# WEKA TOOL

**B.Divya**
192011067

## 1. K MEANS ALGORITHM:

**DATA SET:** diabetes

**ALGORITHM:**

1. The k means clustering algorithm computes the centroids and repeats until optimal centroid is found.
2. Firstly,provide k number of clusters.
3. Choose k data points and assign to each clusters and divide data based on data points.
4. The cluster centroids will be constructed.
5. Iterate steps until ideal centroid is found.
6. The sum of squared distances between data points and clusters should be find.
7. Allocate each data point to cluster which is closest to centroid.
8. Construct the centroids for clusters by averaging all data points of clusters.

**OUTPUT:**

## 2. DECISION TREE:

**DATA SET:** diabetes

**ALGORITHM:**

1. determine the root node.
2. Calculate the entropy of classes.
3. Calculate the entropy of the split of the attribute.
4. Calculate the information gain.
5. Perform split.
6. Perform further split.
7. Compute decision tree.

Entropy= -Σ pi log2 pi

Information gain=entropy of parent node-sum of weights of entropy of child node.

## Output:

## Weka Explorer — Window 1

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Classifier**

Choose    REPTree -M 2 -V 0.001 -N 3 -S 1 -L -1 -I 0.0

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

(Nom) class

Start    Stop

**Result list (right-click for options)**
09:34:31 - trees.REPTree

**Classifier output**

```
|   |   plas < 166.5
|   |   |   preg < 6.5
|   |   |   |   pedi < 0.33 : tested_negative (22/5) [10/5]
|   |   |   |   pedi >= 0.33
|   |   |   |   |   preg < 5.5
|   |   |   |   |   |   insu < 422.5 : tested_positive (28/8) [13/6]
|   |   |   |   |   |   insu >= 422.5 : tested_negative (2/0) [1/0]
|   |   |   |   |   preg >= 5.5 : tested_negative (5/1) [0/0]
|   |   |   preg >= 6.5
|   |   |   |   mass < 29 : tested_negative (6/3) [1/0]
|   |   |   |   mass >= 29 : tested_positive (21/0) [9/4]
|   |   plas >= 166.5 : tested_positive (52/6) [27/5]

Size of the tree : 49

Time taken to build model: 0.01 seconds

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

Correctly Classified Instances         578               75.2604 %
Incorrectly Classified Instances       190               24.7396 %
Kappa statistic                          0.438
Mean absolute error                      0.3272
Root mean squared error                  0.4289
Relative absolute error                 71.9842 %
Root relative squared error             89.9782 %
Total Number of Instances              768
```

**Status**
OK

---

## Weka Explorer — Window 2

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

**Classifier**

Choose    REPTree -M 2 -V 0.001 -N 3 -S 1 -L -1 -I 0.0

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

(Nom) class

Start    Stop

**Result list (right-click for options)**
09:34:31 - trees.REPTree

**Classifier output**

```
Time taken to build model: 0.01 seconds

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

Correctly Classified Instances         578               75.2604 %
Incorrectly Classified Instances       190               24.7396 %
Kappa statistic                          0.438
Mean absolute error                      0.3272
Root mean squared error                  0.4289
Relative absolute error                 71.9842 %
Root relative squared error             89.9782 %
Total Number of Instances              768

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 0.846    0.422    0.789      0.846   0.817      0.441  0.766     0.833     tested_negative
                 0.578    0.154    0.668      0.578   0.620      0.441  0.766     0.607     tested_positive
Weighted Avg.    0.753    0.328    0.747      0.753   0.748      0.441  0.766     0.754

=== Confusion Matrix ===

   a   b   <-- classified as
 423  77 |   a = tested_negative
 113 155 |   b = tested_positive
```

**Status**
OK

# 3. BAYESIAN CLASSIFICATION:

## DATASET:diabetes

## ALGORITHM:
1. convert given dataset into frequency table.
2. Construct livelihood tables by calculating the probabilities.
3. Use the bayes formula for calculating probabilities.

$$P(A|B) = [P(B|A) P(A)]/ P(B), \text{ where } P(B) \neq 0$$

4. now calculate the probability for all possible choices.
5. Then compare all the outputs.
6. Determine the probability which is more efficient by checking outputs.
**7.** Finally,compute the probability using bayesian classification.
## OUTPUT:

**Weka Explorer** — Classifier: NaiveBayes

Classifier output:

```
                    precision        1.0625          1.0625

plas
    mean                        109.9541        141.2581
    std. dev.                    26.1114         31.8728
    weight sum                      500             268
    precision                    1.4741          1.4741

pres
    mean                         68.1397         70.718
    std. dev.                    17.9834         21.4094
    weight sum                      500             268
    precision                    2.6522          2.6522

skin
    mean                         19.8356         22.2824
    std. dev.                    14.8974         17.6992
    weight sum                      500             268
    precision                    1.98            1.98

insu
    mean                         68.8507        100.2812
    std. dev.                    98.828         138.4883
    weight sum                      500             268
    precision                    4.573           4.573

mass
    mean                         30.3009         35.1475
    std. dev.                     7.6833          7.2537
    weight sum                      500             268
```

Status: OK



**Weka Explorer** — Classifier: NaiveBayes

Classifier output:

```
    weight sum                      500             268
    precision                    0.2717          0.2717

pedi
    mean                          0.4297          0.5504
    std. dev.                     0.2986          0.3715
    weight sum                      500             268
    precision                    0.0045          0.0045

age
    mean                         31.2494         37.0808
    std. dev.                    11.6059         10.9146
    weight sum                      500             268
    precision                    1.1765          1.1765


Time taken to build model: 0.01 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0.02 seconds

=== Summary ===

Correctly Classified Instances         201              77.0115 %
Incorrectly Classified Instances        60              22.9885 %
Kappa statistic                          0.4631
Mean absolute error                      0.266
```

Status: OK



**Weka Explorer** — Classifier: NaiveBayes

Classifier output:

```
=== Evaluation on test split ===

Time taken to test model on test split: 0.02 seconds

=== Summary ===

Correctly Classified Instances         201              77.0115 %
Incorrectly Classified Instances        60              22.9885 %
Kappa statistic                          0.4631
Mean absolute error                      0.266
Root mean squared error                  0.3822
Relative absolute error                 58.9747 %
Root relative squared error             81.6432 %
Total Number of Instances              261

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 0.843    0.386    0.824      0.843   0.833      0.463  0.854     0.918     tested_negative
                 0.614    0.157    0.646      0.614   0.630      0.463  0.854     0.760     tested_positive
Weighted Avg.    0.770    0.313    0.767      0.770   0.769      0.463  0.854     0.868

=== Confusion Matrix ===

   a    b   <-- classified as
 150   28 |   a = tested_negative
  32   51 |   b = tested_positive
```
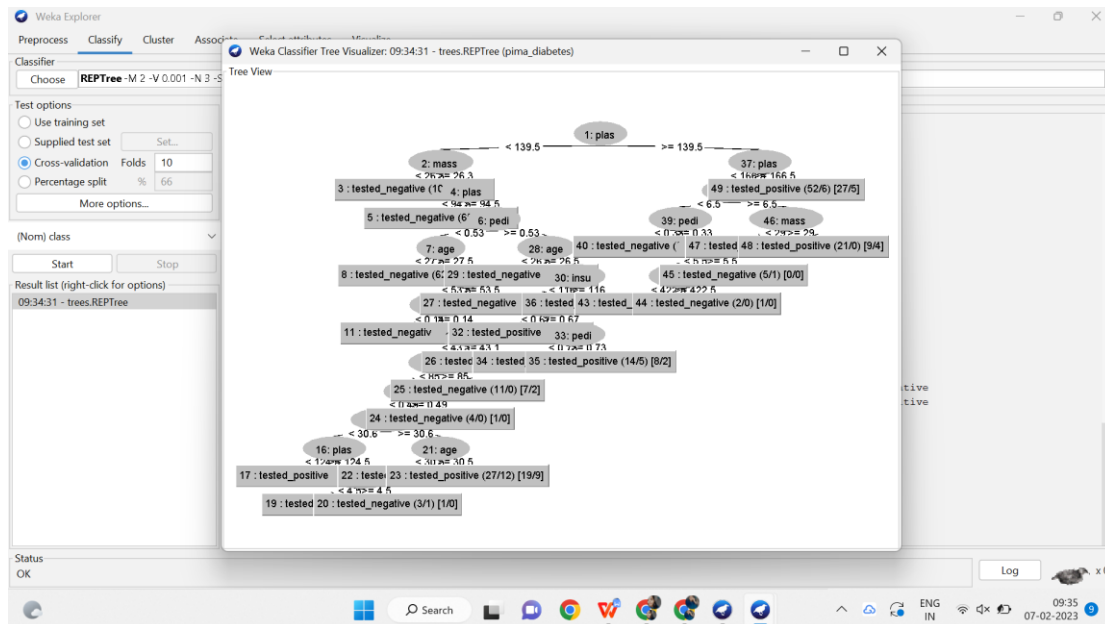
Status: OK

# 4. APRIORI ALGORITHM:

## DATASET: supermarket

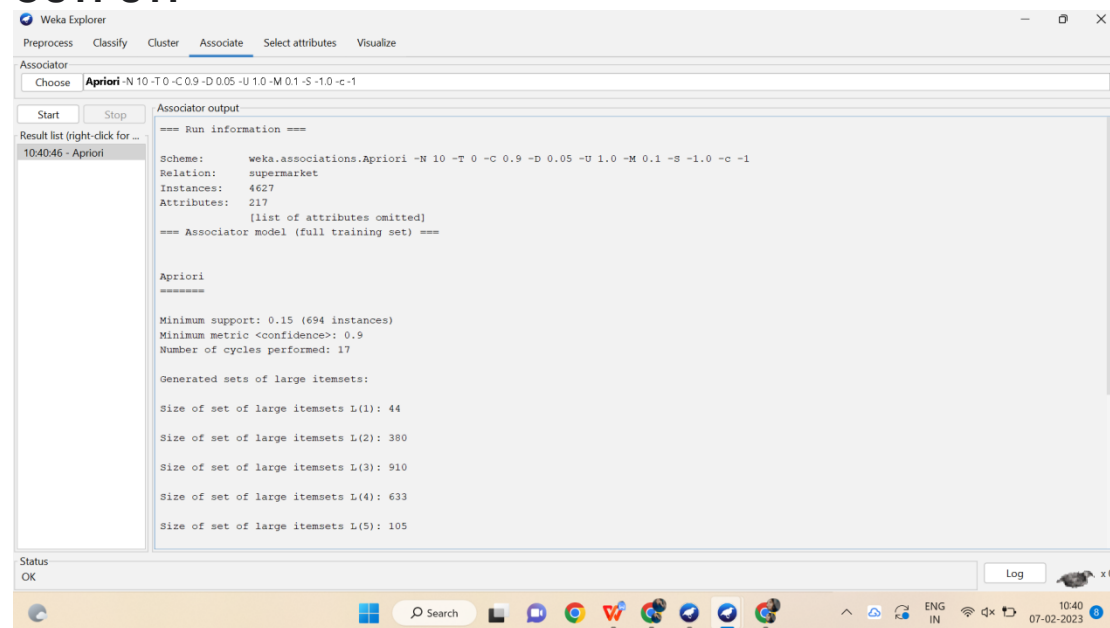## ALGORITHM:
1. firstly,convert the given transactional database into an frequency table.
2. Assign any minimum support to the frequency table,in which contains item sets and suppor count.
3. The item sets and support count is combinely called as candidate set.
4. Now,check the support count with the minimum support.
5. Remove the support count which is less than minimum support and write the remaining item sets in descending order.
6. Again checking by combining two itemsets.
7. iterate the steps until the support count should be equal to minimum support.

   Confidence=support(A∩B)/support(A)
8. calculate the confidence and convert it into percentage.
9. Finally,check which is more efficient.

## OUTPUT:

## 5. FP GROWTH ALGORITHM:

**DATASET:** supermarket

**ALGORITHM:**
1. firstly,convert the given transactional database into an frequency table.
2. Assign any minimum support to the frequency table,in which contains itemsets and support count.
3. The item sets and support count is combinely called as candidate set.
4. Now,check the support count with the minimum support.
5. Remove the support count which is less than minimum support and write remaining items in descending order.
6. Find the ordered item set using frequency table.
7. Construct the FP gowth using the ordered item set.
8. Then compute the conditionally pattern using FP grpowth.
9. Again find the conditionally frequency pattern.
10. Finally compute the FP gtowth algorithm.

**OUTPUT:**

Preprocess   Classify   Cluster   **Associate**   Select attributes   Visualize

**Associator**

Choose   **FPGrowth** -P 2 -I -1 -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1

Start   Stop

**Result list (right-click for ...**

10:40:46 - Apriori
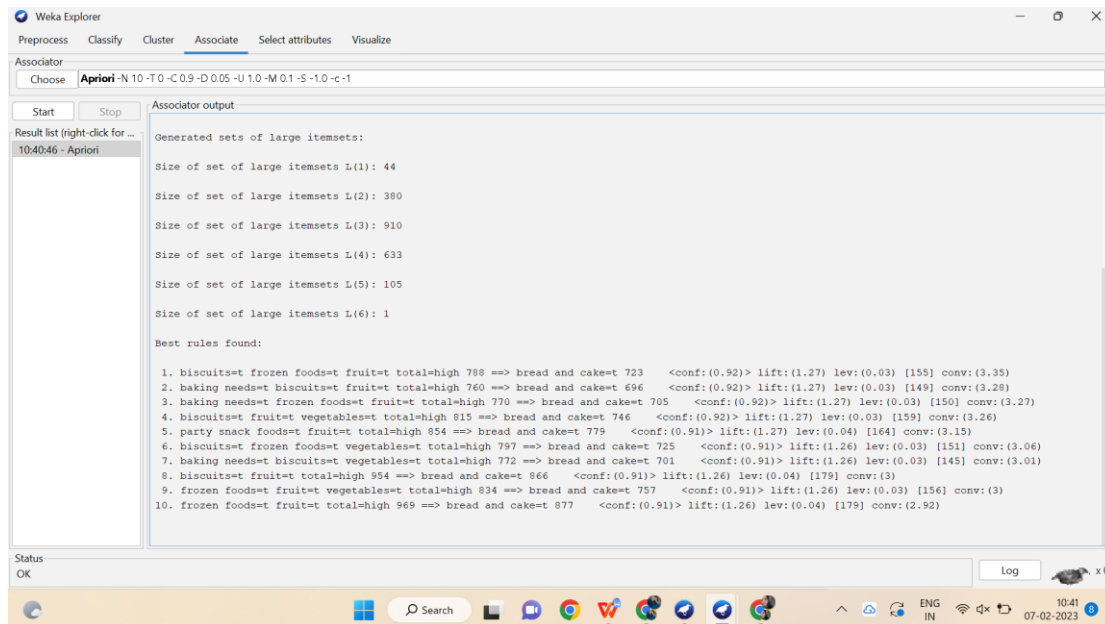10:55:56 - FPGrowth

**Associator output**

```
=== Run information ===

Scheme:        weka.associations.FPGrowth -P 2 -I -1 -N 10 -T 0 -C 0.9 -D 0.05 -U 1.0 -M 0.1
Relation:      supermarket
Instances:     4627
Attributes:    217
               [list of attributes omitted]
=== Associator model (full training set) ===


FPGrowth found 16 rules (displaying top 10)

 1. [fruit=t, frozen foods=t, biscuits=t, total=high]: 788 ==> [bread and cake=t]: 723    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.35)
 2. [fruit=t, baking needs=t, biscuits=t, total=high]: 760 ==> [bread and cake=t]: 696    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.28)
 3. [fruit=t, baking needs=t, frozen foods=t, total=high]: 770 ==> [bread and cake=t]: 705    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.27)
 4. [fruit=t, vegetables=t, biscuits=t, total=high]: 815 ==> [bread and cake=t]: 746    <conf:(0.92)> lift:(1.27) lev:(0.03) conv:(3.26)
 5. [fruit=t, party snack foods=t, total=high]: 854 ==> [bread and cake=t]: 779    <conf:(0.91)> lift:(1.27) lev:(0.04) conv:(3.15)
 6. [vegetables=t, frozen foods=t, biscuits=t, total=high]: 797 ==> [bread and cake=t]: 725    <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3.06)
 7. [vegetables=t, baking needs=t, biscuits=t, total=high]: 772 ==> [bread and cake=t]: 701    <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3.01)
 8. [fruit=t, biscuits=t, total=high]: 954 ==> [bread and cake=t]: 866    <conf:(0.91)> lift:(1.26) lev:(0.04) conv:(3)
 9. [fruit=t, vegetables=t, frozen foods=t, total=high]: 834 ==> [bread and cake=t]: 757    <conf:(0.91)> lift:(1.26) lev:(0.03) conv:(3)
10. [fruit=t, frozen foods=t, total=high]: 969 ==> [bread and cake=t]: 877    <conf:(0.91)> lift:(1.26) lev:(0.04) conv:(2.92)
```

**Status**

OK   Log   x 0