3. (a) Suppose there is a rule such that,

$$X \rightarrow Y$$

Then the support of this rule is defined as the

$$Support = \frac{(X \cup Y). Count}{N}$$

where N is total number of samples i.e. it is the ratio of samples containing $X \cup Y$ divided by total No. of instances.

and the item set is referred to as frequent if the it's min support is greater than or equal to the minimum support and.

An item is referred to as important rule if it it's frequent and the confidence of item set is greater than or equal to the minimum confidence.

④ Confidence = $\frac{(X \cup Y). Count}{X. Count}$

6.) Let $I_1 =$ Bread

$I_2 =$ Butter

$I_3 =$ Milk

$I_4 =$ Jelly

$I_5 =$ Coke

where:

$T_1 :$ $I_1, I_2$

$T_2 :$ $I_1 I_3 I_2$

$T_3 :$ $I_1 I_4 I_2$

$T_4 :$ $I_1 I_5$

$T_5 :$ $I_1 I_3$

$T_6 :$ $I_3 I_5$

Min.- Confidence $= 0.8$

Min - Support $= 0.3$

So, $C_1$ are

| | Confidence support |
|---|---|
| $I_1 :$ | 5/6 |
| $I_2 :$ | 3/6 |
| $I_3 :$ | 3/6 |
| $I_4 :$ | 1/6 |
| $I_5 :$ | 2/6 |

So, $F_1$ are

$I_1$ $I_2$ $I_3$ $I_5$

So, $C_2$ are

$I_1 \quad I_2 \quad - \quad 3/6$

$I_1 \quad I_3 \quad - \quad 2/6$

$I_1 \quad I_5 \quad - \quad 1/6$

$I_2 \quad I_3 \quad - \quad 1/6$

$I_2 \quad I_5 \quad - \quad 0/6$

$I_3 \quad I_5 \quad - \quad 1/6$

So $F_2$ are

$I_1 \quad I_2$

$I_1 \quad I_3$

$\frac{1}{2} = \frac{3}{6}$

So, $C_3$ are

$I_1 \quad I_2 \quad I_3 \quad - \quad 1/6$

So, $F_3$ is None:

Rules $\quad x \rightarrow y \quad$ Possible are

| X | Y | confidence |
|---|---|---|
| $I_1$ | $I_2$ | $3/5$ |
| $I_2$ | $I_1$ | $3/3$ |
| $I_1$ | $I_3$ | $2/5$ |
| $I_3$ | $I_1$ | $2/3$ |

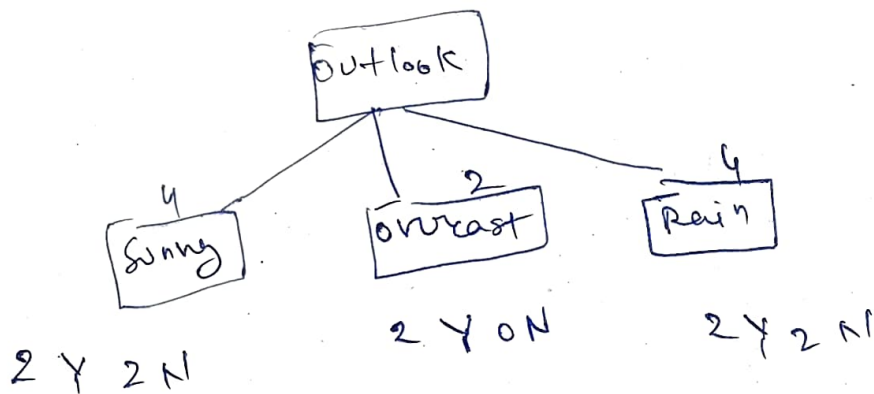So, only Rule having confidence greater than $0.8$ is

$$\boxed{I_2 \rightarrow I_1}$$

c) The major drawback of a-priori algorithm is also it's computational complexity. Sometimes the possible No. of cases may be very large and hence it is slower.

There is always some ambiguity in providing the result and hence it may lead to some errors as well. The result machine is not very good and hence may need improvement.

4.

a) If we split upon outlook



So, entropy upon splitting on outlook is

For 2Y 2N

Entropy is

$$-\frac{2}{4} \log \frac{2}{4} - \frac{2}{4} \log \frac{2}{4}$$

$$-\log \frac{2}{4} \Rightarrow -\log \frac{1}{2} = \log 2 = 1$$

For 2Y 0N  entropy is

$$-\frac{2}{2} \log \frac{2}{2} - \frac{0}{2} \log \frac{0}{2}$$

Again For Rain entropy is = 0

Again for Rain entropy is 1

So, weighted entropy is

$$\frac{4}{10} \times 1 + \frac{2}{10} \times 0 + \frac{4}{10} \times 1$$

$$\Rightarrow \frac{8}{10} = 0.8$$

Now we split upon Humidity

③½

Humidity



Normal          High   = Entropy = 0.82

4 Yes 1 No      2 Yos 3 NO

$-\frac{4}{5} \log \frac{4}{5} -\frac{1}{5} \log \frac{1}{5}$     $-\frac{2}{5} \log \frac{2}{5} -\frac{1}{3} \log \frac{1}{3}$

$\Rightarrow -\frac{4}{5}(\log 4 - \log 5) + \log 5 \times \frac{1}{5}$     $\Rightarrow \frac{2}{5}(\log 3 - \log 2) + \frac{1}{3} \log 3$

$\Rightarrow \frac{4}{5} \log(5 - \log 4) + \log 5 \times \frac{1}{5}$     $\Rightarrow \frac{2}{3}(0.8) + \frac{1}{3} \times 1.6$

$\frac{4}{5}(2.3 - 2) + 2.3 \times \frac{1}{5}$     $\Rightarrow \frac{1.2}{3} + \frac{6}{3} = \frac{2.8}{3}$

$\frac{4}{5} \times 0.3 + \frac{2.3}{5}$     $= 0.93$

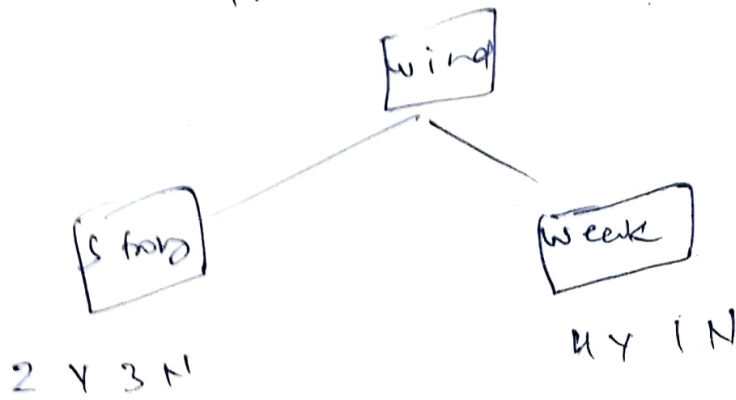$\Rightarrow \frac{1.2}{5} + \frac{2.3}{5} = \frac{3.5}{5} = 0.7$     $= -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5}$

$\frac{2}{5}(1.3) + \frac{3}{5}(0.7)$

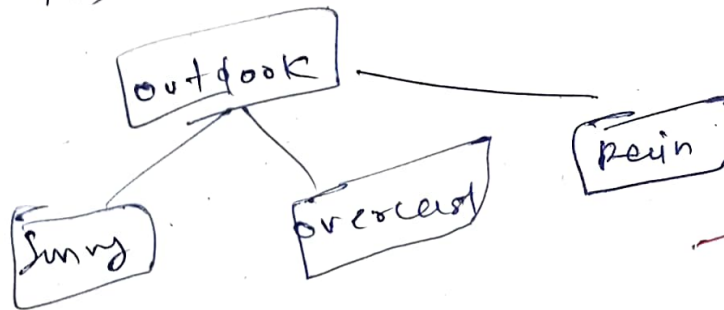So, weighted     $\Rightarrow \frac{2.6}{5} + \frac{2.1}{5} \Rightarrow \frac{4.7}{5}$

So Entropy is $\frac{1}{2} \times 0.7 + \frac{1}{2} \times 0.94$     $= 0.94$

$\Rightarrow 0.35 + 0.47$

$\Rightarrow 0.82$

If we split on wind

[wind]

[strong]                    [weak]

2 Y 3 N                    4Y 1N

It will be seme as Humidity
because both are splitting caually
So, it's entropy is 0.82

So, it's root would be outlook

[outlook]                    [rain]

[sunny]    [overcast]

b) The possible terminality criterion of a decision tree algorithm is

i) If all the values in the Node have same labels then we should terminate that

ii) If the entropy gain upon splitting the node is less than a threshold value then also no splitting

$$\left(\frac{1}{2}\right)$$

iii) If upon splitting the depth of tree will become greater than a threshold depth then also we will not split.

C) The causes of model overfitting can be beacause the data given had too many outliers and it was not properly segmented.

Overfitting corresponds to when a decison tree work very well on training data but worst on the test data.

It can be solved by mostly 2 ways.

i) pre-pruning
ii) post-pruning

③

i) pre-pruning : The pre-pruning a decision tree means is to set some sort of rules so that you further do not split the node and terminate that. For Example by defining a minimum information guin or the max depth of the tree you can terminate a node

ii) post-pruning : you First develop a decison tree and prune the splits or nodes whose split does not seem useful. In this way you save the decision tree from overfitting and hence it provide good resut.

6.

| Outlook | Humidity | Wind | Plg Tennis (Class variable) | Plg Tennis (predicted) |
|---|---|---|---|---|
| Sunny | Normal | Strong | Yes | Yes |
| Overcast | Normals | Strong | No | Yes |
| Rain | Htim | Strong | Yes | No |
| Sunny | Htmn | Weak | No | No |
| Rain | High | Strong | No | No |

④

The confussion Matrix Correspond to



TP | FN
FP | TN

| 1 | 1 |
| 1 | 2 |

i) precision $= \dfrac{TP}{TP + FP} = \dfrac{1}{1+2} = \dfrac{1}{3} = 0.33 \approx 0.5$

ii) Recall $\dfrac{TP}{TP+FN} = \dfrac{1}{1+1} = \dfrac{1}{2} = 0.5$

iii) F-score is Harmonic mean of precision and Recall
$= 2 \times \dfrac{Precision \times Recall}{Precision + Recall}$

$= \dfrac{2 \times \frac{1}{2} \times \frac{1}{2}}{\frac{1}{2} + \frac{1}{2}} \Rightarrow 0.5$

$= \dfrac{1}{2}$

$\dfrac{2}{6} = \dfrac{5}{6} = 0.4$

5) i) <u>Holdout method</u> :-

The Holdout method used for estimatib classification model is that we hold the classifier ho model model so that it gives some correct output

ii) <u>Cross-Validation</u> :- In the cross validation method the dataset is deo devided into k sub sets so that (k-1) datasets are used for trainib and I 's used for testing

iii) <u>Bootstrap</u> : In the bootstrap techniqlle all the 1/2 doo mo classifiers ar bootstrapped so that it gives correct output upon predicting the data

c) The purpose of ensemblers of Classifiers is that is the majority voting. It leads to avoid the overfitting of the classification model. So like instead of predictib on 1 strong classifier, we have more than one Classifier so that we have many weak models and we predict the result on the basis of the majority voting.

The Ada boost algorithm works like this when we are training the data set on the classification algorithm, it ~~on training data~~ some instance of training data if it gives correct output then, we multiply the weight of that Instance by $\frac{e}{1-e}$ given

e is the error rate of classifier so that it's weight decreases and next time it focus more on instances where it predicts wrong.

on testing if some classifier predicts some class then the weight of that class is increased by

$-\log\left(\frac{e}{1-e}\right)$ ..

In the end the model will predict the class with the highest weight.

2.2

2. b) Intercluster distance: The Inter cluster distance is defined as the distance between elements or data points of the different clusters. For example Minimum Inter cluster distance is defined as the minimum distance between two points both belonging to different clusters.

Min - inter cluster distance = $\text{Min} \sum_{q=1}^{N} \sum_{j=1}^{N} \delta(x_i x_j)$

③ where $x_i$ belongs to cluster 1 and $x_j$ belongs to cluster 2.

where $\delta(x_i x_j)$ is the distance between point $x_i$ and $x_j$

Also, the maximum inter cluster distance is defined as the maximum distance between two data points both belonging to different cluster. So,

Max - intercluster distance = $\text{Max} \sum_{i=}^{N} \sum_{j=1}^{m} \delta(x_i x_j)$

$\Delta (x_k) \text{ MAX}$

Intra Cluster distance corresponds to the distance between the points belonging to the same cluster.

Maximum intra-cluster distance is the distance between points belonging to the same cluster.

$\text{Min} \quad i \times j = \dfrac{N \times (N-1)}{2} \quad \delta(x_i, x_j)$

The average intra cluster distance is the average distance between the points belonging to the same cluster.

$$(Avg) \quad \frac{\sum_{i,j=\frac{N\times(N-1)}{2}} \delta(x_i, x_j)}{\frac{N\times(N-1)}{2}}$$

c) The Dunn's cluster validation index is defined as

$$Min \quad \sum_{i=N}^{N} \sum_{j=1}^{M} \frac{\delta(x_i x_j)}{\left(\max_{k=1}^{N} \Delta(x_k)\right)}$$

②

with the help of dunn's cluster validation index we try to maximise the inter cluster distance and minimize the intra-cluster distance. The maximum is the dunn's index the good the closer cluster result is. In this way it helps us in evaluating the cluster.

2. a)

(i) _Purity_ :- Purity corresponds to validate the cluster index in a different way so that the it measures the Purity of clusters. The more distant the different clusters are the more greater it is

O

(ii) _Rand index_ :- Rand index Corresponds to the random cluster formation by the clustering algorithms so that when the clusters are evaluated we get the best out of them.