

Q.
(a.)

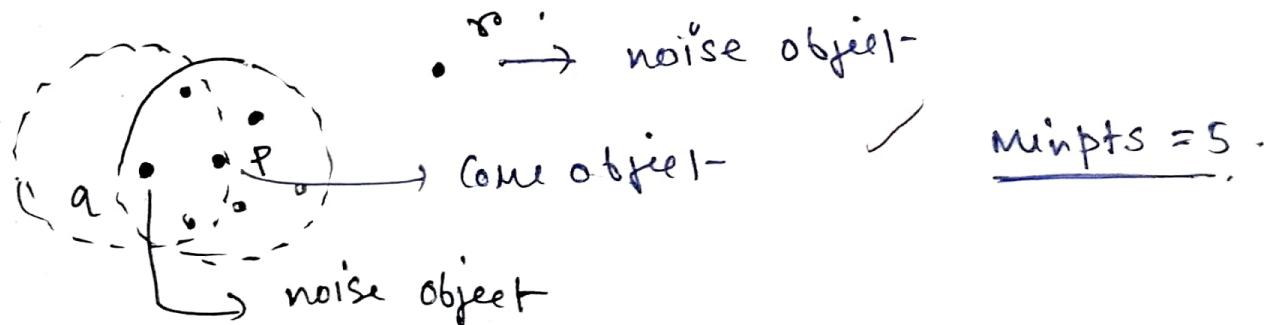
(i) Core object :- An object in an object-database is defined to be core object if there exist at least Minpts number of points within its ϵ -neighborhood.

where, Minpts and ϵ are the parameters (taken as input) in DBSCAN clustering algorithm.

(ii) Border object :- An object in an object-database is defined to be border object if there are less than Minpts number of ~~other~~ objects within its ϵ -neighborhood.

(iii) Noise object : \rightarrow In DBSCAN clustering an object is defined as noise object if neither it is a core object nor a border object.

The core object, border object and noise objects are shown below: — (w.r.t Minpts and ϵ) .

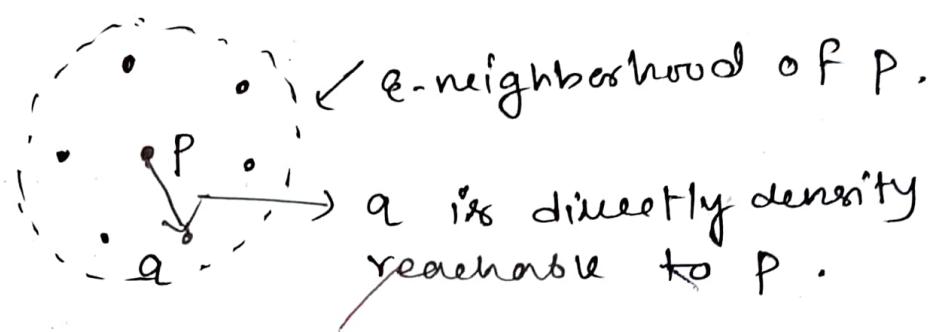


(ii) Density reachability :-

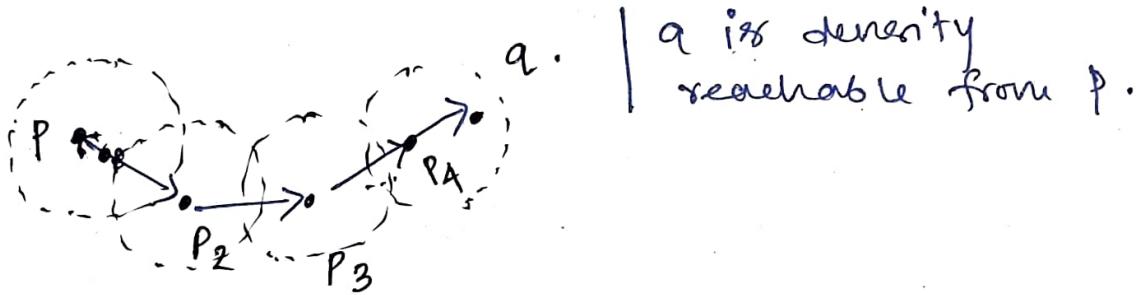
In DBSCAN clustering density reachability is defined as how in terms of reaching from an object (P) to another object (a) directly or via some other paths (indirectly).

It can be divided in two parts —

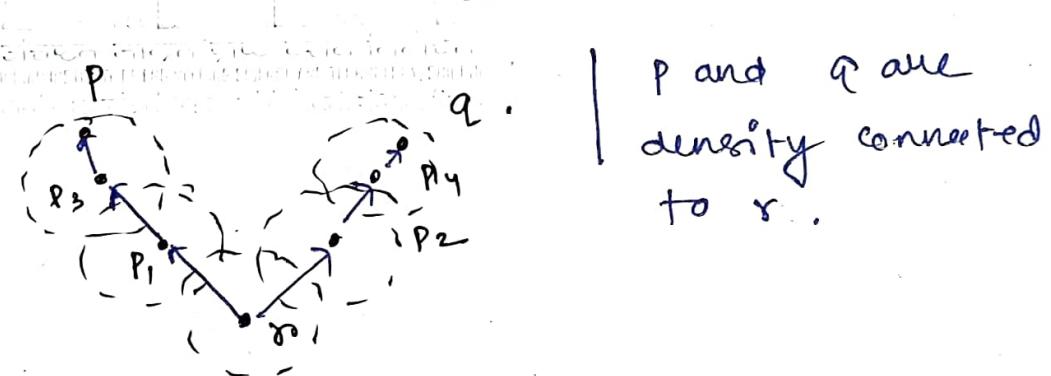
(i) Directly density reachable :— A point q is said to be directly density reachable from another point p if p is a core object and q is within the ϵ -neighborhood of p .



(ii) Density reachable :— A point q is said to be density reachable from another point p , w.r.t ϵ and minpts, if there exist a chain of points p_1, p_2, \dots, p_n such that $p_1 = p$ and $p_n = q$ and p_{i+1} is directly density reachable from p_i .



Density Connectivity :- An object q is density w.r.t ϵ and minpts connected to another object p , if there exist another object r such that both p and q are density reachable from r .



Maximality Condition :- The maximality condition in DBSCAN algorithm states that if there exist an object p which belongs to some cluster C , that is $p \in C$, then all other points that are density reachable from p also belongs to C , that is if q is in ϵ -neighborhood of p , ~~and~~ then $q \in C$.

(ii) Connectivity condition :- If there exist a cluster C_i , such that two object P and Q is in C_i , that is, $P, Q \in C_i$ then P and Q are density connected.

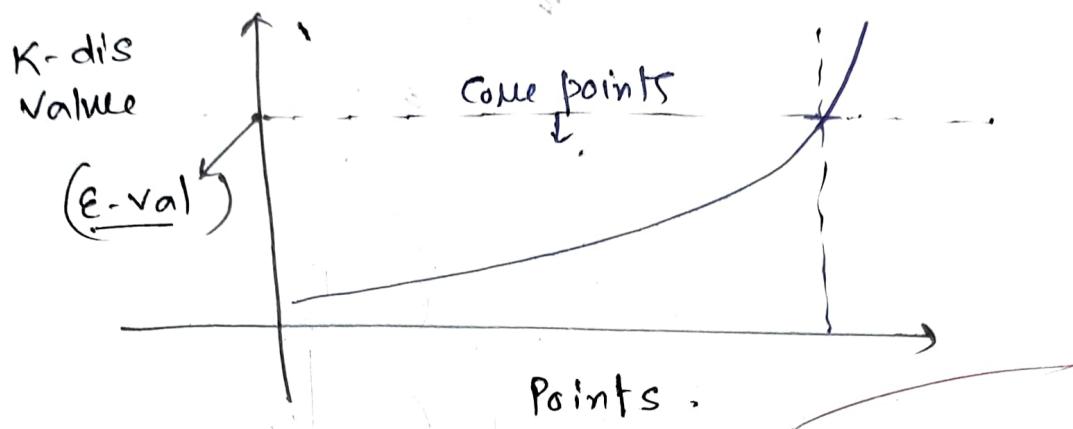
③ Parameters of DBSCAN algorithm :-

(*) In order to determine the parameters of the DBSCAN algorithm, we will use the concept of K^{th} nearest neighbor. It is assumed that the K^{th} nearest neighbor of every points within a cluster are at a same distance. Let's call it $k\text{-dist}$.

If so, $k\text{-dist}$: K^{th} nearest neighbor distance of any points.

- (i) First, ~~take~~ select a K value.
- (ii) find the K^{th} nearest neighbor distance of all points.
- (iii) sort these points based on the K^{th} nearest neighbor distance.
- (iv) plot these distances.
- (v) Carefully, examine the point where $k\text{-dist}$ value drastically changes.

(ii) The K-dist value of that point will be considered as the Eps (ϵ) value and the initial K value will be considered as Minpts.



②

(i)

Purity :— The purity of a cluster is defined as the ratio of the number of points in dominant class to the total number of points.

Mathematically the purity can be defined,

$$\text{Purity } (c_i) = \frac{1}{n_i} \max (n_{ji})$$

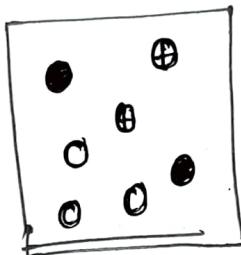
↓

Purity of the
ith cluster.

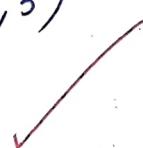
where, $n_i^o \rightarrow$ total number of point in the cluster
 $n_{ji}^o \rightarrow$ number of points that belongs to
 j^{th} class in the c_i^{th} cluster.

For example :-

consider the following cluster,



$$\text{Purity}(c) = \frac{1}{7} \max(2, 2, 3)$$
$$= \frac{3}{7}$$



(ii)

Rand Index :- The Rand Index is defined as the ratio of the number of points ~~in~~ in the correct cluster ~~with~~ with the total no of points in the cluster.

3

It is denoted by the following matrix,

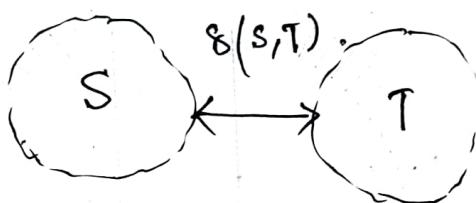
Number of Points	Points that belongs to cluster	Points that belongs to different cluster
Points that show same class labels	A	C
Points that show different class labels	B	D

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

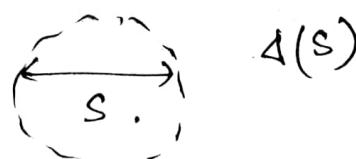


(b) Intercluster distance :- The intercluster distance is the distance between any two objects belonging to two different clusters. High value of the distance means better clusters. As shown below,

Suppose, there exist two clusters S and T. The distance between these two clusters are called intercluster distance denoted by $\delta(S, T)$.



→ Intracluster distance :- The intracluster distance is the distance ~~between~~ of the cluster itself. It denotes the compactness of a cluster, i.e. low value of the distance indicates more compactness hence good quality cluster.



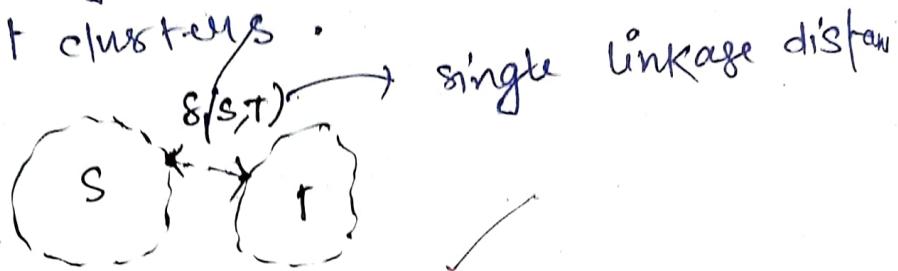
→ In both the cases, the distance that are typically used are given below —

- (i) Euclidean distance
- (ii) Chebyshew distance
- (iii) Manhattan distance etc.

PODS

□ Two intercluster distances : -

(i) single linkage distance : - The single linkage distance is the minimum distance between any two objects belonging two different clusters.



The equation is as follows,

$$\text{single linkage distance, } d_1(S, T) = \min \left\{ \begin{array}{l} d(x, y) \\ x \in S \\ y \in T \end{array} \right\}$$

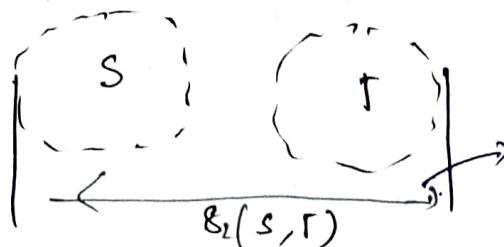
where, $S, T \rightarrow$ two different clusters.

$d(x, y) \rightarrow$ distance from x to y .

$x \in S$

$y \in T$.

(ii) complete linkage distance : - the complete linkage distance is the maximum distance of any two objects belonging two different clusters.



complete linkage distance .

The equation is as follows,

$$S_2(S, T) = \max \left\{ d(x, y) \middle| \begin{array}{l} x \in S \\ y \in T \end{array} \right\}$$

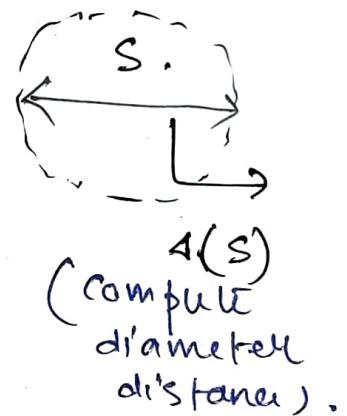
distance between
x and y objects.

compute linkage
distance

Two introcluster distances :-

(i) compute diameter distance :— The complete diameter distance is defined as the maximum distance of any two objects belonging to the same cluster.

$$d_1(S) = \max \left\{ d(x, y) \middle| x, y \in S \right\}$$



(ii) Average diameter distance :— The average diameter distance is the average distance of all the data points belonging to the same cluster.

It is denoted by,

$$d_2(S) = \frac{1}{|S|(|S|-1)} \cdot \sum_{x, y \in S} d(x, y)$$

(c) Dunn's cluster validation index : →

The Dunn's cluster validation index is defined by the following equation,

$$D\text{Index}(U) = \min_{1 \leq i \leq c} \left\{ \min_{\substack{1 \leq j \leq c \\ j \neq i}} \left\{ \frac{s(x_i, x_j)}{\max_{1 \leq k \leq c} d(x_k)} \right\} \right\}$$

where,

~~s(x)~~ \rightarrow any data points. (set of cluster)

$x_i, x_j \rightarrow$ different cluster.
(i th cluster, j th cluster)

$c \rightarrow$ number of clusters.

$s(x_i, x_j) \rightarrow$ intercluster distance
between x_i and x_j .

$d(x_k) \rightarrow$ intracluster distance of
cluster x_k .

❑ Usefulness of Dunn's Index in cluster evaluation:-

(i) Dunn's validation index can be used to evaluate the cluster quality.

(ii) We know that, the condition for a good cluster system is —

(a) inter cluster distance should be more.

(b) intra cluster distance should be less.

Now, in the equation of Dunn's Validation index, in the numerator we have intercluster distance which is supposed to be more whereas in the denominator we have intraclass cluster distance which is supposed to be less.

Hence, higher value of Dunn's index signify better quality cluster.

(a) Let, D = following labelled database.

~~info(D)~~ Here, class variable = Play Tennis,

- Number of tuples having value "yes" = 6.
- Number of tuples having value "no" = 4.
- Total number of tuples = 10.

$$\begin{aligned}
 \therefore \text{info}(D) &= -\frac{6}{10} \text{wg}_2\left(\frac{6}{10}\right) - \frac{4}{10} \text{wg}_2\left(\frac{4}{10}\right) \\
 &= -\frac{3}{5} \text{wg}_2\left(\frac{3}{5}\right) - \frac{2}{5} \text{wg}_2\left(\frac{2}{5}\right) \\
 &= -\frac{3}{5} [\text{wg}_2^3 - \text{wg}_2^5] - \frac{2}{5} [\text{wg}_2^2 - \text{wg}_2^5] \\
 &= -\frac{3}{5} [1.6 - 2.3] - \frac{2}{5} [1 - 2.3] \\
 &= 0.94
 \end{aligned}$$

→ Calculating Gain of Outlook :-

$$\text{Gain}(\text{outlook}) = \text{Info}_{\text{outlook}}(D) - \frac{\text{info}(D)}{\text{outlook}} \quad \longrightarrow (i)$$

Note $\text{info}_{\text{outlook}}(D)$ outlook

outlook	Count	Yes	No
Sunny	4	2	2
Rain	4	2	2
overcast	2	2	0

$$\therefore \text{info}_{\text{outlook}}(D) = \frac{4}{10} \text{info}_{\text{outlook}}(D_{\text{sunny}}) + \frac{4}{10} \text{info}_{\text{outlook}}(D_{\text{rain}}) + \frac{2}{10} \text{info}_{\text{outlook}}(D_{\text{overcast}}).$$

$$= \frac{4}{10} \left[-\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right] + \frac{4}{10} \left[-\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right] + \frac{2}{10} \left[-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - \frac{0}{2} \log_2 \left(\frac{0}{2} \right) \right].$$

$$\begin{aligned}
 &= \frac{1}{10} \left[\cancel{-2.1} - 2.1 (w_{g_2^1} - w_{g_2^2}) \right] \times 2 + \frac{2}{10} [0] \\
 &= \frac{1}{10} [+] \times 2 \\
 &= \cancel{\frac{1}{10}} \frac{8}{10} \\
 &= \text{req } 0.8
 \end{aligned}$$

\therefore from (i),

$$\begin{aligned}
 \text{Gain(outlook)} &= 0.94 - 0.88 \\
 &= \cancel{0.06} 0.14
 \end{aligned}$$

\rightarrow calculating Gain of Humidity :-

$$\text{Gain(Humidity)} = \text{info}(D) - \text{info}_{\text{Humidity}}(D)$$

Humidity	Count	yes	no
Normal	5	4	1
High	5	2	3

$$\begin{aligned}
 \text{info}_{\text{Humidity}}(D) &= \frac{5}{10} \left[\text{info}_{\text{humidity}}(D_{\text{Normal}}) + \frac{5}{10} \text{info}_{\text{humidity}}(D_{\text{High}}) \right]
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{5}{10} \left[-\frac{1}{5} w_{g_2} \left(\frac{4}{5} \right) - \frac{1}{5} w_{g_2} \left(\frac{1}{5} \right) \right] + \\
 &\quad \frac{5}{10} \left[-\frac{2}{5} w_{g_2} \left(\frac{2}{5} \right) - \frac{3}{5} w_{g_2} \left(\frac{3}{5} \right) \right]
 \end{aligned}$$

$$= -\frac{5}{10} \left[-\frac{1}{5} (wg_2^4 - wg_2^5) - \frac{1}{5} (wg_2^1 - wg_2^5) \right] +$$

$$\frac{5}{10} \left[-\frac{2}{5} (wg_2^2 - wg_2^5) - \frac{3}{5} (wg_2^3 - wg_2^5) \right]$$

$$= -\frac{5}{10} \left[-\frac{4}{5} \times (2 - 2.3) - \frac{1}{5} (0 - 2.3) \right] +$$

$$\frac{5}{10} \left[-\frac{2}{5} (1 - 2.3) - \frac{3}{5} (1.6 - 2.3) \right]$$

$$= 0.35 + 0.47$$

$$= 0.82$$

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

$$= 0.12$$

→ Calculating Gain for Wind :-

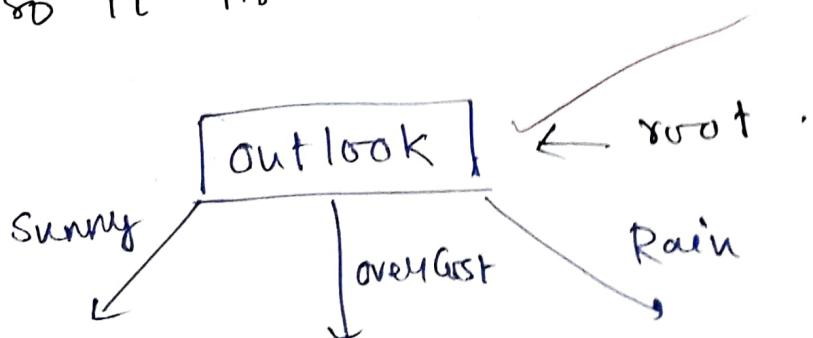
$$\text{Gain (wind)} = \text{info}(D) - \text{info}_{\text{Humidity}}(D).$$

wind	count	yes	no.
strong	5	2	3
weak.	5	4	1

$$\begin{aligned}
 \text{info}_{\text{wind}}(D) &= \frac{5}{10} \left[-\frac{2}{5} \text{wg}_2\left(\frac{2}{5}\right) - \frac{3}{5} \text{wg}_2\left(\frac{3}{5}\right) \right] + \\
 &\quad \frac{5}{10} \left[-\frac{4}{5} \text{wg}_2\left(\frac{4}{5}\right) - \frac{1}{5} \text{wg}_2\left(\frac{1}{5}\right) \right] \\
 &= \frac{5}{10} \left[-\frac{2}{5} (\text{wg}_2^2 - \text{wg}_2^5) - \frac{3}{5} (\text{wg}_2^3 - \text{wg}_2^5) \right] \\
 &\quad + \frac{5}{10} \left[-\frac{4}{5} (\text{wg}_2^4 - \text{wg}_2^5) - \frac{1}{5} (\text{wg}_2^1 - \text{wg}_2^5) \right] \\
 &= \frac{5}{10} \left[-\frac{2}{5} (1 - 2 \cdot 3) - \frac{3}{5} (1 \cdot 6 - 2 \cdot 3) \right] + \\
 &\quad \frac{5}{10} \left[-\frac{4}{5} (2 - 2 \cdot 3) - \frac{1}{5} (0 - 2 \cdot 3) \right] \\
 &= 0.47 + 0.35 \\
 &= 0.82
 \end{aligned}$$

$$\text{Gain}(\text{wind}) = 0.94 - 0.82 = 0.12.$$

Hence, the maximum Gain, we get for the attribute outlook, so it is the root node of the decision tree.



b. Terminating criteria of a decision tree :-

There are two possible terminating criteria of a decision tree —

(i) If the attribute list is empty then we assign the class labels to the corresponding node by using majority voting.

(ii) If all the tuples of a given attribute belongs to the same class labels, then we make it a leaf node by assign the ~~comes~~ pending class label to it.

c. Causes of model overfitting :-

The term 'overfitting' in case of Data mining is defined as below, if a model performs well on the training dataset but perform poorly on a test (unknown) dataset then the model is said to be overfitted. Here are some of the possible reasons of model overfitting —

- (i) Anomalies or noise in the dataset :- If there exist anomalies or noisy data within the training dataset the model will be overfitted.
- (ii) Inconsistent data :- If there exist an inconsistent tuple in the training dataset itself, it causes overfitting. For inconsistent data, the leaf of the decision tree assigns probability to each class labels.
- (iii) Not generalized training sample :- If training sample are not generalized well, it can cause overfitting.

Solution of overfitting in Decision Tree :-

There are two method to solve the overfitting problem

- (i) Pre pruning :- In this method, during the decision tree construction itself, we remove some of the branches of the decision tree.
→ The removal of branch depends upon the initial threshold value.

- Each splitting criteria measures the goodness score which if below threshold, causes the decision tree to be pruned.
- and the pruned attribute is set to the majority class in that subtree.

(ii)

Post Pruning: — In this method, we prune the tree after construction is done.

4

→ The ~~pruned~~ nodes would be pruned branch that would be pruned, will be replaced by the majority class present in the pruned subtree of the decision tree.

5
(a)

Applying Naive Bayes classification algorithm:-

Total number of tuples = 10

~~Total number~~ = 10

Number of classes = 2.

Let, c_1 = when ($\text{playTennis} = \text{"yes"}$).

c_2 = when ($\text{playTennis} = \text{"no"}$).

Number of tuples with, $c_1 = 6$.

Number of tuples with $c_2 = 4$.

$$\therefore P(C_1) = \frac{6}{10} = 0.6$$

$$P(C_2) = \frac{4}{10} = 0.4.$$

Now,

let, $x = \langle \text{outlook} = \text{Overcast}, \text{Humidity} = \text{High}, \text{Wind} = \text{Weak} \rangle.$

→ To predict the class label of x , we need to use the ~~as~~ maximum posteriority,

$$P(c/x) = \frac{P(x/c) \cdot P(c)}{P(x)}$$

~~As~~ As, we need to maximize $P(c/x)$ we ~~will~~ only consider the numerator as denominator is constant over all class labels.

→ For, $c = c_1$ (Play Tennis = "yes") .

Maximum posteriority,

$$\cancel{P(c_1/x)} \quad P(x/c_1) \cdot P(c_1) \cancel{\neq (i)}.$$

Now,
let,

$$P(x/c_1) = P(\text{outlook} = \text{Overcast} / \text{play tennis} = \text{'yes'})$$
$$= \frac{2}{6} = 0.34$$

$$\text{Let } P(x_2/c_1) = P(\text{Humidity} = \text{"High"} / \text{playTennis} = \text{"yes"}).$$

$$= \frac{2}{6} = 0.34$$

$$\text{Let, } P(x_3/c_1) = P(\text{Wind} = \text{"Weak"} / \text{playTennis} = \text{"yes"})$$

$$= \frac{4}{6}$$

$$= 0.67.$$

$$\therefore P(x/c_1) = P(x_1/c_1) \cdot P(x_2/c_2) \cdot P(x_3/c_3)$$
$$= 0.34 \times 0.34 \times 0.67$$
$$= 0.077$$

from (i).

$$\therefore P(x/c_1) \cdot P(c_1) = 0.077 \times 0.6$$
$$= 0.046 \quad \text{--- (2)}$$

That is, $P(x / \text{playTennis} = \text{"yes"}) = 0.046$. --- (2)

NOW, For $c = c_2$ [playTennis = "No"],

let,

$$P(x/c_2) = P(\text{outlook} = \text{"overcast"} / \text{playTennis} = \text{"no"})$$
$$= 0.$$

As, one term is zero,

hence,

$$P(X_1/c_2) = 0.$$

and, ~~P(X/c_2)~~ $\cdot P(c_2)$ is also zero.

So, the Naive Bayes' classification algorithm will predict the class label ~~c_2~~ as playTennis = "yes".

- (b.) In the above, Naive Bayes' classifier algorithm when we tried to calculate the maximum value of $P(X/\text{playTennis} = \text{"no"})$, we got one probability value as zero, hence it causes the whole value to be equal to zero. This is an error. [zero probability error].

→ The error can be corrected using Laplace's estimator which states that if there is a zero probability problem, adding 1 to all the occurrences of the discrete values of that attributes causes the probability to get some value closer to the actual value.

→ Applying this correction to the above Naïve Bayes' classification,

→ For, $c = c_2$ [Play Tennis = "no"]

let, $P(x_1/c_2) = P(\text{outlook} = \text{"overcast"}) / P(\text{playTennis} = \text{"no"})$

$$= \frac{0+1}{4}$$

$$= 0.25$$

(4) let, $P(x_2/c_2) = P(\text{humidity} = \text{"high"}) / P(\text{playTennis} = \text{"no"})$

$$= \frac{3+1}{4}$$

$$= 1$$

let, $P(x_3/c_2) = P(\text{wind} = \text{"weak"}) / P(\text{playTennis} = \text{"no"})$

$$= \frac{1+1}{4}$$

$$> 0.5.$$

$$\therefore P(x/c_2) = 0.25 \times 1 \times 0.5 = 0.125$$

$$\therefore P(x/c_2) \cdot P(c_2) = 0.125 \times 0.4 = 0.05.$$

→ Hence after correction, the actual class label for x is, (Play Tennis = "no").

G.S.
11.5.23

- C. For any feature, if it is continuous valued attribute, then calculation of the likelihood is done by using the Gaussian's distribution as shown below,

$$\text{P}(\text{Humidity} / \text{value } (x_k)) = g(\text{Humidity}, \mu_c, \sigma_c^2)$$

where, μ_c = mean value of the class label.
 σ_c^2 = variance of the class label.

Now,
$$g(\text{Humidity}, \mu_c, \sigma_c^2) = \frac{1}{\sqrt{2\pi\sigma_c^2}} \cdot e^{-\frac{1}{2} \left(\frac{x_k - \mu_c}{\sigma_c} \right)^2}$$

In the above algorithm, we need to replace, the given attribute value probability using the above procedure.

Ex 3a

support :- The support of an association rule $(X \rightarrow Y)$ is defined as the ratio of the number of transaction line Re (XUY) item set occurs with the total number of transactions,

$$\text{Support } (X \rightarrow Y) = \frac{\text{Count}}{n} = \frac{(X \cup Y) \cdot \text{Count}}{n}$$

confidence :- The confidence of an association rule $(X \rightarrow Y)$ is defined as $P(Y/X)$.

$$\text{confidence } (X \rightarrow Y) = P(Y/X) = \frac{P(X \cup Y)}{P(X)}$$

→ An itemset is referred to as the frequent itemset if, its support has at least that of the minimum support.

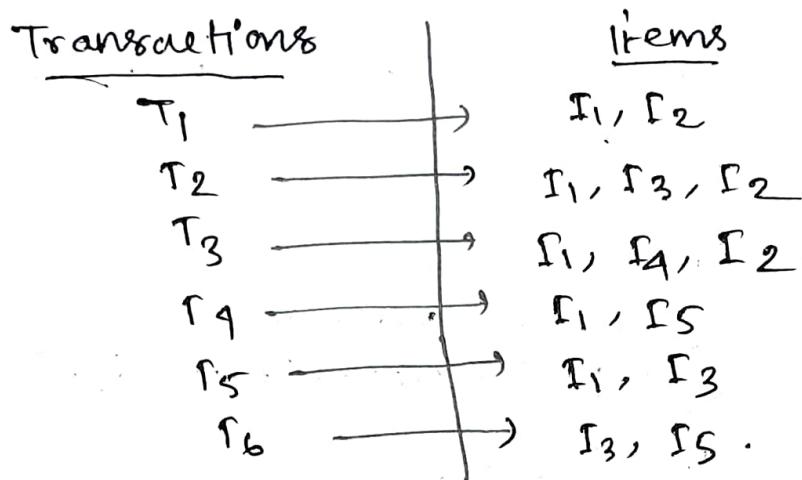
→ An association rule is referred to as the important rule, IF its support $\geq \text{minsup}$ and confidence $\geq \text{minconf}$.

Given, $\text{minsup} = 30/100 = 0.3$
 $\text{minconf} = 80/100 = 0.8$

→ Using Apriori Algorithm,

Let, Bread = I_1 , Jelly = I_4
 Butter = I_2 , Coke = I_5
 Milk = I_3 ,

The transaction database is,



C₁:

items	Support
I ₁	5/6 = 0.84 ✓
I ₂	3/6 = 0.5 ✓
I ₃	3/6 = 0.5 ✓
I ₄	1/6 = 0.17 ✗
I ₅	2/6 = 0.33 ✓

frequent itemset
→

f₁:

items	support
I ₁	0.84
I ₂	0.5
I ₃	0.5
I ₅	0.33

C₂:

items	Support
(I ₁ , I ₂)	3/6 = 0.5 ✓
(I ₁ , I ₃)	2/6 = 0.33 ✓
(I ₁ , I ₅)	1/6 = 0.17 ✗
(I ₂ , I ₃)	1/6 = 0.17 ✗
(I ₂ , I ₅)	0/6 = 0 ✗
(I ₃ , I ₅)	1/6 = 0.17 ✗

frequent itemset
→

f₂:

items	support
(I ₁ , I ₂)	0.5
(I ₁ , I ₃)	0.33

C_3 : NULL

f_3 = NULL

hence, frequent itemset, $f_2 = \{ (I_1, f_2), (I_1, f_3) \}$.

→ Association rules :-

For,

rules	confidence
$I_1 \rightarrow I_2$	$\frac{\text{support}(I_1 \cup I_2)}{\text{support}(I_1)} = 0.5 / 0.84$ = 0.59 ×
$I_2 \rightarrow I_1$	$\frac{\text{support}(I_1 \cup I_2)}{\text{support}(I_2)} = 0.5 / 0.25$ = 0.5 1 ✓

for, f_1, f_3 ,

Rule	Confidence
$I_1 \rightarrow I_3$	$0.34 / 0.84 = 0.40$ ✗
$I_3 \rightarrow I_1$	$0.34 / 0.5 = 0.68$ ✗

So, important association rules, $I_2 \rightarrow I_1$.

Drawbacks of a-priori algorithm:-

- (i) The time complexity is exponential.
- (ii) Space requirement is huge.
- (iii) Not efficient enough.