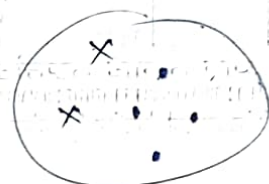② 

(a)

(i) <u>Purity</u> : Purity of a cluster determines the goodness of a cluster. It is the ratio of max objects of a class in the cluster to the size of the cluster.

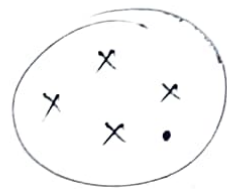$$Purity\ (\omega_i) = \frac{1}{n_i}\ max_j\ \{n_{ij}\}$$

$n_i$ = size of ith cluster

$n_{ij}$ = size of jth class in ith cluster.

<u>Example</u> ,



cluster 1

$purity\ (\omega_1) = \frac{4}{6}$

cluster 2

$purity\ (\omega_2) = \frac{4}{5}$ .

(ii) <u>Rand Index</u> : It determines the percentage of correct classification out of the total number of classifications by the model,

$$Rand\ Index\ (RI) = \frac{TP + TN}{TP + FN + FP + TN}$$

where, TP = True Positive , TN = True Negative
FP = False Positive , FN = False Negative ,

# Confusion Matrix



Predicted

|  | same class | different class |
|---|---|---|
| **same class** | TP | FN |
| **different class** | FP | TN |

Ground truth

---

(F) **Intercluster distance ($\delta$)**

Intercluster distance is the measure of the distance between two (closest) clusters generated by any clustering algorithm.

1. **Single linkage distance**

$$\delta(S,T) = \min \left\{ \begin{array}{c} d(x,y) \\ x\in S, y\in T \end{array} \right\} \qquad \left[ \begin{array}{c} d(x,y) = \text{distance} \\ \text{between two objects} \\ x \text{ and } y. \end{array} \right]$$

It is the closest distance between two clusters S and T which is equal to the minimum distance between an object from S and another object from T.

2. **Complete linkage distance**

$$\delta(S,T) = \max \left\{ \begin{array}{c} d(x,y) \\ x\in S, y\in T \end{array} \right\}$$

It is the farthest distance between any two objects in the two clusters S and T respectively.

\* Since, a good clustering have large S value and small Δ value the Dunn's index helps in determining the goodness of clustering using S and Δ values.

## Intra cluster distance (Δ)

Intra cluster distance is the measure of the distance between two objects within the same cluster.

1. ## Complete diameter distance

(A)
$$\Delta(s) = \max_{x, y \in s} \{ d(x, y) \}.$$ ✓

It is the max distance between two objects in a cluster S.

2. ## Average diameter distance

$$\Delta(s) = \frac{1}{|s| \cdot (|s|-1)} \left\{ \sum_{\substack{x, y \in s \\ x \neq y}} d(x, y) \right\}$$ ✓

It is defined as the average of the distances between all pairs of objects within a cluster, S.

(c) ## Dunn's Index

(2)
$$D_{index}(U) = \min_{1 \leq i \leq c} \left\{ \min_{\substack{1 \leq j \leq c \\ j \neq i}} \left\{ \frac{\delta(X_i, X_j)}{\max_{1 \leq k \leq c} \{ \Delta(X_k) \}} \right\} \right\}$$

A good cluster has large value of Dunn's index,

$\delta(X_i, X_j)$ = inter cluster distance between the clusters $X_i$ and $X_j$

$\Delta(X_k)$ = intra cluster distance between of $X_k$.

\*

(a) <u>Support</u> : It is defined as the ratio of the frequency of an itemset in a transaction database to the total number of transactions in the database.

$$\text{Support} = \frac{(X \cup Y) \cdot \text{count}}{n} \qquad \begin{bmatrix} (X \cup Y) = \text{itemset} \\ n = \text{no. of transactions} \end{bmatrix}$$

<u>Confidence</u> : It is defined as the ratio of the frequency of an itemset to the frequency of a proper non-empty subset of the itemset in the transaction database.

$$\text{Confidence} = \frac{(X \cup Y) \cdot \text{count}}{X \cdot \text{count}} \qquad [X \subseteq X \cup Y]$$

An itemset is referred to as a frequent itemset when the frequency <sup>support</sup> of the itemset $\geqslant$ minsup.

(i.e, minsupport).

An association rule is referred as an important rule when the confidence of the association rule $\geqslant$ minconf (i.e, min confidence or threshold).

(b)

| Transactions | Items |
|---|---|
| T1 | Bread, Butter |
| T2 | Bread, Milk, Butter |
| T3 | Bread, Jelly, Butter |
| T4 | Bread, Coke |
| T5 | Bread, Milk |
| T6 | Milk, Coke |

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

then,

| Transactions | Items |
|---|---|
| $T_1$ | $I_1, I_2$ |
| $T_2$ | $I_1, I_2, I_3$ |
| $T_3$ | $I_1, I_2, I_4$ |
| $T_4$ | $I_1, I_5$ |
| $T_5$ | $I_1, I_3$ |
| $T_6$ | $I_3, I_5$ |

$C_1$ :  $\{I_1\}$, $\{I_2\}$, $\{I_3\}$, $\{I_4\}$, $\{I_5\}$

$s = 5/6$ ✓  $s = 3/6$ ✓  $s = 3/6$ ✓  $s = 1/6$ ✗  $s = 2/6$ ✓

minsup = 30%,
(s = support)

$F_1$ :  $\{I_1\}$, $\{I_2\}$, $\{I_3\}$, $\{I_5\}$

$C_2$ :  $\{I_1, I_2\}$,  $\{I_1, I_3\}$,  $\{I_1, I_5\}$
$s = 3/6$ —  $s = 2/6$ ✓  $s = 1/6$ ✗

$\{I_2, I_3\}$,  $\{I_2, I_5\}$
$s = 1/6$ ✗  $s = 0/6$ ✗

$\{I_3, I_5\}$  $s = 1/6$  ✗

$F_2 : \{I_1, I_2\}, \{I_1, I_3\}$ ✓

$C_3 : \{I_1, I_2, I_3\}$

$\quad\quad s = \frac{1}{6}$ ✗

$F_3 : \quad \phi$

∴ Frequent itemset $F = F_1 \cup F_2 \cup F_3$

$$F = \{I_1\}, \{I_2\}, \{I_3\}, \{I_5\}, \{I_1, I_2\}, \{I_1, I_3\}$$

## Association rules

For $\{I_1\}, \{I_2\}, \{I_3\}, \{I_5\}$ itemsets there are no association rules as they are singleton sets.

For $\{I_1, I_2\}$, and $\{I_1, I_3\}$.

(6)

minconf = 80%.
c = confidence.

| $A \longrightarrow B$ | |
|---|---|
| $I_1 \longrightarrow I_2$ | $c = \frac{3}{5} = 0.6$ ✗ |
| $I_2 \longrightarrow I_1$ | $c = \frac{3}{3} = 1$ ✓ |
| $I_1 \longrightarrow I_3$ | $c = \frac{2}{5} = 0.4$ ✗ |
| $I_3 \longrightarrow I_1$ | $c = \frac{2}{3} = 0.66$ ✗ |

Important association rules :

$$I_2 \longrightarrow I_1$$

Butter $\longrightarrow$ Bread.

(c) Drawbacks of a-priori algorithm:

1. The space complexity for generating all association rules is exponential i.e, $O(2^m)$ where $m$ is the number of items in the itemset $I$.

2. The algorithm exploits the sparseness of data, high minsup and high minconf values.

3. Single minsup value problem: The algorithm considers all itemsets of same nature and similar frequency greater than a single minsup value. In practical, this is not always true as the nature and frequency of items in dataset may vary.

4. The algorithm generates high number of association rules which are difficult to interpret.

④

(a) Attributes: Outlook, Humidity, Wind.

Output variable: Play Tennis.

The output variable has 2 class ; Yes and No.

Let, $p$ = Yes count , = 6

$n$ = No . count = 4    ∴ $p+n = 10$.

∴ Information for the whole table,

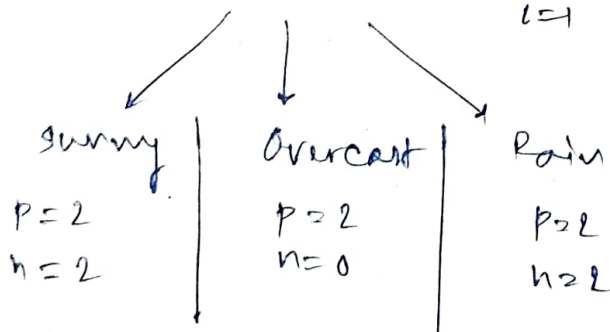$$I(p,n) = - \frac{p}{p+n} \log_2 \left(\frac{p}{p+n}\right) - \frac{n}{p+n} \log_2 \left(\frac{v}{p+n}\right)$$

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

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

$$\boxed{I(p,n) = 0.94}$$

Now entropy for each attribute,

$$E(\text{outlook}) = \sum_{i=1}^{3} \frac{p_i + n_i}{p+n} \cdot I(p_i, n_i)$$

Sunny | Overcast | Rain

$p = 2$    $p = 2$    $p = 2$
$n = 2$    $n = 0$    $n = 2$

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

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

$$= \frac{4}{10} + 0 + \frac{4}{10} = 0.8$$

$$[I(2,0) = 0 , I(2,2) = 1]$$

$$E(\text{Humidity}) = \sum_{i=1}^{2} \frac{p_i + n_i}{p+n} \cdot I(p_i, n_i)$$



Normal        High

$p = 4$        $p = 2$

$n = 1$        $n = 3$

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

$$= 0.5 \left[ -\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \right] +$$

$$0.5 \left[ -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right]$$

⑤

$$= 0.5 (0.24 + 0.46) + 0.5 (0.52 + 0.42)$$

$$= 0.82$$

$$\boxed{E(\text{Humidity}) = 0.82}$$

$$E(\text{Wind}) = \sum_{i=1}^{2} \frac{p_i + n_i}{p+n} \cdot I(p_i, n_i)$$

Strong        Weak

$p = 2$        $p = 4$

$n = 3$        $n = 1$

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

$$= 0.82$$

$$\boxed{E(\text{Wind}) = 0.82}$$

Therefore, information gain for each attributes are,

$$G(\text{Outlook}) = I(p,n) - E(\text{Outlook}) = 0.94 - 0.8 = 0.14$$

$$G(\text{Humidity}) = I(p,n) - E(\text{Humidity}) = 0.94 - 0.82 = 0.12$$

$$G(\text{Wind}) = I(p,n) - E(\text{Wind}) = 0.94 - 0.82 = 0.12$$

Since attribute 'Outlook' has the maximum information gain among all attributes, the root of the decision tree is the 'Outlook' attribute.

(b) The decision tree algorithm terminates when all instances of the dataset are classified using the attributes and conditions (classes).

(c) The main cause of model overfitting is noise learning which results in high accuracy in the training data samples and low accuracy in the test data samples. In overfitting, the model train itself with outliers or noise which have very low correlation coefficient.

In decision tree, the overfitting problem can be solved by,

1. Prepruning (Early stop): The growth of the decision tree is stopped before the completion of the whole tree.

2. Post pruning : The tree is pruned after the decision tree is completely built.

(5)

(a) $X:$ ⟨Outlook = Overcast, Humidity = High, Wind = Weak⟩

According to Naive Bayes Classification algorithm,

$$P(C_i/x) = P(C_i) \cdot P(x/C_i)$$

where,
$$P(x/C_i) = \prod_{j=1}^{k} P(x_j/C_i)$$

$C_i$ = $i$th class of output variable

$x_j$ = $j$th ~~class of~~ term in test sample,

$k$ = number of terms in the test sample.

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

$$P(C_2) = P(C = No) = \frac{4}{10} = 0.4$$

| $C = Yes$ | $C = No$ |
|---|---|
| $P\left(\dfrac{Outlook = Overcast}{C = Yes}\right)$ | $P\left(\dfrac{Outlook = Overcast}{C = No}\right)$ |
| $= \dfrac{2}{6} = 0.33$ | $= \dfrac{0}{4} = 0$ |
| $P\left(\dfrac{Humidity = High}{C = Yes}\right)$ | $P\left(\dfrac{Humidity = High}{C = No}\right)$ |
| $= \dfrac{2}{6} = 0.33$ | $= \dfrac{3}{4} = 0.75$ |
| $P\left(\dfrac{Wind = Weak}{C = Yes}\right)$ | $P\left(\dfrac{Wind = Weak}{C = No}\right)$ |
| $= \dfrac{4}{6} = 0.66$ | $= \dfrac{1}{4} = 0.25$ |

$$\therefore P(x/c_1) = 0.33 \times 0.33 \times 0.66 = 0.0718$$

$$P(x/c_2) = 0 \times 0.75 \times 0.25 = 0$$

**4.2** $P(c_1/x) = 0.6 \times 0.0718 = 0.043$

$$P(c_2/x) = 0.4 \times 0 = 0$$

$\therefore$ According to Naive Bayes algorithm, the predicted class of the sample $x$ is Yes, as, $P(c_1/x)$ is greater,

i.e. PlayTennis = Yes, for sample $x$,

(b) Yes, there is an error in the prediction of the sample class, as $P(\text{Outlook} = \text{overcast}/c = No) = 0$,

**2.2** the error can be corrected using Laplacian correction.

1. We inject a tuple, with 'Outlook' attribute set to Overcast and 'PlayTennis' variable set to No, this will make the corresponding probabilities non zero,

2. The new probabilities (corrected) will be,

$P(\text{Outlook} = \text{Overcast}/c = No) = \dfrac{1}{5} = 0.2$

$P(\text{Humidity} = \text{High}/c = No) = \dfrac{3}{5} = 0.6$

$P(\text{wind} = \text{Weak}/c = No) = \dfrac{1}{5} = 0.2$

3. $P(C_2/x) = P(C_2) \cdot P(x/C_2)$

$$= 0.4 \times (0.2 \times 0.6 \times 0.2)$$

$$= 0.0096.$$

(4) For any continuous value attribute, the algorithm will work after discretization of the attribute. The gini-index will help in finding the prediction of a particular sample only after the attribute is discretized.

Discretization can be done in different ways including quartile divisions and range classification methods.

The humidity attribute can be clas. discretized and then the sample can be predicted using Naive Bayes algorithm.

✓