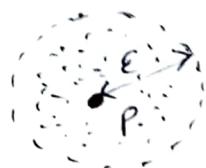


① (a) Terms related to DBSCAN clustering algorithm

(i) Cone object — An object or point is said to be a cone object if there are at least "minpts" No. of objects in the ϵ -neighbourhood of an object or point.

Given, ϵ = parameter or radius of a circle drawn with any object as the centre point.

minpts = Minimum No. of points (an integer value)



p is core point/object if it has more than minpts No. of points in its ϵ -neighbourhood.

(ii) Border object — An object is said to be a border point/object if there are less than minpts No. of points in its ϵ -neighbourhood but the point itself is in the vicinity of a core object or a core point.



p = core obj

q = border object

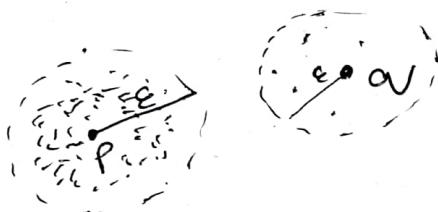
(p, T, o)

(iii) Noise object \rightarrow A point or an object that is neither a border ~~object~~ nor a cone object is a noise point.

In other words, it has less than minpts no. of objects in its ϵ -neighbourhood and it is also not near to any cone point.

(Also, no point is directly density reachable from a cone point)

p = cone object
 a = noise object.



(iv) Density reachability \rightarrow Point p is said to be density reachable from another point q if there are points such as $p_1, p_2, p_3, \dots, p_n$ in between p and q such that those n points are directly density reachable with one other, then point p is said to be density reachable from point q .

Point p and point q may ~~be cone~~ be cone points or border points.

p_1 should be directly density reachable from p .
 p_n should be directly density reachable from q .

(P.T.O)

$P = \{p_1, p_2, p_3, p_4, \dots, p_n\}$

~~reach to each other~~

point p is density
reachable to q .

(v) Density connectivity is two points, p and q are said to be density connected if there is another point r , such that p is ~~density~~ density reachable from r and q is density reachable from r , then points p and q are density connected to each other.



p and q are density connected to each other.

(P.T.O)

(b) The conditions to be satisfied by DBSCAN clustering algorithm are

(i) Maximality condition — clusters formed by DBSCAN clustering algorithm should have as maximum no. of points possible in the clusters formed, so as to have dense clusters. Also, the algorithm should include as many points as possible in the clusters.

(ii) Connectivity condition — points that are produced/put together in a cluster should be density reachable or density connected to each other, such that the algorithm produces ~~separate~~ dense regions as a cluster and sparse regions as noise. This ensures that density based clustering is done properly.

DBSCAN algorithm:

- (i) Select a point p . Density reachable from p .
- (ii) Traverse all points. Points form a cluster.
- (iii) If $p = \text{core object}$,
else if $p = \text{border object}$ or p does not have
minpts no. of points in its ϵ -neighbourhood,
then skip p and goto next point.
- (iv) Continue doing so till all points are exhausted.
 $(f.T.O)$

(c) parameters ↗
 $\text{eps} (\epsilon)$ →

the chosen radius value that dictates how big or small the ϵ -neighbourhood of a point should be.

minpts →

chosen value that dictates the min. no. of points that should be in a ϵ -neighbourhood to form a dense region.

Choosing minpts . ↗

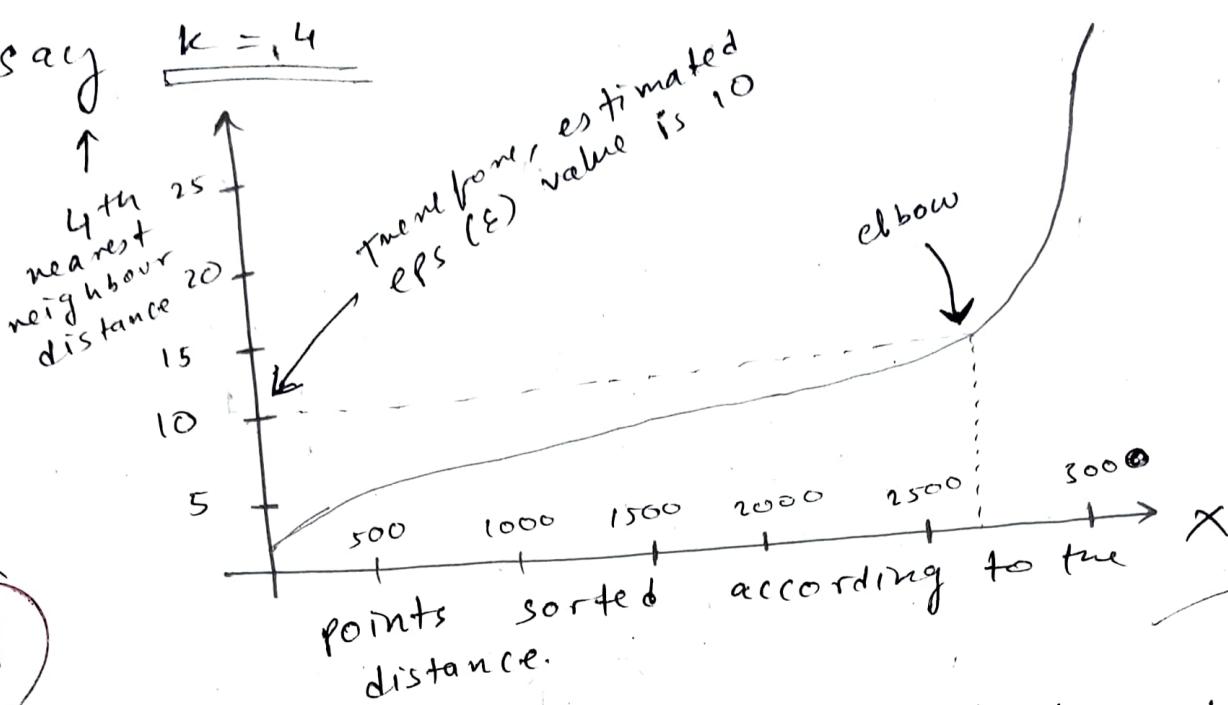
An arbitrary minpts value is chosen by the user and given to the DBSCAN algorithm such that we can estimate an $\epsilon(\text{eps})$ value corresponding to a given minpts value.

Choosing $\text{eps} (\epsilon)$ ↗ from given minpts , we assign $k = \text{minpts. (chosen)}$.

Now, we try to find the k^{th} nearest neighbour distance ~~for~~ all given points in the point set and try to ~~graph~~ plot the sorted graph vs the No. of points according to the distance.

(P.T.O)

say $k = 4$



Now, we plot the k th nearest neighbour distance vs the points sorted according to their distance.

Now, then we use the elbow finding method in the graph to select the point of sharp change in the graph.

corresponding y value to the point of sharp change in the elbow → finding method will give us the appropriate $\text{eps} (\epsilon)$ value, that will result in optimal clustering in the DBSCAN algorithm.

(please turn over)

(P.T.O)

② (a) External Cluster Validation

(i) Purity :- purity is the simplest measure of cluster validation wherein we find the dominant class in a cluster and find its cardinality. Then we divide the cardinality of the dominant cluster to the no. of points in the cluster to give us the purity measure for an ith cluster.

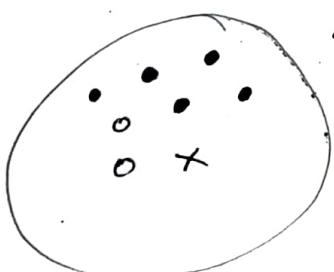
$$\text{purity (i)} = \frac{1}{n} \times \max_j(n_{ij}) ;$$

\uparrow
purity
measuring
of
ith cluster

$n =$ No. of points in
ith cluster

Example :-

- → class A
- → class B
- × → class C



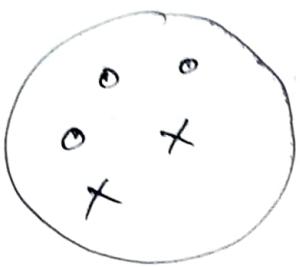
$\max_j(n_{ij}) =$ cardinality
of dominant
class in the ith
cluster.

$$n = 8$$

$$\text{purity} = \frac{1}{8} \times \max(5, 2, 1)$$

$$\text{purity} = \frac{1}{8} \times 5 = \frac{5}{8}$$

dominant class is class A
cardinality of class A is 5 -

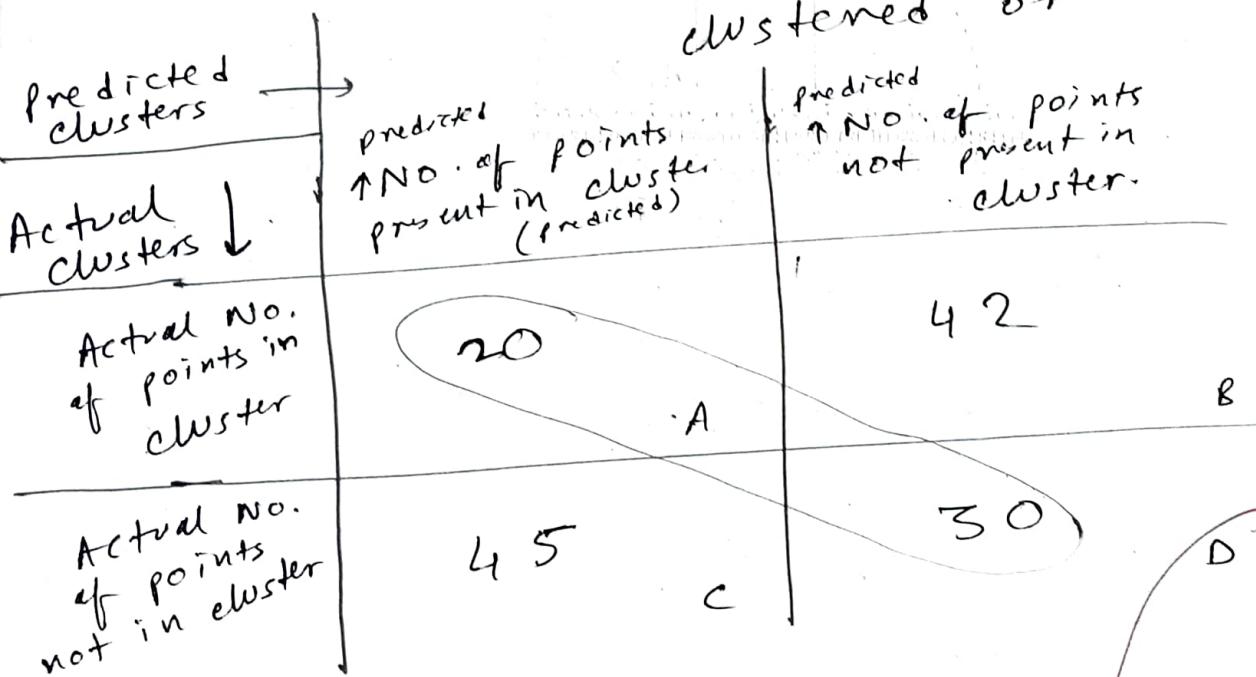


cluster 2

$$\text{parity} = \frac{1}{5} \times \max(0, 3, 2) \\ = \frac{3}{5}$$

dominant class is class B.
so, cardinality of B = 3

(i) Rand index → if it is a measure, similar to the confusion matrix wherein we plot the no. of points in actual clusters to the predicted clusters and try to analyse how much of it is correctly clustered or classified.



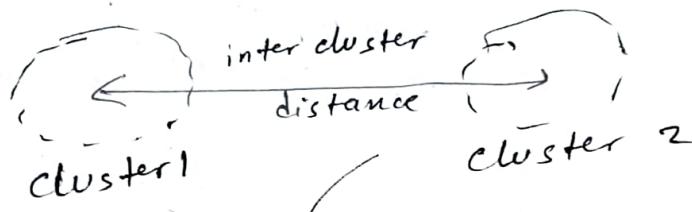
A = True positive
B = False negative
C = False positive
D = True negative

RAND index

$$\text{Rand index (R-I)} = \frac{A + D}{A + B + C + D} = \frac{20 + 30}{20 + 42 + 45 + 30} = 0.365$$

(P-T.O)

(b) inter cluster distance refers to the measure of the distance between two distinct clusters. (S)



Two inter cluster distance are as follows:

(i) Centroid Linkage Distance :
Distance between the centroids of two clusters, given as

$$S_1 = d(\bar{v}_S, \bar{v}_T) \text{ where}$$

$$\bar{v}_S = \frac{1}{|S|} \sum \phi(x)$$

$$\bar{v}_T = \frac{1}{|T|} \sum \phi(y)$$

where

S = Cluster 1

T = Cluster 2

\bar{v}_S = Centroid of cluster 1

\bar{v}_T = Centroid of cluster 2

x = points in cluster 1 ($x \in S$)

y = points in cluster 2 ($y \in T$)

$(P-T^{\circ})$

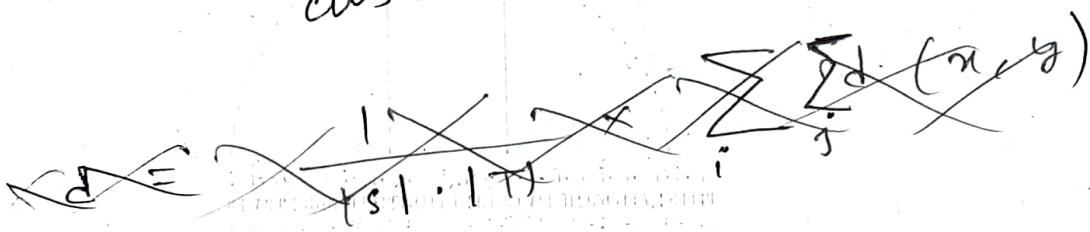
(ii) ~~Simple Linkage dist~~ The smallest intercluster distance in between ~~two~~ two clusters among their nearest points, given by

$$\delta_2 = \min \left(d(x, y) \right)$$

where $x \in S, y \in T.$

(iii) Average Linkage distance is Average distance computed across / in between all the points in cluster 1 and cluster 2.

5



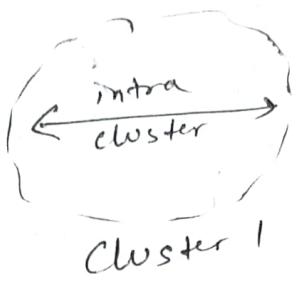
$$\delta_3 = \frac{1}{|S| \times |T|} \times \sum_{x \in S} \sum_{y \in T} d(x, y)$$

where
 $S = \text{cluster 1}$
 $T = \text{cluster 2}$
 $x \in S, y \in T.$

Total no. of points in cluster 1 and cluster 2 computed each-one-one, ~~is~~ is $|S| \times |T|.$

$d(x, y) \rightarrow$ distance between object x and y . (P, T, O)

Intra cluster distance refers to the width of a single cluster.



Two types of intra cluster distance ↗

(i) Complete diameter distance ↗

Maximum distance between any 2 given points in a cluster. which is equal to the ^{largest} diameter of the cluster.

Given by,
 ~~$D_1 = \max(d(x, y))$~~ ,

where $x \in S, y \in S$

(ii) Average diameter distance ↗ from every

point x in cluster S ,

Total $|S| \times (|S|-1)$
No. of points } calculate distance of every other point y and find the average.

Given by
$$D_2 = \frac{1}{|S| \cdot (|S|-1)} \times \sum_{x,y} d(x, y)$$

$d(x, y) \rightarrow$
distance between
 x and y .
(P-T.O)

(iii) Centroid diameter distance is distance of all points in from the centroid \bar{v} , divided by half of the cardinality of cluster S .

$$\text{Distance} = \frac{2}{|S|} \times \sum d(x, \bar{v})$$

(A3)

$$\text{where } \bar{v} = \frac{1}{|S|} \sum x$$

(c) Dunn's cluster validation index is an index that tries to maximize the intercluster distance (δ) and tries to minimize the intracluster distance (Δ), so that a cluster is well formed and distinct.

Dunn's index is basically, the ratio of the intercluster distance (δ) to the maximum intracluster distance (Δ).

Dunn's index: (Mathematically,)

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

where $x_i = i^{\text{th}}$ cluster
 $x_j = j^{\text{th}}$ cluster

$x_k = k^{\text{th}}$ cluster.
 $c = \text{Total No. of clusters.}$

The maximum intra cluster distance among all the clusters is fixed and then it is used as denominator and the numerator is the current inter cluster distance between

32 x_i and x_j :

Since (δ) is larger and (Δ) is smaller,

so, Dunn's index produces value $\gg 1$.

Larger the better is the clustering.

Smaller the value of Dunn's index, worse is the clustering.

This is the usefulness of Dunn's index in clustering problems.

Dunn's index can be used to ~~decide~~ rate / give an index for how good the clustering is and classification or the can be used as a improvement clustering ~~and etc.~~

(P.T.O)

~~(3)~~ (a) Let X and Y be two items.
 Then, the support can be defined as

$$\text{support} = \frac{(X \cup Y) \cdot \text{count}}{n}$$
 where, n is the total no. of transactions,
 and $(X \cup Y) \cdot \text{count}$ is the no. of transactions in which X and Y occur together.

Support can be defined as the ratio of the no. of transactions in which both item X and Y appear together to the total no. of transactions.

Confidence can be defined as,

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

where $X \cdot \text{count} =$ ~~No. of transactions~~ in which only X appears.

$(X \cup Y) \cdot \text{count} =$ No. of transactions in which both X and Y appear.

(P.T.O)

Confidence can be defined as the ratio of the No. of transactions in which both item X and Y appear to the No. of transactions in which only item X appears.

(Confidence can be thought of as conditional probability $P(Y/X)$.)

(4)

An item is referred to as an frequent item set, if its support value is greater than or equal to the given minimum support (minsup) value. ($\text{support} \geq \text{minsup}$)

An association rule is referred to as an important rule, when its confidence value is greater than or equal to the minimum confidence (minconf) value. ($\text{conf} \geq \text{minconf}$)

(b) Given Transactions.

- T_1 :- bread, butter
- T_2 :- bread, milk, butter
- T_3 :- bread, jelly, butter
- T_4 :- bread, coke
- T_5 :- bread, milk
- T_6 :- coke, milk.

$$\begin{cases} \text{minsup} = 30\% \\ \text{minconf} = 80\% \end{cases}$$

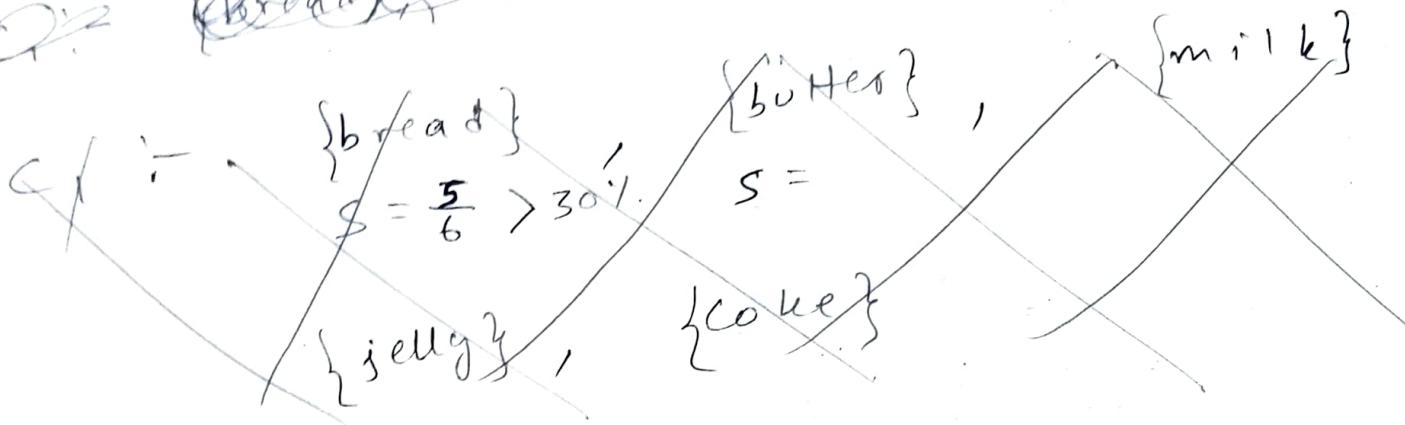
n = Total No. of transactions

$$n = 6$$

Using A priori Algorithm

candidate key / frequent item set generation if
 $s = \text{support value}$.

Q: ~~(bread, butter)~~



$$C_1 := \checkmark \{ \text{bread} \} \rightarrow s = \frac{5}{6} > 30\% \text{ (valid)}$$

$$\checkmark \{ \text{butter} \} \rightarrow s = \frac{3}{6} > 30\% \text{ (valid)}$$

$$\checkmark \{ \text{milk} \} \rightarrow s = \frac{3}{6} > 30\% \text{ (valid)}$$

$$\times \{ \text{jelly} \} \rightarrow s = \frac{1}{6} < 30\% \text{ (invalid)}$$

$$\checkmark \{ \text{coke} \} \rightarrow s = \frac{2}{6} > 30\% \text{ (valid)}$$

$F_1 := \{ \text{bread} \}, \{ \text{butter} \}, \{ \text{milk} \}, \{ \text{coke} \}$
 (based on support values)

Now,

(P.T.O)

C₂ :- $\checkmark \{ \text{bread, butter} \} \rightarrow s = \frac{3}{6} > 30\%.$ (valid)

$\checkmark \{ \text{bread, milk} \} \rightarrow s = \frac{2}{6} > 30\%.$ (valid)

$\times \{ \text{bread, coke} \} \rightarrow s = \frac{1}{6} < 30\%.$ (invalid)

$\times \{ \text{butter, milk} \} \rightarrow s = \frac{1}{6} < 30\%.$ (invalid)

$\times \{ \text{butter, coke} \} \rightarrow s = \frac{0}{6} < 30\%.$ (invalid)

$\times \{ \text{milk, coke} \} \rightarrow s = \frac{1}{6} < 30\%.$ (invalid)

F₂ :- $\{ \text{bread, butter} \}, \{ \text{bread, milk} \}$

C₃ :- $\{ \text{bread, butter, milk} \}$

$$S = \frac{1}{6} < 30\%.$$
 (invalid)

only 1 candidate from
possible candidates
candidate generation

Also, we can see that the present
subset $\{ \text{butter, milk} \}$ is not ~~subset~~
in F₂. So this combination
is not valid in C₃.

So,

F₃ :- \emptyset (null set)

80, the frequent itemset can be written as

$$F = \left\{ \left\{ \text{bread, butter} \right\}, \left\{ \text{bread, milk} \right\} \right\}$$

$$\text{minconf} = 80\%.$$

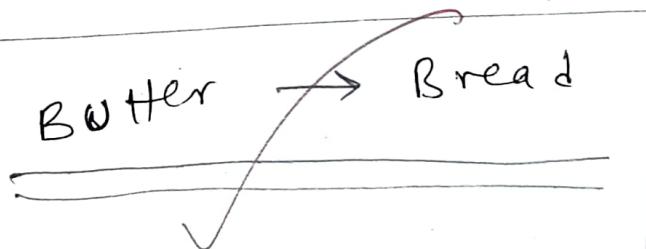
Now,
Association rules :-

$X \rightarrow Y$		confidence
\times	bread	butter
\checkmark	butter	bread
\times	bread	milk
\times	milk	bread

So, as per $\text{minconf} = 80\%$,
only one association rule can be deemed as important rule.

So,

Important
Association
rule



with confidence $= 100\% > \text{minconf} = 80\%$.

- (c) The major drawbacks of apriori algorithm are as follows:
- (i) Lack of data or incomplete data set can result in improper association rules that may not reflect the real scenario.
 - (ii) Improper tuning of min support and min. confidence values can lead to improper association rules.
 - (iii) Too much variations in transactions may lead to improper extraction of frequent itemsets that will result in wrong association rules.
 - (iv) Apriori algorithm can sometimes give wrong association rules when there are conflicting transactions.

(P.T.O)

From the given dataset in the question-paper,
 we can calculate the entropy of the class variable as

$$I(p, n) = - \frac{p}{p+n} \log_2 \left(\frac{p}{p+n} \right) - \frac{n}{p+n} \log_2 \left(\frac{n}{p+n} \right).$$

(5x)
52

Now, $p =$ Yes class value. $\quad p = 6 \quad \} \quad p+n = 10$
 $n =$ No class value. $\quad n = 4 \quad \}$

$$\therefore I(6, 4) = - \frac{6}{10} \log_2 \left(\frac{6}{10} \right) - \frac{4}{10} \log_2 \left(\frac{4}{10} \right)$$

~~$$= - \left[\frac{6}{10} \times (\log_2^2 + \log_2 3) \right]$$~~

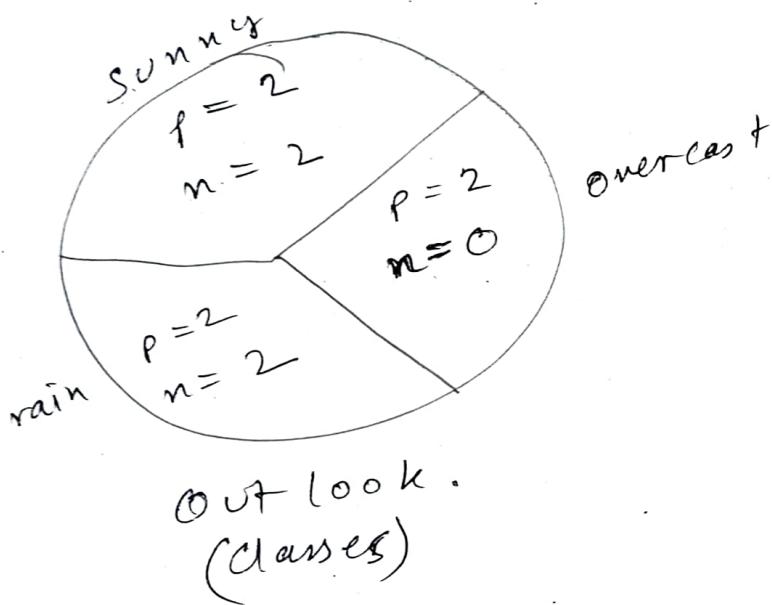
$$= - \left[\frac{6}{10} \times (\log_2^2 + \log_2 3 - \log_2 10) + \frac{4}{10} \times (\log_2^2 - \log_2 10) \right]$$

Now, $\log_2 10 = \log_2^2 + \log_2^5$

$$= 1 + 2 \cdot 3 = 3.3$$

$$\begin{aligned}
 \Rightarrow I(6/4) &= - \left[\frac{6}{10} (1 + 1 \cdot 6 - 3 \cdot 3) \right. \\
 &\quad \left. + \frac{4}{10} \times [2 - 3 \cdot 3] \right) \\
 &= - \left[0.6 \times -0.7 + 0.4 \times -1.3 \right] \\
 &= 0.42 + 0.52 \\
 &= 0.94
 \end{aligned}$$

NOW,
For attribute outlook is



~~PCZ~~

p for sunny \Rightarrow No. of yes outcome for class sunny of attribute outlook

~~on~~ for sunny
 \Rightarrow No. of no outcome for class sunny of attribute outlook

Thus,

$E(\text{outlook}) =$ Expected entropy
for attribute outlook

Similarly for all
classes, we find out
 p and n .

$$\text{Now, } E(\text{outlook}) = \sum_i \frac{p_i + n_i}{p+n} I(p_i, n_i)$$

$$= \frac{4}{10} \times I(2,2) + \frac{2}{10} \times I(2,0) \\ + \frac{4}{10} \times I(2,2)$$

$$\text{Now, } I(2,0) = -\frac{2}{2} \log_2 1 = 0$$

$$I(2,0) = 0$$

$$I(2,2) = -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4}$$

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

$$= -[\log_2 2 - \log_2 4]$$

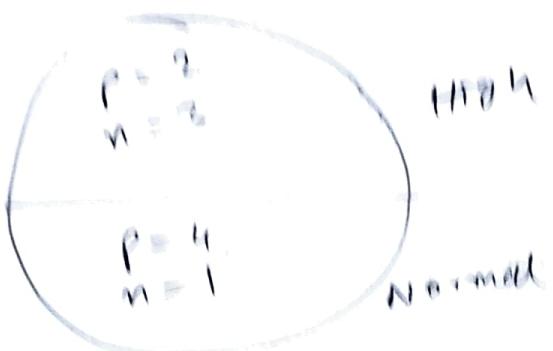
$$= -[1 - 2]$$

$$\Rightarrow I(2,2) = 1$$

$$\Rightarrow E(\text{outlook}) = \frac{4}{10} \times 1 + 0 + \frac{4}{10} \times 1 \\ = \frac{8}{10} = 0.8$$

$$\Rightarrow \boxed{E(\text{outlook}) = 0.8}$$

For Attitude Humidity



humidity
(2 classes)

No. 10

$$E(\text{Humidity}) = \frac{5}{10} \times I(4,1) + \frac{5}{10} \times I(2,3)$$

$$\begin{aligned} \text{Now, } I(4,1) &= -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \\ &= \left[\frac{4}{5} \times \left[2 \log_2 2 - \log_2 5 \right] + \frac{1}{5} \times \left[\log_2 1 - \log_2 5 \right] \right] \\ &= \left[\frac{4}{5} \times [2 - 2.3] + \frac{1}{5} \times -2.3 \right] \\ &= + \frac{4}{5} \times 0.3 + \frac{1}{5} \times 2.3 \end{aligned}$$

0.7

$$\begin{aligned} I(2,3) &= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \\ &= \left[\frac{2}{5} \times (\log_2 2 - \log_2 5) + \frac{3}{5} (\log_3 - \log_5) \right] \\ &= \cancel{\frac{2}{5}} \times 1.3 + \frac{3}{5} \times 0.7 \\ &= 0.94 \end{aligned}$$

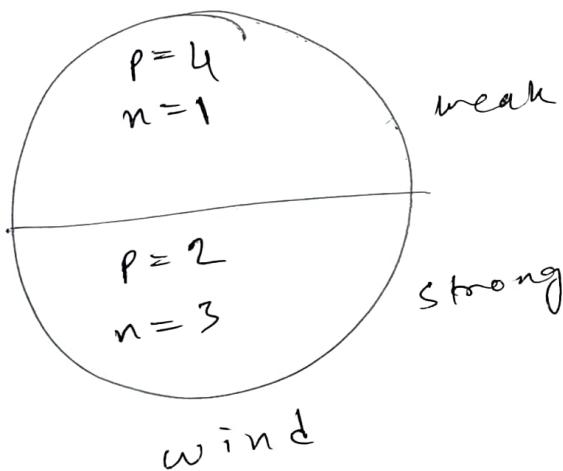


11/8/20

Signature(s) of the Invigilator with date

$$E(\text{humidity}) = 0.5 \times 0.7 + 0.5 \times 0.94 \\ = 0.82$$

for attribute wind,



$$E(\text{wind}) = \frac{5}{10} \times I(2,3) \\ + \frac{5}{10} \times I(4,1)$$

using previously calculated values,

$$E(\text{wind}) = 0.5 \times 0.94 + \\ 0.5 \times 0.7 \\ = 0.82$$

Now,
we have

$$E(\text{outlook}) = 0.8$$

$$E(\text{humidity}) = 0.82$$

$$E(\text{wind}) = 0.82$$

(P-T-O)

Now, information gain for each attribute

$$\text{gain}(\text{outlook}) = I(p, n) - E(\text{outlook})$$

$$= 0.94 - 0.8$$

$\boxed{\text{gain}(\text{outlook}) = 0.14}$

$$\text{gain}(\text{humidity}) = 0.94 - 0.82$$

$$= 0.12.$$

$$\text{gain}(\text{wind}) = I(p, n) - E(\text{wind})$$

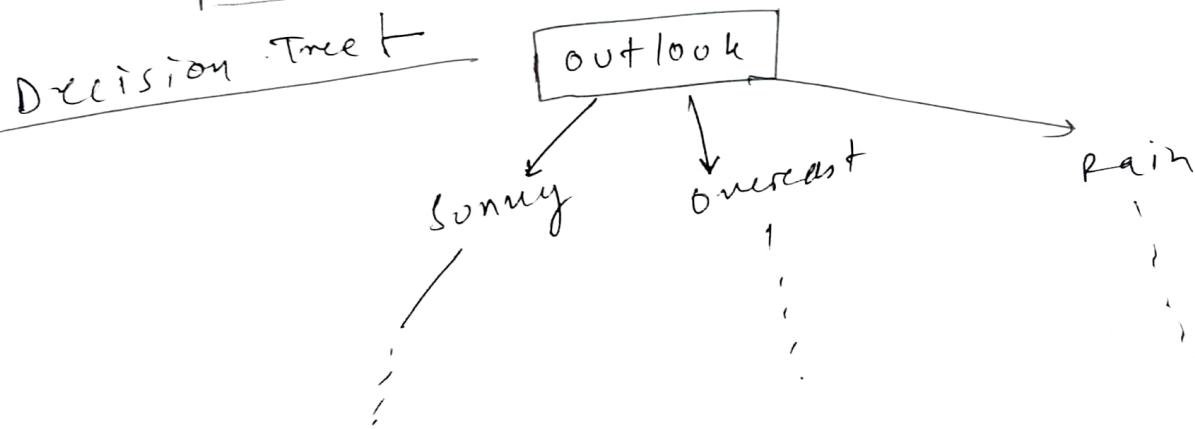
$$= 0.94 - 0.82$$

$$= 0.12.$$

Now, we have maximum information gain for attribute outlook

so,

outlook is the root for the decision tree.



- (b) The possible terminating criteria for a decision tree algorithm are:
- (i) All attributes exhausted and no more attributes are left to be added to tree.
 - (ii) The tree is complete without taking into consideration, all possible attributes, i.e., some attributes influence the final outcome so less that they are not needed in the decision tree.
 - (iii) from given attributes, no complete decision tree is possible, i.e. very few attributes but too many variations in tuples to be factored in.
- (c) The causes of model overfitting are noisy data present in the table or too many variations that do not influence a proper outcome.
- Overfitting can be defined as high training accuracy but low validation accuracy.
- The overfitting problem in decision trees can be solved using
- (i) Pre-pruning
 - (ii) Post-pruning
- (P.T.O)

Pre pruning if we prune the branches in a decision tree, while we are creating it. i.e., we perform pruning and creation simultaneously till we get a final tree.

③

Post pruning if we create a complete decision tree first, then we prune branches, depending on conflicts present in tree and table (miss classifications)

⑤ (a) Bayesian classification

$$P(C_i/X) = P(C_i) \times \underbrace{P(X/C_i)}_{\uparrow}$$

Given tuple,

$$X = (\text{outlook} = \overset{\text{overcast}}{\text{overcast}}, \text{humidity} = \overset{\text{high}}{\text{high}}, \text{wind} = \overset{\text{weak}}{\text{weak}})$$

$$P(C_i = \text{yes}) = \frac{6}{10}$$

$$P(C_i = \text{no}) = \frac{4}{10}$$



Examination Roll No. 510819011 Sheet No. 2

Subject Code CC-4161

skm 11/5/23

Signature(s) of the Invigilator with date

C₁ = No

$$P(\text{outlook} = \text{overcast}/\text{No}) = \frac{0}{4}$$

$$P(\text{humidity} = \text{high}/\text{No}) = \frac{3}{4}$$

$$P(\text{wind} = \text{weak}/\text{No}) = \frac{1}{4}$$

Here we have same zero probability problem.

so we use Laplacian correction - (+1 to all cases)

$$P(\text{overcast}/\text{No}) = \frac{1}{7}$$

$$P(\text{high}/\text{No}) = \frac{4}{7}$$

$$P(\text{weak}/\text{No}) = \frac{2}{7}$$

Total 3 cases = so denominator +3 in all

numerator +1
in all.

Using Laplacian
correction to avoid
zero probability
problem.

$$C_1 = 4\text{e}1$$

$$P(\text{outlook} = \text{overcast}/4\text{e}1) = \frac{2}{6}$$

$$P(\text{humidity} = \text{high}/4\text{e}1) = \frac{2}{6}$$

$$P(\text{wind} = \text{weak}/4\text{e}1) = \frac{4}{6}$$

(4.2)

Now,

$$P(\times/C_1=4\text{e}1) = \frac{2}{6} \times \frac{2}{6} \times \frac{4}{6}$$

$$= \frac{16}{216} = \frac{2}{27}$$

$$= 0.074$$

$$P(\times/C_1=\text{No}) = \frac{1}{7} \times \frac{4}{7} \times \frac{2}{7}$$

$$= 0.023$$

Now, for given tuple X ,

$$\begin{aligned} P(c_i = \text{yes} | X) &= P(c_i = \text{yes}) \times P(X | c_i = \text{yes}) \\ &= \frac{6}{10} \times 0.074 \\ &= 0.044 \end{aligned}$$

$$\begin{aligned} P(c_i = \text{no} | X) &= P(c_i = \text{no}) \times P(X | c_i = \text{no}) \\ &= \frac{4}{10} \times 0.023 \\ &\approx 0.0092. \end{aligned}$$

So, for $c_i = \text{yes}$, we have the higher probability,

So, predicted class Label = Yes

for given tuple $X = (outlook = \text{overcast}, \text{humidity} = \text{high}, \text{wind} = \text{weak})$

(b) The main error in bayesian classification problem is the zero-probability error, where we may get a value 0 corresponding to a probability. But upon multiplication, we can get the total class value probability as 0¹⁰, whereas it may not be true in the real scenario.

This is the zero probability problem. Such a problem can be corrected using Laplacian correction, wherein we add +1 to all possible ~~existing~~ probabilities that may be affected by a zero value.

Suppose,

$\begin{cases} P(A) = \frac{0}{7} \\ P(B) = \frac{3}{7} \\ P(C) = \frac{4}{7} \end{cases}$

To avoid zero probability, we use ~~Laplacian~~ Laplacian correction and add 1 to all numerator probabilities. So denominator becomes (+3) for all cases as we add (+1) to 3 numerators in 3 probabilities.

After Laplacian correction,
 $P(A) = \frac{1}{10} ; P(B) = \frac{4}{10} / P(C) = \frac{5}{10}$

(c) If humidity attribute had continuous values, then we would have had to convert the continuous values to ~~set~~ a set of discrete class values using various discretization techniques in data pre processing. Then we could have applied decision Tree / Bayesian Algorithm as Decision tree works only on discrete values.

