

3) a) Let  $I$  be an itemset,  $I = \{i_1, i_2, \dots, i_m\}$ .

Let  $t$  be a transact<sup>n</sup>  $t \subseteq I$  and  $T$  be the dataset of transactions.

Assuming,  $X$  and  $Y$  to be subsets of transaction dataset such that  $X, Y \subseteq I$  and  $X \cap Y = \{\phi\}$

An association rule is given by:  $X \rightarrow Y$

which denotes if transact<sup>n</sup>  $X$  happens then transact<sup>n</sup>  $Y$  will happen with some probability.

Support: The rule has a support which can be given by prob. of  $X \cup Y$ .

i.e. 
$$\text{Support} = \frac{\{X \cup Y\}_{\text{count}}}{n}; \quad n = \text{no. of transact}^n$$

Confidence: For the rule above, confidence is the probability of occurrence of  $Y$  if transaction  $X$  has already been done.

$$\text{Confidence} = \frac{\{X \cup Y\}_{\text{count}}}{\{X\}_{\text{count}}}$$

→ An item set is referred to as frequent item set if its support is greater than or equal to the minimum support.

→ An association rule is referred to as an important rule if it has both support and confidence more than their minimum values.

b) Given

Min-support = 80%

Min-confidence = 80%

Size-1 itemset	Frequency
Bread	5 ✓
Butter	3 ✓
Milk	3 ✓
Jelly	1 ✓
Coke	2 ✓

Min. support  $\geq 0.3$

$$\Rightarrow \frac{\text{X count}}{6} \geq 0.3$$

$\therefore F_1 = \{\text{Bread}\}, \{\text{Butter}\}, \{\text{Milk}\}, \{\text{Coke}\}$   $\Rightarrow (X_{\text{count}} \geq 1.8 \rightarrow 2)$

$F_1$ : Frequent itemsets of size '1'

Size-2 itemset	Frequency
Bread, Butter	3 ✓
Bread, Milk	2 ✓
Butter, Milk	1
Bread, Jelly	1
Butter, Jelly	1
Bread, Coke	1
Milk, Coke	1

$\therefore F_2 = \{\text{Bread, Butter}\}, \{\text{Bread, Milk}\}$

Size-3 itemset	Frequency
Bread, Milk, Butter	1
Bread, Jelly, Butter	1

$F_3 = \{\}$

None of the itemset satisfies min. support criteria

Subsets of  $\{\text{Bread, Butter}\}$  are  $\{\text{Bread}\}, \{\text{Butter}\}$

Rule 1: Bread  $\rightarrow$  Butter

Support =  $\frac{3}{6} = 0.5$ , Confidence =  $\frac{3}{5} = 0.6 < 0.8$   
(Not important)

Rule 2 : Butter  $\rightarrow$  Bread

$$\text{Support} = \frac{3}{6} = 0.5, \quad \text{Confidence} = \frac{3}{3} = 1 > 0.8$$

(important) ✓

From the set {Bread, milk}

Rule 3 : Bread  $\rightarrow$  Milk

$$\text{Support} = \frac{2}{6} = \frac{1}{3} = 0.33, \quad \text{Confidence} = \frac{2}{5} = 0.4 < 0.8$$

(Not important)

Rule 4 : milk  $\rightarrow$  Bread

$$\text{Support} = \frac{2}{6} = 0.33, \quad \text{Confidence} = \frac{2}{3} = 0.66 < 0.8$$

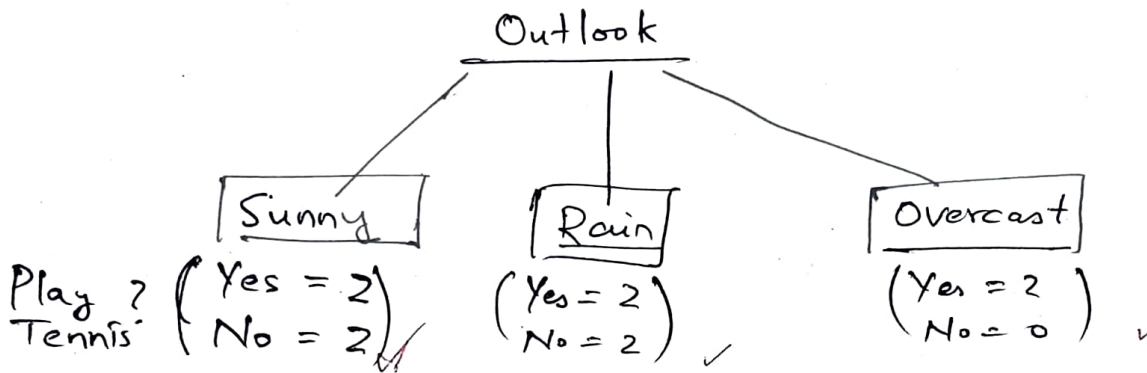
(Not important)

∴ The important rule obtained by Apriori algorithm is :

Butter  $\rightarrow$  Bread (Confidence = 100%  
Support = 50%) ✓

c) ~~Following are the major drawbacks of a-priori algorithm~~  
~~It takes~~

4) a) If decision is taken on Outlook.

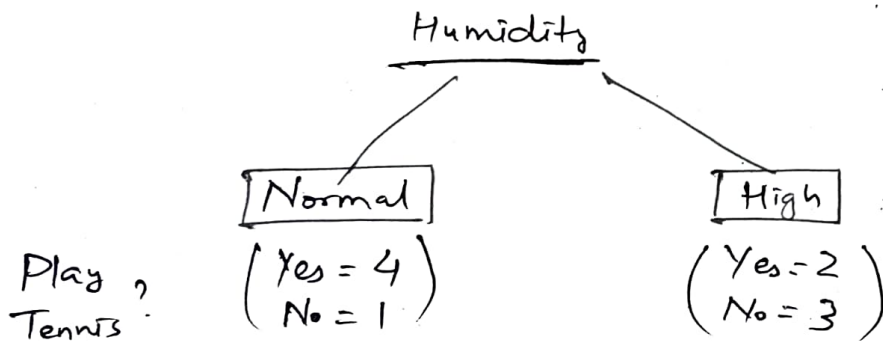


$$\text{Expected info / Entropy} = \sum_{i=1}^k \frac{p_i + n_i}{p+n} I(p_i, n_i)$$

$$E(\text{outlook}) = \frac{4}{10} I(2,2) + \frac{4}{10} I(2,2) + \frac{2}{10} I(2,0)$$

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

$$\therefore E(\text{outlook}) = 0.4 \times 1 + 0.4 \times 1 + 0.2 \times 0 = \underline{0.8}$$



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

$$\begin{aligned}
 I(4,1) &= - \left( \frac{4}{5} \log_2 \left( \frac{4}{5} \right) + \frac{1}{5} \log_2 \left( \frac{1}{5} \right) \right) \\
 &= 0.8 \left[ \log_2(5) - \log_2(4) \right] + 0.2 \log_2(5) \\
 &= \log_2(5) - 0.8 \log_2(2) \\
 &= 2.3 - (1.6 \times 1) = \underline{0.7}
 \end{aligned}$$

min = Humidity

$$I(2,3) = - \left[ \frac{2}{5} \log_2 \left( \frac{2}{5} \right) + \frac{3}{5} \log_2 \left( \frac{3}{5} \right) \right]$$

$$= 0.4 \left[ \log_2 5 - \log_2 (2) \right] + 0.6 \left[ \log_2 (5) - \log_2 (3) \right]$$

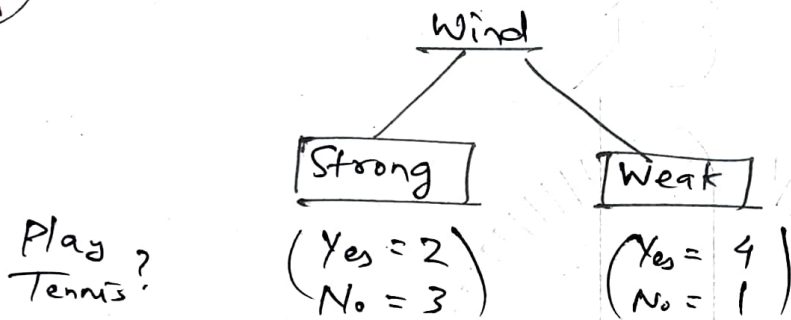
$$= 0.4 \left[ \log_2 5 - \log_2 (2) \right] + 0.6 \left[ \log_2 (5) - \log_2 (3) \right]$$

$$= 0.4 \log_2 (5) - 0.4 \log_2 (2) - 0.6 \log_2 (3)$$

$$= 2.3 - 0.4 - (0.6 \times 1.6) = 1.9 - 0.96 = 0.94$$

$$\therefore E(\text{Humidity}) = \frac{1}{2} [0.7 + 0.94] = 0.82$$

4



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

$$= 0.82$$

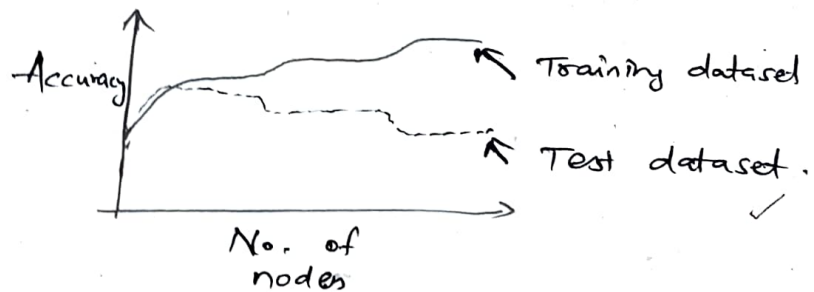
$\therefore$  Lowest entropy is obtained for  $E(\text{outlook}) = 0.8$   
Hence root will be "Outlook" for the decision tree.

4/b) A decision tree can terminate if all the tuples in the dataset have been classified and the scope of any further classification is over.

Also, the decision tree can terminate if any further classification is not possible by any branch of the decision tree.

~~Ques 2)~~ 1)c) Model overfitting can occur by training too much on a limited data set so that the model tends to generalize the training data too much and hence gives a large error for <sup>all</sup> the ~~whole data~~ other instances of data.

It can also happen that the no. of nodes in the decision tree is so high that it is unable to perform well on the test data.



Overfitting can be avoided in the decision trees by either

3) i) Prepruning: Stopping the expansion of tree before getting to the perfect classification results which effectively helps, not to over-generalize on the training dataset.

ii) Post-pruning: A "fully-grown" decision tree can be pruned by removing some of the branches so that it can retain its performance for the testing dataset too!



5) a) The Naive Bayes classifier is given by the following formula

$$P(C_i | X) = \frac{P(X | C_i) \cdot P(C_i)}{P(X)} ; \quad \begin{array}{l} X = \text{data sample with unknown class} \\ C_i = \text{Class label } C_i \end{array}$$

$$\propto P(X | C_i) P(C_i) ; \quad \text{Since } P(X) = \text{const for all classes.}$$

Here, there are two classes i.e "Yes" and "No"

and  $X = \langle \text{Outlook} = \text{Overcast}, \text{Humidity} = \text{High}, \text{Wind} = \text{Weak} \rangle$

$$\begin{aligned} P(\text{Yes} | X) &\propto P(X | \text{Yes}) \cdot P(\text{Yes}) \\ &= P(\text{Overcast} | \text{Yes}) \cdot P(\text{Humidity} = \text{High} | \text{Yes}) \cdot P(\text{Wind} = \text{Weak} | \text{Yes}) \cdot P(\text{Yes}) \end{aligned}$$

$$P(\text{Overcast} | \text{Yes}) = 0$$

$$P(\text{Humidity} = \text{High} | \text{Yes}) = \frac{1}{2}$$

$$P(\text{Yes}) = \frac{2}{5}$$

$$P(\text{Wind} = \text{Weak} | \text{Yes}) = 0$$

$$\Rightarrow P(\text{Yes} | X) = 0.$$

$$\begin{aligned} P(\text{No} | X) &= P(\text{Overcast} | \text{No}) \cdot P(\text{Humidity} = \text{High} | \text{No}) \cdot P(\text{Wind} = \text{Weak} | \text{No}) \cdot P(\text{No}) \\ &= \left(\frac{1}{3}\right) \times \left(\frac{2}{3}\right) \times \left(\frac{1}{3}\right) \times \left(\frac{3}{5}\right) = \frac{2}{45} \end{aligned}$$

Thus, prediction will be Play Tennis = "No"

5) b) There is an error in calculating  $P(X|Yes)$  since two of the probabilities turn out to be zero.

This can be corrected using Laplacian correction method

For  $C_i = \text{"Yes"}$

Outlook
Sunny
Rain

→  
Add  
"overcast"

Outlook
Sunny
Rain
Overcast

Also  $C_i \Rightarrow \text{Play tennis} = \text{"Yes"}$

3

Wind
Strong
Strong

→  
Add 1  
type of "weak"

Wind
Weak
Strong
Strong

Modified probability

$$\therefore P(\text{Overcast} | \text{Yes}) = \frac{1}{3}$$

$$P(\text{Wind} = \text{Weak} | \text{Yes}) = \frac{1}{3}$$

$$\therefore P(\text{Yes} | X) = \frac{1}{3} \times \frac{1}{2} \times \frac{1}{3} \times \frac{2}{5} = \frac{1}{45}$$

$$\text{and as before } P(\text{No} | X) = \frac{2}{45}$$

On normalisation,  $P(\text{Yes} | X) = \frac{1}{3}$  and  $P(\text{No} | X) = \frac{2}{3}$

$\therefore$  Prediction will be Play Tennis = "No."



5)c) If any feature like Humidity has continuous values, we can use this algorithm but in a different manner.

We can use Gaussian Naïve Bayes Classifier, where the probabilities for continuous variables will be given by

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} ; \quad \mu = \text{mean of the variable,}$$

$\sigma^2 = \text{Variance of the variable.}$

Rest process of finding the <sup>class</sup> ~~classifier~~ will still remain the same.

$$P(C_i | x) = P(x | C_i) \cdot P(C_i)$$
$$= P(C_i) \prod_{i=1}^K P(x_k | C_i) \quad \left\{ \begin{array}{l} \text{Assuming independence of} \\ \text{the attributes} \end{array} \right.$$

The max value of  $P(C_k | x)$   $\forall k \in \{1, 2, \dots, n\}$  will give the class label for an instance.

6) a) Following is the confusion matrix for the given model

Actual \ Prediction	Yes	No
Yes	1	1
No	1	2

True positive  $\cdot \left( \begin{matrix} \text{Actual} = Y \\ \text{Predict} = Y \end{matrix} \right) = 1$  (Sunny, Normal, Strong)

False positive  $\left( \begin{matrix} \text{Actual} = N \\ \text{Predict} = Y \end{matrix} \right) = 1$  (Overcast, Normal, Strong)

True negative  $\left( \begin{matrix} \text{Actual} = N \\ \text{Predict} = N \end{matrix} \right) = 2$  (Sunny, High, weak)  
(Rain, High, strong)

False negative  $\left( \begin{matrix} \text{Actual} = Y \\ \text{Predict} = N \end{matrix} \right) = 1$  (Rain, High, Strong)

i) Precision =  $\frac{\text{True positive}}{\text{True positive} + \text{False positive}} = \frac{1}{1+1} = \frac{1}{2} = 50\%$

ii) Recall =  $\frac{\text{True positive}}{\text{True positive} + \text{False negative}} = \frac{1}{1+1} = \frac{1}{2} = 50\%$

iii) F-score = Harmonic mean of Recall & Precision  

$$= \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$= \frac{2 \times \frac{1}{2} \times \frac{1}{2}}{\frac{1}{2} + \frac{1}{2}} = \frac{1}{2}$$

6)b)

i) Holdout method: In this method, the data is divided into two parts for training and testing separately.

If the same data is used for training & testing purpose then there are chances that model might overfit.

Hence to actually evaluate the performance a separate testing data is used.

ii) Cross-validation: The data is divided into  $K$  subsets and each time (for  $K$  iterations) one of the  $K$  subsets is used for testing and the rest  $(K-1)$  subsets are used for training.

$K=3 \rightarrow$

train	train	test
test	train	train
train	test	train

iii) In bootstrap, a fixed no. of tuples are sampled with replacement from the training dataset and the tuples which are left in this process are used for testing the model.

For example <sup>in</sup> 0.632 bootstrap, data is sampled for  $d$  times hence there are  $(1 - \frac{1}{d})^d$  no. of tuples that can be left out, which will further be used for testing.

6) c) The main purpose of ensemble of classifiers is to increase the accuracy the classification accuracy by combining multiple classifiers and then making a decision based on the majority of the classifiers are saying. This classification method can significantly increase the accuracy of models.

### Adaboost algorithm

This algorithm trains its classifiers iteratively in a sequential manner. Here, voting is based not on majority but on the weights assigned to the classifiers according to their performance. Idea here is to make any classifier to train more on those data which were incorrectly classified by its previous classifier.

#### Steps:

- 1) Assign equal weights to all the tuples
- 2) Let the classifier be tested to the data.
- 3) Find the error( $e$ ) by taking weighted avg. of errors in each tuple
- 4) Find the performance index given by  ~~$\frac{1-e}{2}$~~   $\frac{1}{2} \ln \left( \frac{1-e}{e} \right)$
- 5) Update weights of the tuples
  - i) For incorrectly classified ones: multiply by  $e^{\text{perf. index}}$
  - ii) For correctly classified ones: divide by  $e^{\text{perf. index}}$

This helps to put more emphasis on correctly classifying the incorrectly classified data of this classifier for the next classifier.