points as of $MinPts$ within a radius of ϵ of the point p.

Border Object.

For a specific $MinPts$ and ϵ , a border object will be having less no. of points than $MinPts$ within a radius of ϵ , but the point is a neighbor of a core object.

Noise Object.

Neither a Core nor a border object. Noise objects lie far away from core objects or border objects.



(ii) Density reachability.

A point 'q' will be density reachable from point 'p' if —

there exist some points p_1, p_2, \dots, p_n

where $p_1 = p$ and $p_n = q$

then p_2 is directly-density reachable from p_1

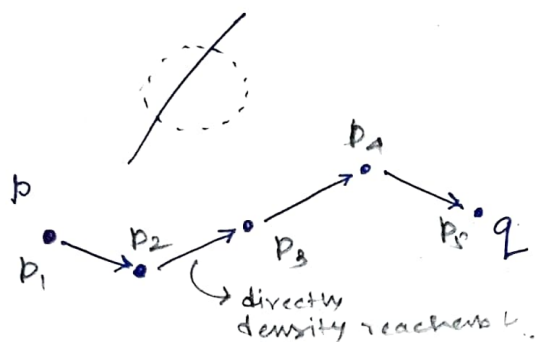
p_3 is " " " " p_2

\vdots

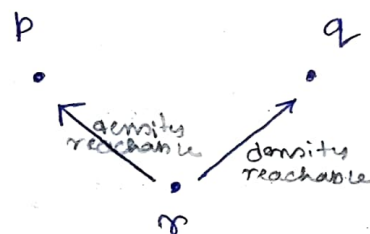
$p_n(q)$ is " " " " p_{n-1}

Density Connectivity.

A point 'q' will be density connected to point 'p' if there exists another point 'r' and both 'p' and 'q' are density reachable from point 'r'.



Density
reachable



Density
Connectivity

(b)

(i) Maximality Condition.

3 It states that if a point 'p' is within a cluster 'c' (i.e. $p \in c$), and another point 'q' which is density reachable from 'p' then point 'q' is also within the cluster i.e. $q \in c$.

(ii) Connectivity Condition.

It states that if two points 'p' and 'q' are within a cluster c, then 'p' & 'q' will be density connected to each other.

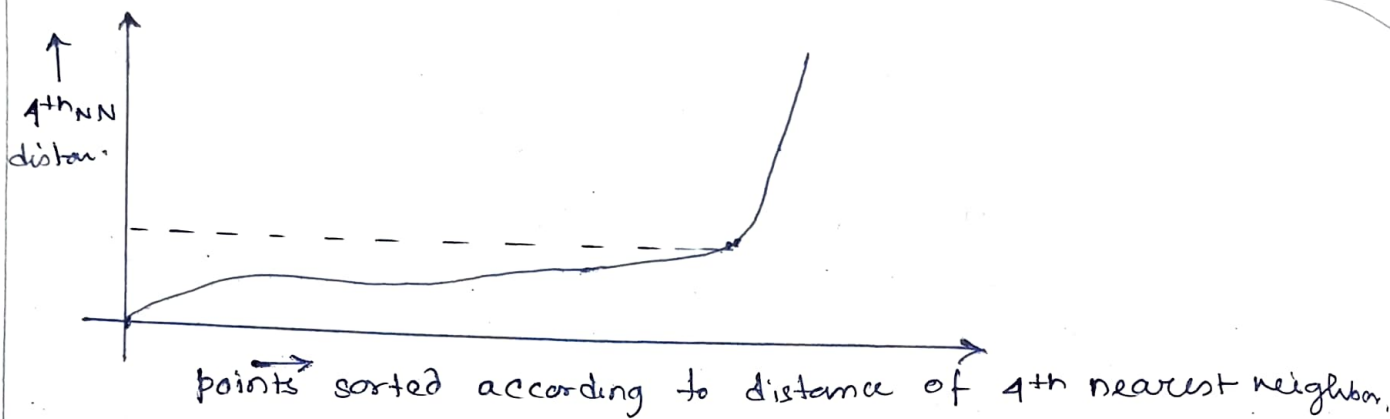
(c)

The main idea of determining Eps and MinPts in DBSCAN algo is —

The k^{th} nearest neighbors of each ~~for~~ object will be approximately at the same distance.

3 For noisy data, the k^{th} nearest neighbors will be at a further ~~for~~ distance.

Hence, for a particular value of 'k', we plot the graph of sorted points according to distance of ~~the~~ k^{th} nearest neighbor with ~~it~~ the k^{th} nearest neighbor distance.



Let's assume $k \geq 4$.

There will be a sharp point in the plot due to the noisy data.

At that point the value of ' k ' can be considered as $MinPts$ and the k^{th} nearest neighbor will be considered as Eps .

2.

(a)

(i) Purity.

purity is defined as —

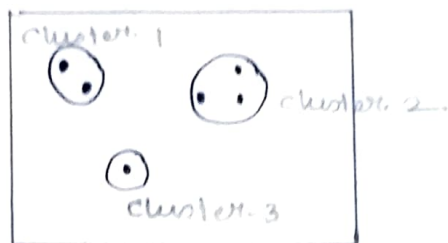
$$\text{purity} = \frac{\text{maximum number of points in a cluster}}{\text{total number points (including all clusters)}}$$

for example, there are total 6 points.

2 points form cluster-1

3 points form cluster-2

1 point form cluster-3



the purity will be = $\frac{\max(2, 3, 1)}{6} = \frac{3}{6} = \frac{1}{2}$

(ii) Rand Index.

Here we compare the clusters formed by the algorithm with the ground truth data.

There can be 4 different possibilities.

Prediction.

		<u>Prediction.</u>	
		Same Class in Clustering	Different Class in Clustering
<u>Ground Truth.</u>	Same class in ground truth	A	B
	Different class in ground truth	C	D

Rand Index = $\frac{A + D}{A + B + C + D}$

(b)

Intercluster distance.

It refers to the distance between two different clusters. It is denoted by ' δ '. There are several ways to calculate intercluster distance —

(i) Single linkage distance.

For any two clusters S, T ; it refers to the minimum distance between two points of S and T .

$$\delta = \min(d(x, y))$$

$x \in S, y \in T$, $d(x, y)$ distance between points x & y .

(ii) Complete linkage distance.

For any two clusters S, T ; it refers to the maximum distance between two points in S and T .

$$\delta = \max(d(x, y))$$

$x \in S, y \in T$, $d(x, y)$ distance between points x & y .

Intra-cluster distance.

It refers to the distance between two points within the same cluster. It is denoted by ' Δ '. There are several ways to calculate intra-cluster distance.

(i) Complete diameter distance.

It is the farthest distance of any two nodes within the same cluster, S .

$$\Delta = \max (d(x, y))$$

$x, y \in S$ and $d(x, y)$ distance between two points x and y .

(ii) Average diameter distance.

It is the average ^{distance} of all points with other points.

$$A = \frac{\sum d(x, y)}{|S|(|S|-1)}$$

$x, y \in S$ but $x \neq y$.

$|S|$ is the ~~see~~ number of points in cluster S .

$d(x, y)$ distance between two points x & y .

(c)

Dunn's Index

$$= \max_{1 \leq i \leq c} \left\{ \max_{\substack{1 \leq j \leq c \\ j \neq i}} \left\{ \frac{d(x_i, x_j)}{\max_{1 \leq k \leq c} (\Delta(x_k))} \right\} \right\}$$

$d(x_i, x_j)$ is the inter-cluster distance between two clusters x_i & x_j .

$\Delta(x_k)$ is intra-cluster distance of the cluster x_k .

A cluster will be 'good' cluster if its intra-cluster distance is less and inter-cluster distance with other clusters is more.

3.2 Hence, in Dunn's index, the aim will be to ^{increase} ~~decrease~~ $\delta(x_i, x_j)$ and decrease $\Delta(x_k)$.

Hence, if the Dunn's index value is more, the cluster is good cluster, the more Dunn's index value, better the cluster is.

3.

(a)

support.

~~At~~ In any set of transactions 't', an itemset $\{x, y\}$.

The support of the itemset will be

$$= \frac{\text{number of transactions where } x, y \text{ appears.}}{\text{total number of transactions.}}$$

$$= \frac{(x \cup y). \text{count}}{n}$$

[$n \rightarrow$ total no. of transaction]

confidence.

In any set of transactions, 't', an association rule $X \rightarrow Y$ is there.

The confidence of the rule will be

(4)

$$= \frac{\text{no. of transactions where } \{X, Y\} \text{ appears.}}{\text{no. of transaction where } X \text{ appears.}}$$

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

■ For a specific minsup (minimum support) and minconf (minimum confidence) value, when the support of itemset will be greater than the minsup value, it will be referred to as a frequency item set. Similarly, for an association rule, if its confidence value is greater than the ~~min~~ minconf, then the association rule will be an important rule.

(b) Let, Bread = I_1 , Butter = I_2 , Milk = I_3 , Jelly = I_4
Coke = I_5 .

Transactions.

Items

T_1

I_1, I_2

T_2

I_1, I_3, I_2

T_3

I_1, I_4, I_2

T_4

I_1, I_5

T_5

I_1, I_3

T_6

I_3, I_5

minsup = 30%.

minconf = 80%.

3) $c_1: \{I_1\} = 5, \{I_2\} = 3, \{I_3\} = 3, \{I_4\} = 1, \{I_5\} = 2$

support $\{I_1\} = \frac{5}{6} > 30\%$.

support $\{I_3\} = \frac{3}{6} > 30\%$.

support $\{I_2\} = \frac{3}{6} > 30\%$.

support $\{I_4\} = \frac{1}{6} < 30\%$.

support $\{I_5\} = \frac{2}{6} > 30\%$.

$\therefore F_1: \{I_1\}, \{I_2\}, \{I_3\}, \{I_5\}$

$c_2: \{I_1, I_2\} = 3, \{I_1, I_3\} = 2, \{I_1, I_5\} = 1$

$\{I_2, I_3\} = 1, \{I_2, I_5\} = 0, \{I_3, I_5\} = 1$

support $\{I_1, I_2\} = \frac{3}{6} > 30\%$.

support $\{I_2, I_3\} = \frac{1}{6} < 30\%$.

support $\{I_1, I_3\} = \frac{2}{6} > 30\%$.

support $\{I_2, I_5\} = 0 < 30\%$.

support $\{I_1, I_5\} = \frac{1}{6} < 30\%$.

support $\{I_3, I_5\} = \frac{1}{6} < 30\%$.

$\therefore F_2: \{I_1, I_2\}, \{I_1, I_3\}$

$\therefore C_3 = \{I_1, I_2, I_3\}$
 but this cannot be considered as $\{I_2, I_3\} \notin F_2$.

$\therefore C_3 = \emptyset$.

Hence no important rules can be generated.

Q. Drawbacks of a-priori algorithm.

- setting the minsup, minconf value incorrectly may not lead to an optimal solution.
- The time taken to train the model is relatively higher.

1.

(a)

Total Yes (p) = 6

Total No (n) = 4

Entropy of entire table -

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

$$= - \frac{6}{10} \log \left(\frac{6}{10} \right) - \frac{4}{10} \log \left(\frac{4}{10} \right)$$

$$= 0.971$$

$$E(A) = \sum \frac{p_i + n_i}{p+n} I(p_i, n_i)$$

Outlook	p_i	n_i	$I(p_i, n_i)$
Sunny	2	2	1
Overcast	2	0	0
Rain	2	2	0.971 1

$$\therefore E(\text{outlook}) = \frac{4}{10} \times 1 + 0 + \frac{3}{10} \times 1 = 0.7$$

<u>Humidity</u>	<u>p_i</u>	<u>n_i</u>	<u>$I(p_i, n_i)$</u>
Normal	4	1	0.722
High	2	3	0.971

$$\therefore E(\text{Humidity}) = \frac{5}{10} \times 0.722 + \frac{5}{10} \times 0.971$$

$$= 0.8465$$

<u>Wind.</u>	<u>p_i</u>	<u>n_i</u>	<u>$I(p_i, n_i)$</u>
Strong	2	3	0.971
Weak	4	1	0.722

$$E(\text{Humidity}) = \frac{5}{10} \times 0.971 + \frac{5}{10} \times 0.722$$

$$= 0.8465$$

$$\therefore \text{Information Gain (outlook)} = I(p, n) - E(\text{outlook})$$

$$= 0.971 - \cancel{0.6005} 0.7$$

$$= \cancel{0.2008} 0.271$$

$$\text{Similarly, Gain (Humidity)} = 0.971 - 0.8465 = 0.1245$$

$$\text{Gain (Wind)} = 0.971 - 0.8465 = 0.1245$$

\therefore Outlook has highest information gain.
It will be the root node.

(b) possible terminating Criteria.

3 (i) All samples of a given node belong to the same class.

(ii) There is no remaining attribute for further partitioning.

(iii) There is no more samples left.

(c) Model Overfitting may occur due to several reasons —

(i) The model can become more complex than required.

~~(ii)~~ The model can train noisy data or outliers which will lead to a overfitting scenario.

■ In decision-tree we solve overfitting by mainly two approaches —

(i) Pre-pruning

2/2 Here, we don't let the tree to grow at its fullest. At a threshold level we prune the tree,

(ii) Post-pruning

Here, we let the ~~good~~ tree to grow at its fullest & then start the pruning.