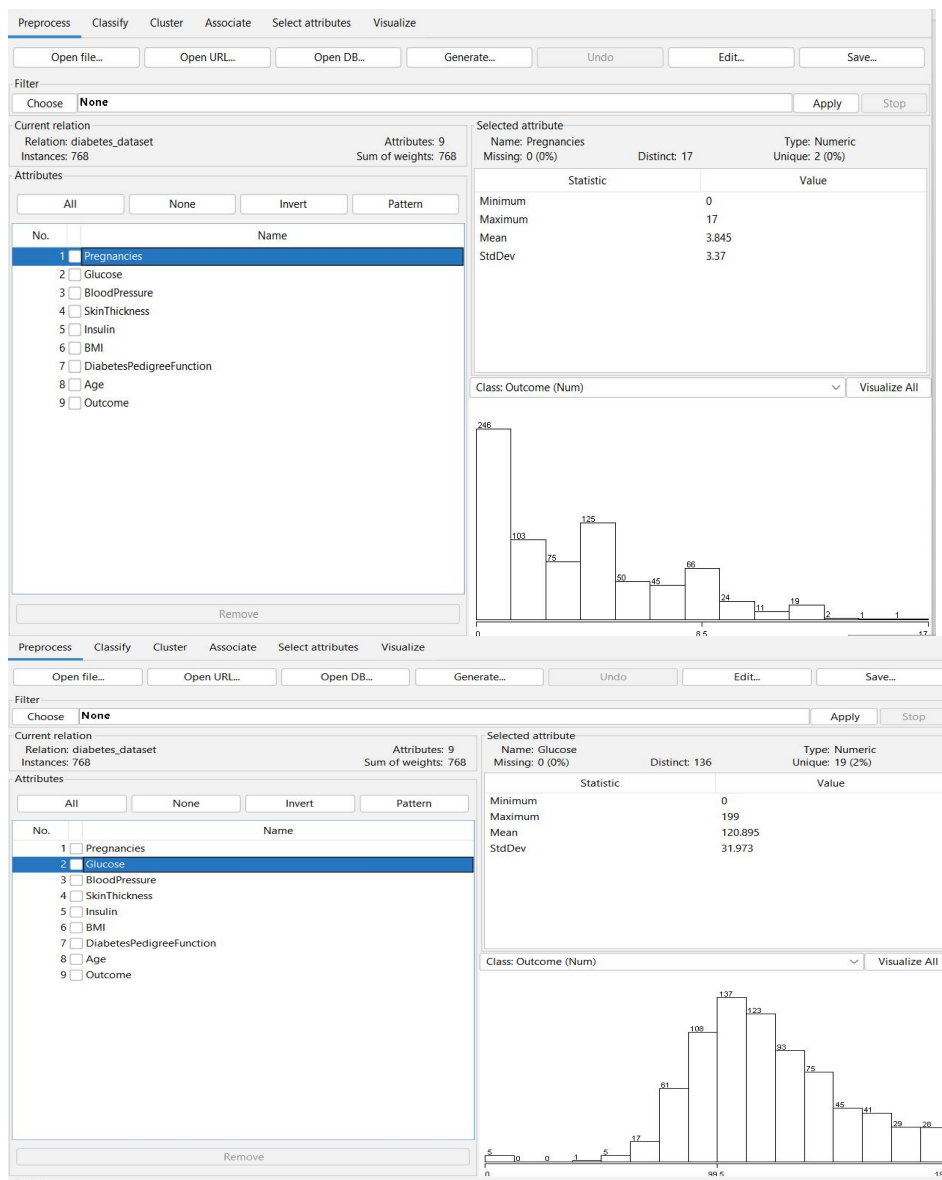**MOHAMMED RAFIK.M**
**192124179**

# CSA1668 - DWDM FOR PATTERN ANALYSIS

# WEKA EXPERIMENTS

## 1.DATA PREPROCESSING AND PREPARATION FOR KNOWLEDGE ANALYSIS USING WEKA.

### DATA PRE-PROCESSING:

It is a data mining technique which is helpful in transforming the raw data into useful and efficent data.

# 2. K-MEANS CLUSTER ANALYSIS USING WEKA.

## K-Means Clusters:

L-Means Clustering is an Unsupervised Learning algorithm, which groups the unlabeled dataset into different clusters. Here K defines the number of pre-defined clusters that need to be created in the process, as if K=2, there will be two clusters, and for K=3, there will be three clusters, and so on.

# 3. DATA ANALYSIS BY EXPECTATION MAXIMISATION ALGORITHM USING WEKA.



# 4.DATA ANALYSIS BY COBWEB-HIERARCHAL CLUSTERING ALGORITHM USING WEKA.

# 5 . KNOWLEDGE MINING USING ASSOCIATION RULE USING WEKA.



```
Preprocess   Classify   Cluster   Associate   Select attributes   Visualize

Associator
 Choose   Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1

 Start    Stop         Associator output
                       --- Run information ---
Result list (right-click for ...)
12:25:46 - Apriori     Scheme:       weka.associations.Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1
                       Relation:     dwdm sample dataset
                       Instances:    9
                       Attributes:   2
                                     T_id
                                     Transaction
                       === Associator model (full training set) ===


                       Apriori
                       =======


                       Minimum support: 0.16 (1 instances)
                       Minimum metric <confidence>: 0.9
                       Number of cycles performed: 17

                       Generated sets of large itemsets:

                       Size of set of large itemsets L(1): 15

                       Size of set of large itemsets L(2): 9

                       Best rules found:

                        1. T_id=T1 1 ==> Transaction=I1 I2 I3 1    <conf:(1)> lift:(4.5) lev:(0.09) [0] conv:(0.78)
                        2. Transaction=I2 I4 1 ==> T_id=T2 1    <conf:(1)> lift:(9) lev:(0.1) [0] conv:(0.89)
                        3. T_id=T2 1 ==> Transaction=I2 I4 1    <conf:(1)> lift:(9) lev:(0.1) [0] conv:(0.89)
                        4. T_id=T3 1 ==> Transaction=I2 I3 1    <conf:(1)> lift:(4.5) lev:(0.09) [0] conv:(0.78)
                        5. Transaction=I1 I2 I4 1 ==> T_id=T4 1    <conf:(1)> lift:(9) lev:(0.1) [0] conv:(0.89)
                        6. T_id=T4 1 ==> Transaction=I1 I2 I4 1    <conf:(1)> lift:(9) lev:(0.1) [0] conv:(0.89)
                        7. T_id=T5 1 ==> Transaction=I1 I3 1    <conf:(1)> lift:(4.5) lev:(0.09) [0] conv:(0.78)
                        8. T_id=T6 1 ==> Transaction=I2 I3 1    <conf:(1)> lift:(4.5) lev:(0.09) [0] conv:(0.78)
                        9. T_id=T7 1 ==> Transaction=I1 I3 1    <conf:(1)> lift:(4.5) lev:(0.09) [0] conv:(0.78)
                       10. Transaction=I1 I2 I3 I5 1 ==> T_id=T8 1    <conf:(1)> lift:(9) lev:(0.1) [0] conv:(0.89)
```

# 6. PREDICTION OF CATEGORICAL DATA USING DECISION TREE ALGORITHM USING  WEKA.

 ABOUT DATASET:

A dataset named Purchase is used in this decision tree which contains attributes of holidays, discount, free delievery and purchase.

# RESULT:





Preprocess   Classify   Cluster   Associate   Select attributes   Visualize

**Classifier**

Choose   **J48** -C 0.25 -M 2

**Test options**
- Use training set
- Supplied test set   Set...
- Cross-validation   Folds   10
- Percentage split   %   66

More options...

(Nom) Purchase

Start   Stop

**Result list (right-click for options)**

12:34:48 - trees.J48

**Classifier output**

```
=== Run information ===

Scheme:       weka.classifiers.trees.J48 -C 0.25 -M 2
Relation:     Purchase_new
Instances:    30
Attributes:   4
              Holiday
              Discount
              Free Delivery
              Purchase
Test mode:    10-fold cross-validation

=== Classifier model (full training set) ===

J48 pruned tree
------------------

Discount = Yes: Yes (20.0/1.0)
Discount = No
|   Free Delivery = Yes: Yes (5.0/1.0)
|   Free Delivery = No: No (5.0/1.0)

Number of Leaves  :     3

Size of the tree :      5


Time taken to build model: 0.01 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances          25               83.3333 %
Incorrectly Classified Instances         5               16.6667 %
Kappa statistic                         0.5098
Mean absolute error                     0.2401
Root mean squared error                 0.3829
Relative absolute error                71.0385 %
```

```
=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances          25               83.3333 %
Incorrectly Classified Instances         5               16.6667 %
Kappa statistic                         0.5098
Mean absolute error                     0.2401
Root mean squared error                 0.3829
Relative absolute error                71.0385 %
Root relative squared error            93.9217 %
Total Number of Instances               30

=== Detailed Accuracy By Class ===
```

|  | TP Rate | FP Rate | Precision | Recall | F-Measure | MCC | ROC Area | PRC Area | Class |
|---|---|---|---|---|---|---|---|---|---|
|  | 0.875 | 0.333 | 0.913 | 0.875 | 0.894 | 0.512 | 0.646 | 0.834 | Yes |
|  | 0.667 | 0.125 | 0.571 | 0.667 | 0.615 | 0.512 | 0.646 | 0.433 | No |
| Weighted Avg. | 0.833 | 0.292 | 0.845 | 0.833 | 0.838 | 0.512 | 0.646 | 0.754 | |

```
=== Confusion Matrix ===

  a  b   <-- classified as
 21  3 |  a = Yes
  2  4 |  b = No
```

# 7. PREDICTION OF CATEGORICAL DATA USING SMO ALGORITHM USING WEKA.



```
=== Detailed Accuracy By Class ===

              TP Rate  FP Rate  Precision  Recall  F-Measure  MCC     ROC Area  PRC Area  Class
              0.875    1.000    0.778      0.875   0.824      -0.167  0.438     0.781     Yes
              0.000    0.125    0.000      0.000   0.000      -0.167  0.438     0.200     No
Weighted Avg. 0.700    0.825    0.622      0.700   0.659      -0.167  0.438     0.664

=== Confusion Matrix ===

  a   b   <-- classified as
 21   3 |  a = Yes
  6   0 |  b = No
```
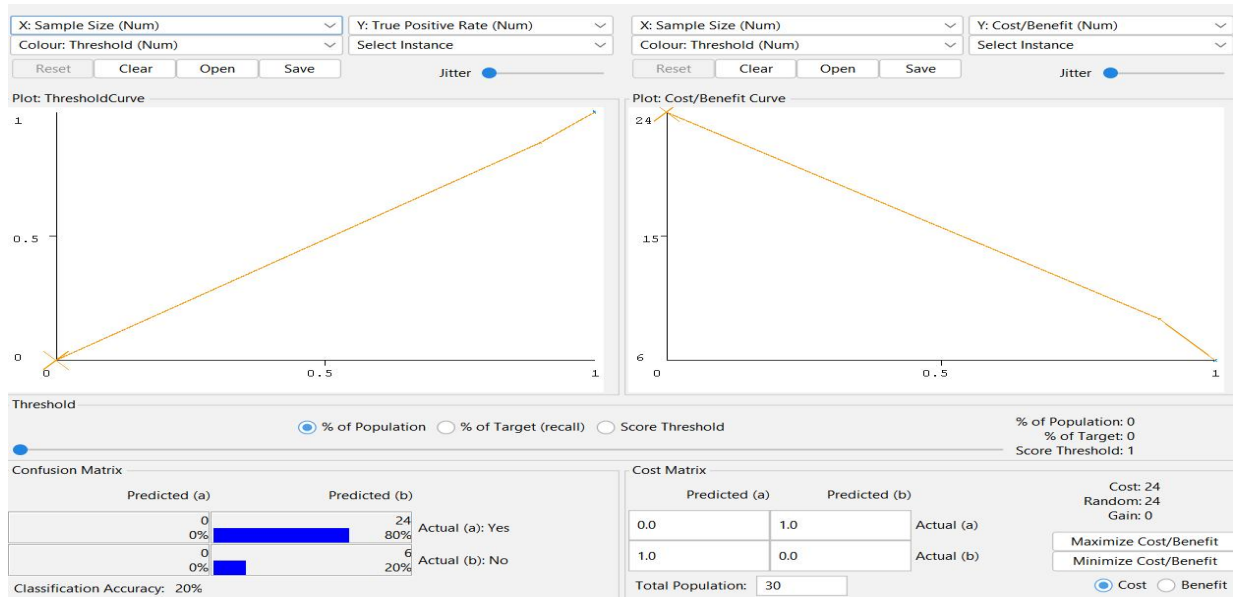
# 8. PREDICTION OF CATERGORICAL DATA USING BAYESIAN ALGORITHM

```
Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances          27               90      %
Incorrectly Classified Instances         3               10      %
Kappa statistic                          0.6667
Mean absolute error                      0.2148
Root mean squared error                  0.3327
Relative absolute error                 63.5572 %
Root relative squared error             81.5945 %
Total Number of Instances               30

=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Class
                0.958    0.333    0.920      0.958   0.939      0.671    0.757     0.867     Yes
                0.667    0.042    0.800      0.667   0.727      0.671    0.757     0.556     No
Weighted Avg.   0.900    0.275    0.896      0.900   0.896      0.671    0.757     0.804

=== Confusion Matrix ===

  a  b   <-- classified as
 23  1 |  a = Yes
  2  4 |  b = No
```

# 9. DATA ANALYSIS BY DENSITY BASED CLUSTERING ALGORITHM USING WEKA.

**10.GIVING THE FOLLOWING DATABASE WITH 5 TRANSACTIONS AND A MINIMUM SUPPORT THRESHOLD OF 60% AND A MINIMUM CONFIDENCE THRESHOLD OF 80%, FIND ALL FREQUENT ITEMSETS USING (A) APRIORI AND (B) FP-GROWTH.**

| TID | Transaction |
|---|---|
| T1 | {A, B, C, D, E, F} |
| T2 | {B, C, D, E, F, G} |
| T3 | {A, D, E, H} |
| T4 | {A, D, F, I, J} |
| T5 | {B, D, E, K} |

**RESULT:**

```
Associator output

=== Run information ===

Scheme:      weka.associations.Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1
Relation:    EXP-20
Instances:   5
Attributes:  2
             TID
             TRANSACTION
=== Associator model (full training set) ===


Apriori
=======

Minimum support: 0.3 (1 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 14

Generated sets of large itemsets:

Size of set of large itemsets L(1): 10

Size of set of large itemsets L(2): 5

Best rules found:

 1. TRANSACTION=A B C D E F 1 ==> TID=T1 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 2. TID=T1 1 ==> TRANSACTION=A B C D E F 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 3. TRANSACTION=B C D E F G 1 ==> TID=T2 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 4. TID=T2 1 ==> TRANSACTION=B C D E F G 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 5. TRANSACTION=A D E H 1 ==> TID=T3 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 6. TID=T3 1 ==> TRANSACTION=A D E H 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 7. TRANSACTION=A D F I J 1 ==> TID=T4 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 8. TID=T4 1 ==> TRANSACTION=A D F I J 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
 9. TRANSACTION=B D E K 1 ==> TID=T5 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
10. TID=T5 1 ==> TRANSACTION=B D E K 1    <conf:(1)> lift:(5) lev:(0.16) [0] conv:(0.8)
```

**11. THE 'DATABASE' BELOW HAS NINE TRANSACTIONS. WHAT ASSOCIATION RULES CAN BE FOUND IN THIS SET, IF THE MINIMUM SUPPORT (I.E COVERAGE) IS 60% AND THE MINIMUM CONFIDENCE (I.E. ACCURACY) IS 80% ?**
**TRANS_ID ITEMLIST**

| TID | List of Items |
|---|---|
| T100 | I1, I2, I5 |
| T100 | I2, I4 |
| T100 | I2, I3 |
| T100 | I1, I2, I4 |
| T100 | I1, I3 |
| T100 | I2, I3 |
| T100 | I1, I3 |
| T100 | I1, I2 ,I3, I5 |
| T100 | I1, I2, I3 |

# RESULT:

Associator

Choose | **Apriori** -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1

Start | Stop

Result list (right-click for ...)
12:25:46 - Apriori
13:19:01 - Apriori

Associator output

```
=== Run information ===


Scheme:        weka.associations.Apriori -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1 -S -1.0 -c -1
Relation:      apprior
Instances:     9
Attributes:    2
               T_ID
               List_of_items
=== Associator model (full training set) ===



Apriori
=======

Minimum support: 0.11 (1 instances)
Minimum metric <confidence>: 0.9
Number of cycles performed: 18

Generated sets of large itemsets:

Size of set of large itemsets L(1): 7


Size of set of large itemsets L(2): 6

Best rules found:

 1. List_of_items=t2 t3 2 ==> T_ID=T100 2     <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 2. List_of_items=t1 t3 2 ==> T_ID=T100 2     <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 3. List_of_items=t1 t2 t3 t5 2 ==> T_ID=T100 2     <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 4. List_of_items=t1 t2 t5 1 ==> T_ID=T100 1     <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 5. List_of_items=t2 t4 1 ==> T_ID=T100 1     <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
 6. List_of_items=t1 t2 t4 1 ==> T_ID=T100 1     <conf:(1)> lift:(1) lev:(0) [0] conv:(0)
```