

## 1) Decision Tree :

### ALGORITHM 3.2

Input:

$T$  //Decision tree  
 $D$  //Input database

Output:

$M$  //Model prediction

DTProc algorithm:

//Simplistic algorithm to illustrate prediction technique using DT

for each  $t \in D$  do

$n$  = root node of  $T$ ;

while  $n$  not leaf node do

Obtain answer to question on  $n$  applied to  $t$ ;

Identify arc from  $t$ , which contains correct answer;

$n$  = node at end of this arc;

Make prediction for  $t$  based on labeling of  $n$ ;

## 2) Naïve Bayes :

Input:

$T$  // Decision tree

$D$  // Input database

Output:

$M$  // Model prediction

NaiveBayesProc Algorithm:

// Simplified Naive Bayes prediction algorithm

// Step 1: Calculate Class and Feature Probabilities

CalculateProbabilities( $D$ );

// Step 2: Make Predictions

for each data point  $t$  in  $D$  do

CalculateClassProbabilities( $t$ );

$M[t] = \text{Argmax}(P(Y | X))$  // Select the class with the highest probability

// Output the final predictions

Output  $M$

### 3) K-Means :

#### ALGORITHM 5.6

**Input:**

$D = \{t_1, t_2, \dots, t_n\}$  //Set of elements

$k$  //Number of desired clusters

**Output:**

$K$  //Set of clusters

**K-means algorithm:**

assign initial values for means  $m_1, m_2, \dots, m_k$ ;

**repeat**

    assign each item  $t_i$  to the cluster which has the closest mean;

    calculate new mean for each cluster;

**until** convergence criteria is met;

The  $K$ -means algorithm is illustrated in Example 5.4.

### 4) K-Medoids(PAM) :

#### ALGORITHM 5.8

**Input:**

$D = \{t_1, t_2, \dots, t_n\}$  //Set of elements

$A$  //Adjacency matrix showing distance between elements

$k$  //Number of desired clusters

**Output:**

$K$  //Set of clusters

**PAM algorithm:**

arbitrarily select  $k$  medoids from  $D$ ;

**repeat**

**for each**  $t_h$  not a medoid **do**

**for each** medoid  $t_i$  **do**

            calculate  $TC_{ih}$ ;

        find  $i, h$  where  $TC_{ih}$  is the smallest;

**if**  $TC_{ih} < 0$ , **then**

            replace medoid  $t_i$  with  $t_h$ ;

**until**  $TC_{ih} \geq 0$ ;

**for each**  $t_i \in D$  **do**

    assign  $t_i$  to  $K_j$ , where  $\text{dis}(t_i, t_j)$  is the smallest over all medoids;

### 5) Hierarchical Clustering – Single :

### ALGORITHM 5.2

**Input:**

$$D = \{t_1, t_2, \dots, t_n\} \quad // \text{Set of elements}$$

```
A //Adjacency matrix showing distance between elements
```

**Output:**

```
DE // Dendrogram represented as a set of ordered triples
```

MST single link algorithm:

$$d = 0$$
$$k = n$$
$$K = \{\{t_1\}, \dots, \{t_n\}\}$$

```
DE=(d,k,K); // Initially dendrogram contains each element in
             its own cluster.
```

$$M = MST(A);$$

repeat

$$\text{old}k = k; \quad i$$

$K_i, K_j$  = two clusters closest together in MST;

$$K = K - \{K_i\} - \{K_j\} \cup \{K_i \cup K_j\};$$

```
k = oldk - 1;
```

$$d = \text{dis}(K_i, K_j);$$
$$DE = DE \cup \{d, k, K\}; \quad // \text{ New set of clusters added to dendrogram.}$$
$$\text{dis}(K_i, K_j) = \infty$$

```
until k = 1
```

### 6) Hierarchical Clustering – Average :

### ALGORITHM 5.3

**Input:**

$$D = \{t_1, t_2, \dots, t_n\} \quad // \text{Set of elements}$$

```
A //Adjacency matrix showing distance between elements
```

**Output:**

```
DE // Dendrogram represented as a set of ordered triples
```

**Average link algorithm:**

$$d = 0;$$
$$k = n;$$
$$K = \{\{t_1\}, \dots, \{t_n\}\};$$

```
DE = (d, k, K); // Initially dendrogram contains each element
                // in its own cluster.
```

**repeat**

$$\text{old}k = k;$$
$$d = d + 0.5;$$

for each pair of  $K_i, K_j \in K$  do

ave = average distance between all  $t_i \in K_i$  and  $t_j \in K_j$ ;

if  $\text{ave} \leq d$ , then

$$K = K - \{K_i\} - \{K_j\} \cup \{K_i \cup K_j\};$$
$$k = \text{old}k - 1;$$
$$DE = DE \cup \{d, k, K\}; \quad // \text{ New set of clusters added}$$

to dendrogram.

```
until k = 1
```

## 7) Hierarchical Clustering – Complete :

```
Input:
D = {t1, t2, ..., tn} // Set of elements
A // Adjacency matrix showing distance between elements

Output:
DE // Dendrogram represented as a set of ordered triples

Complete Linkage Hierarchical Clustering Algorithm:
d = 0
k = n
K = {{t1}, ..., {tn}}
DE = (d, k, K); // Initially, dendrogram contains each element in its own cluster

// Calculate the complete-linkage distance matrix
D_complete = CalculateCompleteLinkageDistanceMatrix(A)

repeat
    oldk = k

    // Find two clusters with the smallest complete-linkage distance
    (Ki, Kj) = FindClustersWithMinCompleteLinkageDistance(D_complete)

    // Merge the two clusters
    K = K - {Ki} - {Kj} U {Ki U Kj}
    k = oldk - 1
    d = dis(Ki, Kj)

    // Update the dendrogram
    DE = DE U (d, k, K)

    // Set the distance between merged clusters to infinity to avoid re-merging
    dis(Ki, Kj) = ∞

until k = 1
```

## 8) Apriori :

### ALGORITHM 6.3

#### Input:

```
I      //Itemsets
D      //Database of transactions
s      //Support
```

#### Output:

```
L      //Large itemsets
```

#### Apriori algorithm:

```
k = 0; //k is used as the scan number.
L = ∅;
```

```

 $C_1 = I;$            //Initial candidates are set to be the items.
repeat
     $k = k + 1;$ 
     $L_k = \emptyset;$ 
    for each  $I_i \in C_k$  do
         $c_i = 0;$     // Initial counts for each itemset are 0.
    for each  $t_j \in D$  do
        for each  $I_i \in C_k$  do
            if  $I_i \in t_j$  then
                 $c_i = c_i + 1;$ 
        for each  $I_i \in C_k$  do
            if  $c_i \geq (s \times |D|)$  do
                 $L_k = L_k \cup I_i;$ 
     $L = L \cup L_k;$ 
     $C_{k+1} = \text{Apriori-Gen}(L_k)$ 
until  $C_{k+1} = \emptyset;$ 

```

## 9) FP-Tree

Algorithm: **FP.growth**. Mine frequent itemsets using an FP-tree by pattern fragment growth.

Input:

- $D$ , a transaction database;
- $min\_sup$ , the minimum support count threshold.

Output: The complete set of frequent patterns.

Method:

1. The FP-tree is constructed in the following steps:

- (a) Scan the transaction database  $D$  once. Collect  $F$ , the set of frequent items, and their support counts. Sort  $F$  in support count descending order as  $L$ , the list of frequent items.
- (b) Create the root of an FP-tree, and label it as "null." For each transaction  $Trans$  in  $D$  do the following. Select and sort the frequent items in  $Trans$  according to the order of  $L$ . Let the sorted frequent item list in  $Trans$  be  $[p|P]$ , where  $p$  is the first element and  $P$  is the remaining list. Call **insert\_tree** ( $[p|P]$ ,  $T$ ), which is performed as follows. If  $T$  has a child  $N$  such that  $N.item\_name = p.item\_name$ , then increment  $N$ 's count by 1; else create a new node  $N$ , and let its count be 1, its parent link be linked to  $T$ , and its node-link to the nodes with the same *item-name* via the node-link structure. If  $P$  is nonempty, call **insert\_tree** ( $P$ ,  $N$ ) recursively.

2. The FP-tree is mined by calling **FP.growth** ( $FP\_tree, null$ ), which is implemented as follows.

procedure **FP.growth** ( $Tree, \alpha$ )

- (1) if  $Tree$  contains a single path  $P$  then
- (2)   for each combination (denoted as  $\beta$ ) of the nodes in the path  $P$
- (3)     generate pattern  $\beta \cup \alpha$  with *support\_count* = minimum support count of nodes in  $\beta$ ;
- (4) else for each  $a_i$  in the header of  $Tree$  {
- (5)   generate pattern  $\beta = a_i \cup \alpha$  with *support\_count* =  $a_i.support\_count$ ;
- (6)   construct  $\beta$ 's conditional pattern base and then  $\beta$ 's conditional FP-tree  $Tree_\beta$ ;
- (7)   if  $Tree_\beta \neq \emptyset$  then
- (8)     call **FP.growth** ( $Tree_\beta, \beta$ ); }

## 10) Page Rank :

Input:

W // Web pages and their links represented as a directed graph

q // Query page

Output:

A, H // Sets of authority and hub pages

PageRank Algorithm:

// Simplified PageRank algorithm

// Step 1: Initialize PageRank Scores

InitializePageRankScores(W);

// Step 2: Perform Iterative PageRank Calculation

for i = 1 to max\_iterations do

    for each page p in W do

$PR(p) = (1 - d) + d * \sum (PR(q) / L(q))$  for all pages q linking to p

// Step 3: Identify Authority and Hub Pages

A = FindTopPagesByPageRank(PR, k); // Select top k pages based on PageRank scores as authority pages

H = FindTopPagesByPageRank(1 / PR, k); // Select top k pages based on the inverse of PageRank scores as hub pages

// Output the final sets of authority and hub pages

Output A, H

## 11) HITS

### ALGORITHM 7.1

#### Input:

$W$  //WWW viewed as a directed graph  
 $q$  //Query  
 $s$  //Support

#### Output:

$A$  //Set of authority pages  
 $H$  //Set of hub pages

#### HITS algorithm

$R = SE(W, q)$   
 $B = R \cup \{\text{pages linked to from } R\} \cup \{\text{pages that link to pages in } R\};$   
 $G(B, L) = \text{Subgraph of } W \text{ induced by } B;$   
 $G(B, L^1) = \text{Delete links in } G \text{ within same site};$   
 $x_p = \sum_q \text{ where } (q, p) \in L^1 Y_q; \quad // \text{ Find authority weights};$   
 $y_p = \sum_q \text{ where } (p, q) \in L^1 x_q; \quad // \text{ Find hub weights};$   
 $A = \{p \mid p \text{ has one of the highest } x_p\};$   
 $H = \{p \mid p \text{ has one of the highest } y_p\};$

For Information package and Designing schema refer Exp. 2

For OLAP operations refer Exp. 7