

2) a) Purity: It is defined as the ratio between the dominant class of a cluster w_i to the size of the ~~the~~ cluster.

$$\text{Purity} = \frac{\max(w_i)}{\text{Size of } |w_i|}$$

Rand Index: It can be defined as the ratio of ~~sum~~ of ~~the~~ samples that are correctly ^{clustered} ~~classified~~ as the ground truth ~~and~~ to the total number of samples.

3

let us take a matrix of size 2×2

	Same clusters	Different clusters
Same cluster in ground truth	A	B
Different cluster in ground truth	C	D

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

b) Inter cluster distance can be defined as the distance between two objects in different clusters

Intraccluster distance can be defined as the distance between two objects in same clusters.

Two Intercluster distances are :-

Single linkage :- It is the distance between any two closest objects placed in different clusters

$$\Rightarrow \min_{\substack{x_i \in S \\ x_j \in T}} (d(x_i, x_j)) \quad S \text{ \& \& } T \text{ are two clusters.}$$

Complete linkage :- It is the distance between any two farthest objects ~~in~~ placed in different clusters.

$$\Rightarrow \max_{\substack{x_i \in S \\ x_j \in T}} (d(x_i, x_j))$$

Two Intraccluster distances are :-

Diameter linkage is the distance between any two farthest objects within ~~in~~ the same cluster.

$$\Rightarrow \max_{\substack{x_i \in S \\ x_j \in S}} \delta(x_i, x_j)$$

Centroid Average linkage is defined as the two times of average distance of each object ~~in~~ within a cluster to its centroid.

$$\Rightarrow 2 \times \frac{\sum_{i=1}^n \delta(x_i, v_s)}{|S|} \quad \text{where } S \text{ is cluster.}$$

g) Dunn's cluster Validation Index is given as :

$$D.I = \min_{1 \leq i \leq K} \left\{ \min_{1 \leq j \leq K} \left\{ \frac{\delta(x_i, x_j)}{\max_{1 \leq p \leq K} (\Delta x_p)} \right\} \right\}$$

where, K = total number of clusters,

x_i represents ~~a~~ the i th cluster.

$\delta(x_i, x_j)$ = Intercluster distance between the cluster x_i & x_j

Δx_p = Intracluster distance of the cluster x_p .

The Dunn's Index is defined in such a way that it tries to maximize the intracluster distance and minimize the intercluster distance.

The usefulness of the Dunn's Index is that it is used in the cluster evaluation. It gives us a measure to determine whether the ^{number of} clusters that we have made is good or not. The Partition of the clusters that maximizes Dunn's Index is the best.

3/4] An association rule is given as :

$$X \rightarrow Y$$

Support of the rule is given as : $\frac{(X \cup Y). \text{count}}{\text{Total number of transactions.}}$

i.e. the ratio of the number of transaction that contains $(X \cup Y)$ to the total number of transactions.

Confidence = ~~1~~ - X

5a) Tuple to classify =

X = < Outlook = Overcast, Humidity = High, Wind = Weak >

Total number of distinct classes = 2.

5) ~~Plus~~ Total number of tuples = 10.

$$P(\text{Yes}) = \frac{6}{10} = \frac{3}{5}$$

$$P(\text{No}) = \frac{4}{10} = \frac{2}{5}$$

In Naïve Bayes theorem the tuple x belongs to class C_i only if

$$P(C_i/x) > P(C_j/x) \quad \text{for } 1 \leq i \leq m \text{ and } i \neq j$$

where m = total number of clusters.

$$P(C_i/x) = \frac{P(x/C_i) \times P(C_i)}{P(x)}$$

$P(x)$ is constant so we maximize $P(x/C_i) \times P(C_i)$ only.

Since Naive Bayes assumes attributes as class Independence
then,

$$P(x|c_i) = \prod_{i=1}^n P(x_i|c_i) \quad \text{where } n = \text{total number of attributes.}$$

~~$$P(\text{Sunny}/\text{Yes}) = \frac{2}{6} = \frac{1}{3}, \quad P(\text{Sunny}/\text{No}) = \frac{2}{4} = \frac{1}{2}$$~~

~~P~~

$$P(\text{Outlook} = \text{Overcast} / c_i = \text{Yes}) = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{Outlook} = \text{Overcast} / c_i = \text{No}) = \frac{0}{4} = 0.$$

$$P(\text{Humidity} = \text{High} / c_i = \text{Yes}) = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{Humidity} = \text{High} / c_i = \text{No}) = \frac{3}{4}$$

$$P(\text{Wind} = \text{Weak} / c_i = \text{Yes}) = \frac{4}{6} = \frac{2}{3}$$

$$P(\text{Wind} = \text{Weak} / c_i = \text{No}) = \frac{1}{4}$$

Hence,

~~P(x)~~

$$\therefore P(x/\text{Yes}) = \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} = \frac{2}{27}$$

$$P(x/No) = 0 + \frac{3}{4} + \frac{1}{4} = 1.$$

$$P(x/yes) + P(yes) = \frac{2}{27} \times \frac{3}{105}$$

clearly ~~$P(x/yes) > P(x/No)$~~

$P(x/yes) + P(yes)$ is greater than $P(x/No) + P(No)$.

Hence, the tuple x is classified as ~~yes~~ Play Tennis = Yes

b) Yes, there is an error in this prediction, there is not a single value for

~~Overcast~~ Outlook = Overcast & Play Tennis = No
 which is why $P(\text{Outlook} = \text{Overcast} / C = \text{No}) = 0$.

This error can be corrected using

Laplacean correction in which for each

distinct values in the attribute outlook we will add
 an extra entry in the table. Total number of samples

= 10

We add an entry for Outlook = Sunny, Rain & Overcast.
 with Play Tennis = No
 so, ~~$P(\text{Sunny}) = \frac{5}{13}$, $P(\text{Rain}) = \frac{5}{13}$ & $P(\text{Overcast}) =$~~

$$P(\text{Outlook} = \text{Sunny} / C_i = N_0) = \frac{3}{7}$$

$$P(\text{Outlook} = \text{Rain} / C_i = N_0) = \frac{3}{7}$$

$$+ P(\text{Outlook} = \text{Overcast} / C_i = N_0) = \frac{1}{7}$$

This will remove the Zero value error and will ~~can~~ give correct classification.

c) The naïve Bayes algorithm is used for both categorical and continuous values. For continuous value we can use the Gaussian distribution to calculate $P(x_i / C_i)$ where x_i is the i th attribute that have continuous values.

Gaussian distribution is given as:

$$\textcircled{3} \quad G(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

where μ = mean of the continuous range of x
 σ = standard deviation.

So, If the attribute Humidity has a continuous value then we can find the mean ~~and~~ (μ_{Humidity}) as well as Standard deviation (σ_{Humidity})

then for any value x_k in the given tuple

$$\begin{aligned} \text{we can find } P(x_k/c_i) &= g(x_k, \mu_{\text{Humidity}}, \sigma_{\text{Humidity}}) \\ &= \frac{1}{\sqrt{2\pi} \sigma_{\text{Humidity}}} e^{-\frac{(x_k - \mu_{\text{Humidity}})^2}{2\sigma_{\text{Humidity}}^2}} \end{aligned}$$

4/a) let @ positive class is Yes then

$$P = 6 \quad N = 4.$$

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

$$= -\frac{3}{5} [\log_2(3) - \log_2(5)] - \frac{2}{5} [\log_2(2) - \log_2(5)]$$

$$= -\frac{3}{5} \times [1.6 - 2.3] - \frac{2}{5} \times [1 - 2.3]$$

$$= \frac{3 \times 0.7}{5} + \frac{1.3 \times 2}{5}$$

$$= 0.94$$

$$E(A) = \sum_{i=1}^k \frac{P_i + n_i}{P + n} I(P_i, n_i)$$

For Attribute outlook,

with value Sunny \div $P_i = 2$
 $n_i = 2$

$$\begin{aligned} \therefore I(P_i, n_i) &= -\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) \\ &= -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \\ &= \cancel{\frac{1}{2}} \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

value = Overcast,

$P_i = 2$ & $n_i = 0$.

$$\therefore I(P_i, n_i) = -\frac{2}{2} \log_2\left(\frac{2}{2}\right) = -1 \times \log_2(1) = 0$$

value = Rain, $P_i = 2$ & $n_i = 2$

$$\therefore I(P_i, n_i) = 1$$

$$\begin{aligned} \text{Hence, } E(\text{outlook}) &= \frac{2+2}{10} \times 1 + 0 + \frac{2+2}{10} \times 1 \\ &= \frac{8}{10} = \frac{4}{5} \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{outlook}) &= \cancel{0.94} 0.94 - 0.8 \\ &= 0.14 \end{aligned}$$

For Attribute Humidity,

Value = Normal, $P_i = 4$ & $n_i = 1$

$$\therefore I(P_i, n_i) = -\frac{4}{5} \log_2\left(\frac{4}{5}\right) - \frac{1}{5} \log_2\left(\frac{1}{5}\right)$$

$$= -\frac{4}{5} \times 0.3 + \frac{1}{5} \times 2.3$$
$$= 0.7$$

Value = High, $P_i = 2$ & $n_i = 3$

$$\therefore I(P_i, n_i) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$= \frac{2}{5} \times 1.3 + \frac{3}{5} \times 0.7$$

$$= 0.94$$

~~Value~~

$$E(\text{Humidity}) = \frac{4+1}{10} \times 0.7 + \frac{2+3}{10} \times 0.94$$
$$= 0.82$$

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

For attribute Wind,

with value = Strong, $P_i = 2$ & ~~4~~ $n_i = 3$

$$I(P_i, n_i) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right)$$

$$= 0.94$$

value = weak, $P_i = 4$ + $m_i = 1$

$$I(P_i; m_i) = -\frac{4}{5} \log_2\left(\frac{4}{5}\right) - \frac{1}{5} \log_2\left(\frac{1}{5}\right)$$

$$= 0.7$$

$$\therefore E(\text{wins}) = \frac{5}{10} \times 0.94 + \frac{5}{10} \times 0.7$$
$$= 0.82$$

$$\therefore \text{gain}(\text{wins}) = 0.94 - 0.82$$
$$= 0.12$$

We can see that the maximum gain is for the attribute outlook with value 0.14

Hence the root of the decision tree is outlook.

b) The possible terminating criteria for the decision tree algorithm are:

i) When all the samples belong to the same class.

ii) When no samples are left.

iii) When there are no other attributes present. In this case we determine the leaf

class by majority.

c) Model overfitting occurs in the presence of the outliers and inconsistent data present in the training sample. Model tries to fit the training example hence resulting in lower accuracy in test data. The overfitting in Decision Tree is avoided by Pruning:

- i) Pre Pruning
- ii) Post Pruning.

3 In Pre Pruning, we stop the growth of the ~~the~~ tree before it perfectly classify the training data. In ~~each~~ decision tree growth for every split we determine the gain of the split, if it is ~~&~~ less than some threshold value then we stop the growth of the tree.

In post pruning, we ~~let~~ the tree grow to its maximum size then we remove the branches because of which overfitting may

occur. The branch is ~~replaced~~ replaced by its leaf. In case of multiple classes we determine the new class with majority rule.

6/6) i) Holdout Method : In this method we simply divide the sample data into training and test data generally the partition is done as $\frac{2}{3}$ data is taken as the training data and $\frac{1}{3}$ data is taken as the test data. We train the model using training data and calculate the accuracy using test data.

ii) Cross-Validation : In this method the sample data is divided into K datasets. It may use stratification as well in which the distribution of classes in each dataset is approximately equal. In the i^{th} iteration the model is trained with all the dataset except the i^{th} data set which is used as the test data to determine accuracy.

- the overall accuracy can be determined by taking the average of all accuracies.

iii) Bootstrap: In this method the ~~data~~ training data is sampled with replacement from

4 the original dataset, i.e., same training data can appear more than once and some tuples might not even be in the training data. The tuples that doesn't end up in training data will be used as test data.

~~accuracy~~

c) The main purpose of ensembles of classifiers is to increase the overall accuracy and decrease the error. It trains multiple base models and the final output is determined ~~and~~ by taking the output from each base model. Since we use multiple base models then the ~~the~~ probability of error made by ~~multiple~~ majority of the models is less.

The Adaboost algorithm works in the following way-

- i) Initially all the samples are given equal weight $\frac{1}{N}$, where N = total number of samples.
- ii) At the i th iteration we train the i th Model with weighted samples.
- iii) We determine the error of the model. If it is greater than 0.5 then we train the model again otherwise we continue.
- iv) We multiply $\frac{\text{error}(M_i)}{1 - \text{error}(M_i)}$ to ~~each~~ the weight of each sample which was correctly classified.
- v) Then we normalize the weights by dividing the new weights with old weights. for all samples. ~~of the~~ In this way the weight of the misclassified samples are increased for the $i+1$ th iteration.
- vi) We Repeat the same Process again.

How to classify a given tuple = ~~case~~

i) Calculate the weight of each model with
 $\log\left(\frac{1 - \text{error}(M_i)}{\text{error}(M_i)}\right)$.

- 4) ii) Find the class predicted by the model M_i ;
iii) ~~Assign~~ ^{Add} the weight w_i to class predicted by Model M_i ;

Now, the class with highest weight will be taken as the final prediction.

6a) For tuple, $x = \langle \text{Sunny, Normal, Strong} \rangle$, Actual ~~Prediction~~ ^{Class} = Yes

Model Prediction = Yes.

Similarly, $x = \langle \text{Overcast, Normal, Strong} \rangle$

Actual ~~Prediction~~ ^{Class} = No, Model Prediction = Yes.

$x = \langle \text{Rain, High, Strong} \rangle$

Actual ~~Prediction~~ ^{Class} = Yes, Model Prediction = No.

$x = \langle \text{Sunny, High, Weak} \rangle$

Actual class = No, Model Prediction = No

$X = \{ \text{Rain, High, Strong} \}$

Actual class = No, Model Prediction = No.

Hence, $TP = 1$, $FP = 1$

$TN = 2$, $FN = 1$

Confusion matrix is

④

	Predicted as Positive	Predicted as negative
Actual Positive class	1	1
Actual negative class	1	2

i) Precision = $\frac{TP}{TP+FP} = \frac{1}{1+1} = 0.5$

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

iii) F-Score = $\frac{\text{Precision} \times \text{Recall}}{2 \times \text{Precision} + \text{Recall}}$

$$= \frac{0.5 \times 0.5}{1 + 0.5} = \frac{0.5 \times 0.5}{1.5} = \frac{0.5}{3}$$