

## **Practical: 1**

**Aim:** Installation of WEKA Tools.

**Code:**

**Data Mining tools**

- Data Mining is the set of techniques that utilize specific algorithms, statical analysis, artificial intelligence, and database systems to analyses data from different dimensionsand perspectives.
- Picturing data mining as a machine, the raw data becomes the input, the data mining activity becomes the task the machine is designed to do and the output from the machine is actionable data, in other words, data that can be used to take strategic or tactical decisions, positively impacting the bottom line.
- So, what does the machine itself in the below figure represent? The machine in this oversimplified model is the tool that is used to execute the various methods and techniques used in data mining.
- Different types of Data Mining Tools are:
  1. Rapid Miner
  2. Oracle Data Mining
  3. IBM SPSS Modeler
  4. Knime
  5. Python
  6. Orange
  7. Kaggle
  8. Rattle
  9. Weka
  10. Teradata
  11. H2O
  12. Apache Spark
  13. Sisense
  14. Xplenty

## WEKA

(Waikato Environment for Knowledge Learning)

- Everybody talks about Data Mining and Big Data nowadays. Weka is a powerful, yet easy to use tool for machine learning and data mining.
  - This course provides a deeper account of data mining tools and techniques.
  - The emphasis is on principles and practical data mining using Weka, rather than mathematical theory or advanced details of particular algorithms.
  - Students will work with multimillion-instance datasets, classify text, experiment with clustering, association rules, etc.
- ❖ Step to Install WEKA:

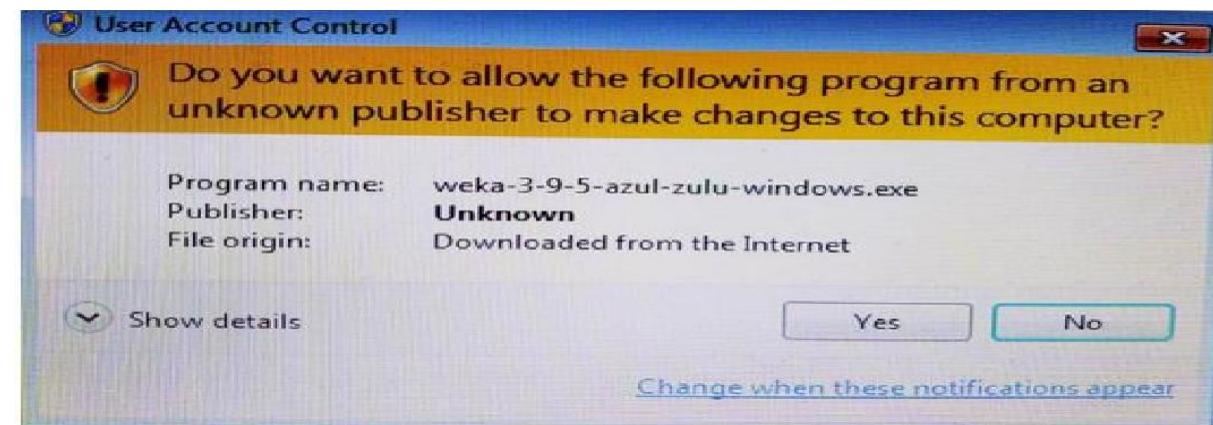
The screenshot shows a web browser displaying the WEKA website at [www.cs.waikato.ac.nz/ml/weka/downloading.html](http://www.cs.waikato.ac.nz/ml/weka/downloading.html). The page features the WEKA logo (a kiwi bird) and the text 'WEKA The University of Waikato'. A navigation bar includes links for Project, Software (which is underlined), Book, Publications, People, and Related. The main content section is titled 'Downloading and installing Weka'. It contains text about the two versions of Weka (3.8 stable and 3.9 development), the package management system, and notes for upgrading users. There are also two notes for users upgrading from 3.7 to 3.8.

There are two versions of Weka: Weka 3.8 is the latest stable version, and Weka 3.9 is the development version. For the bleeding edge, it is also possible to download nightly snapshots.

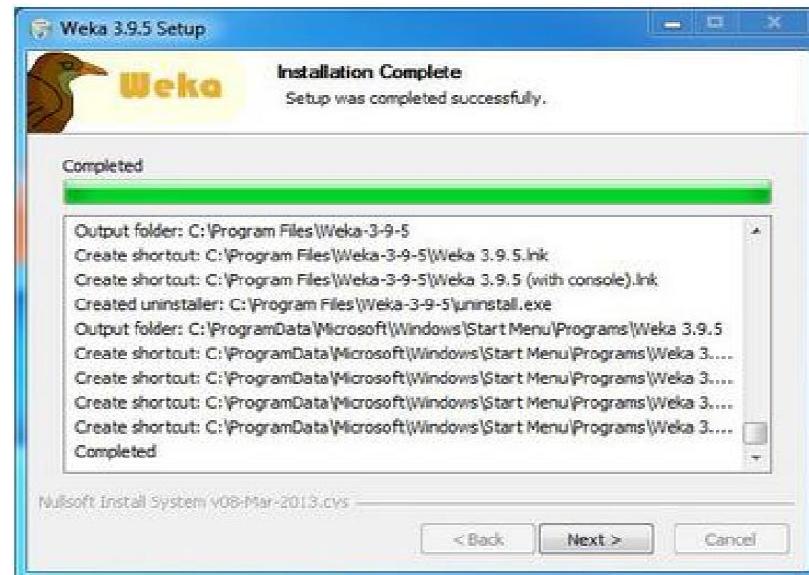
Stable versions receive only bug fixes, while the development version receives new features. Weka 3.8 and 3.9 feature a package management system that makes it easy for the Weka community to add new functionality to Weka. The package management system requires an internet connection in order to download and install packages.

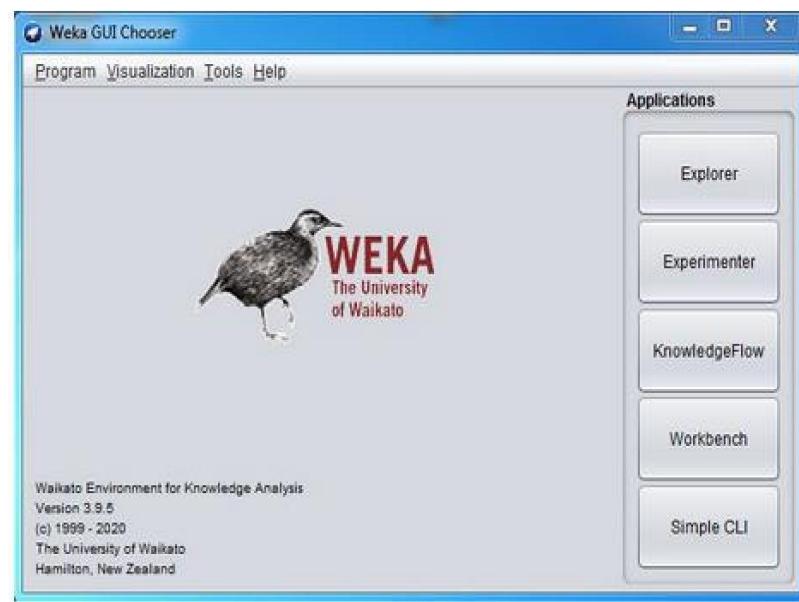
Note (1) for users upgrading from Weka 3.7 to Weka 3.8 or later: if the Weka 3.8 package manager does not start up, please delete the file `installedPackageCache.ser` in the `packages` folder that resides in the `wekafiles` folder in your user home.

Note (2) for users upgrading from Weka 3.7 to Weka 3.8 or later: serialized models created in 3.7 are not compatible with 3.8. We have a **model migrator** tool that can migrate some models to be compatible with 3.8.0. One exception is RandomForest, which can be migrated up to 3.7.13 but no further. Usage is as follows:









❖ **Introduction to WEKA:**

Weka contains a collection of visualization tools and algorithms for data analysis and predictive modelling, together with graphical user interfaces for easy access to these functions.

The original non-Java version of Weka was a Tcl/Tk front-end to (mostly third-party) modelling algorithms implemented in other programming languages, plus data preprocessing utilities in C and a make file-based system for running machine learning experiments.

This original version was primarily designed as a tool for analyzing data from agricultural domains. Still, the more recent fully Java-based version (Weka 3), developed in 1997, is now used in many different application areas, particularly for educational purposes and research. Weka has the following advantages, such as:

- Free availability under the GNU General Public License.
- Portability, since it is fully implemented in the Java programming language and thus runs on almost any modern computing platform.
- A comprehensive collection of data preprocessing and modelling techniques.
- Ease of use due to its graphical user interfaces.

Weka supports several standard data mining tasks, specifically, data preprocessing, clustering, classification, regression, visualization, and feature selection. Input to Weka is expected to be formatted according to the Attribute-Relational File Format and filename with the .arff extension.

## ❖ WEKA's Features

### **1. Preprocess**

The preprocessing of data is a crucial task in data mining. Because most of the data is raw, there are chances that it may contain empty or duplicate values, have garbage values, outliers, extra columns, or have a different naming convention. All these things degrade the results.

To make data cleaner, better and comprehensive, WEKA comes up with a comprehensive set of options under the *filter* category. Here, the tool provides both supervised and unsupervised types of operations.

### **2. Classification**

Classification is one of the essential functions in machine learning, where we assign classes or categories to items. The classic examples of classification are: declaring a brain tumor as “malignant” or “benign” or assigning an email to “spam” or “not\_spam” class.

After the selection of the desired classifier, we select test options for the training set.

### **3. Cluster**

In clustering, a dataset is arranged in different groups/clusters based on some similarities. In this case, the items within the same cluster are identical but different from other clusters. The examples of clustering include, but are not limited to, identifying customers with similar behaviors, organizing the regions according to homogenous land use.

When we apply the EM algorithm, the tool shows the mean and standard deviation of the clusters and attributes:

### **4. Associate**

Association rules highlight all the associations and correlations between items of a dataset. In short, it is an *if-then* statement that depicts the probability of relationships between data items. A classic example of association refers to a connection between the sale of milk and bread.

### **5. Visualize**

In the visualize tab, different plot matrix and graphs are available to show the trends and errors identified by the model.

## Practical: 2

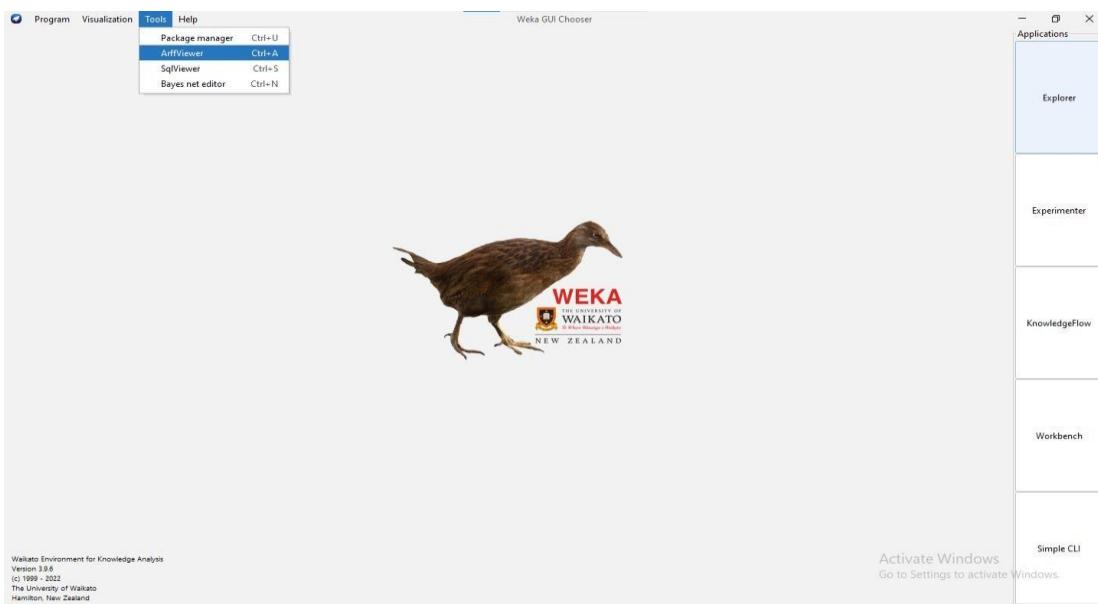
**Aim:** Create your own excel file. Convert the excel file to .csv format and prepare it as ARFF files.

**Code:**

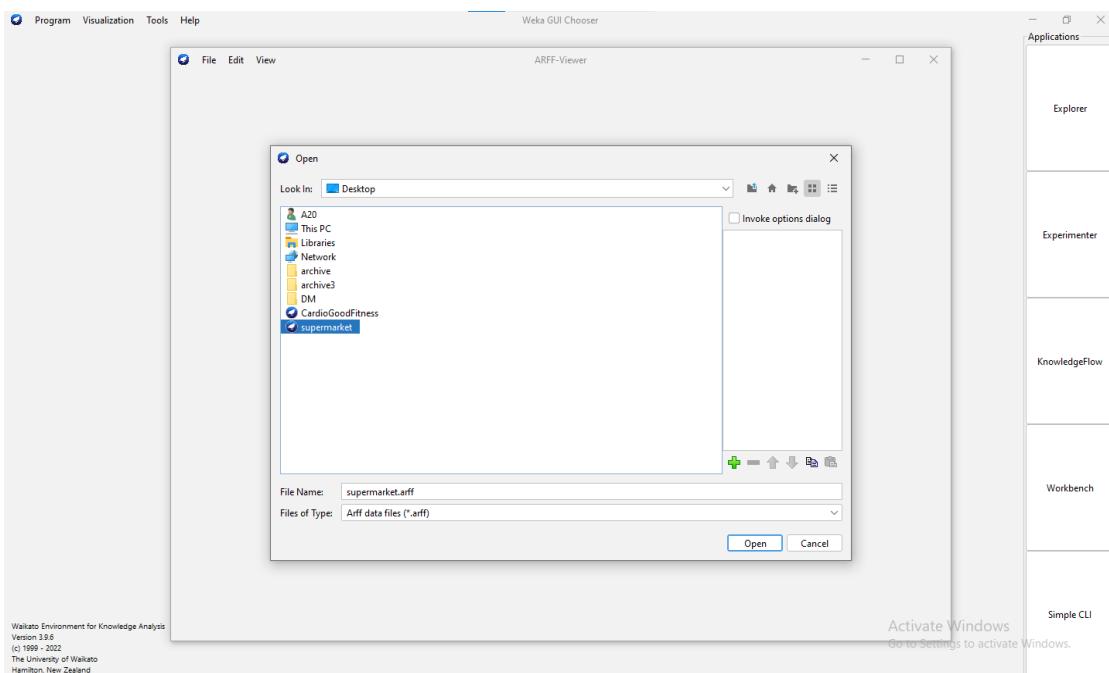
**Steps:**

### 1) Open WEKA

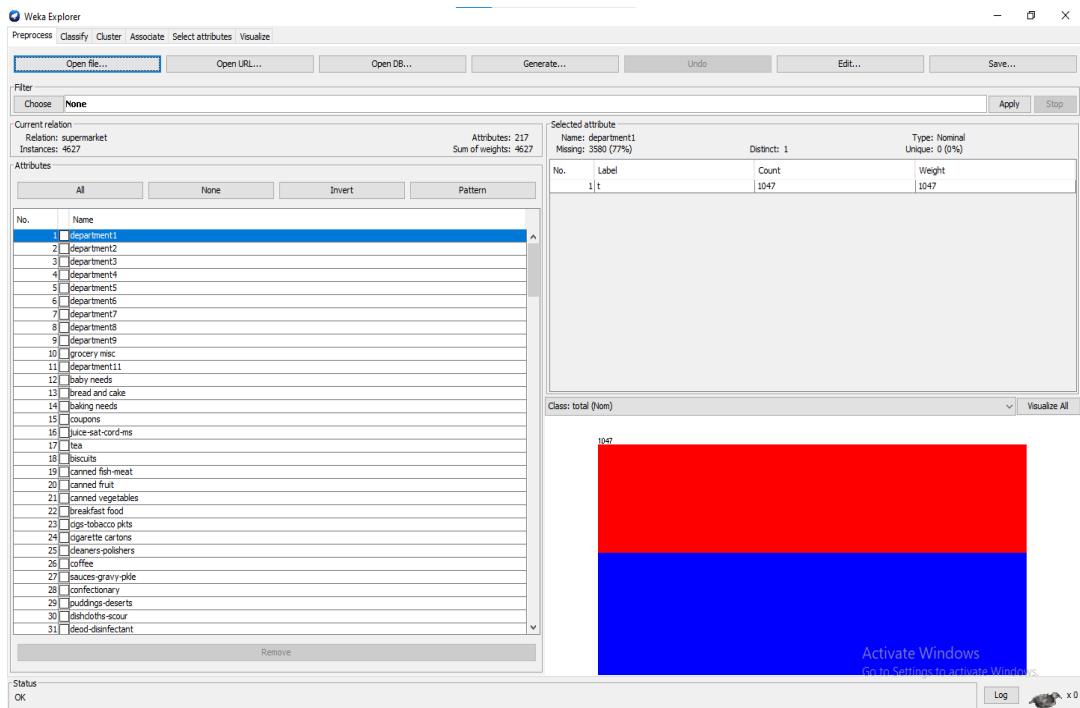
Select Tools and ARFF viewer



### 2) Select the file



### 3) Open the file in weka



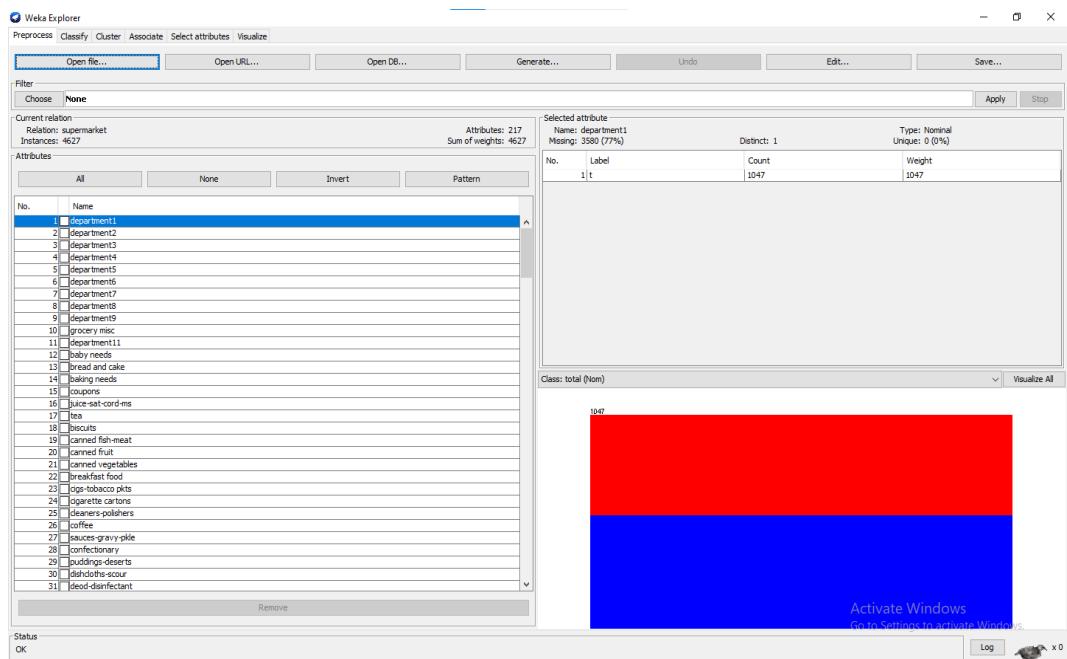
## Practical: 3

**Aim:** Preprocess and classify your dataset.

**Code:**

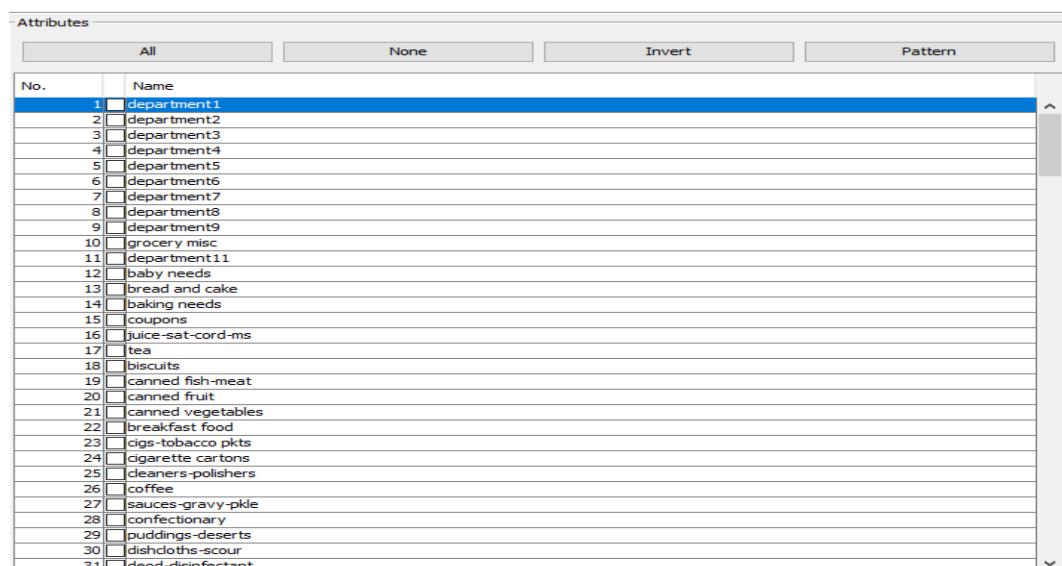
**Steps:**

### 1) Open data in WEKA



### 2) Understanding Data

Let us first look at the highlighted Current relation sub window. It shows the name of the database that is currently loaded. You can infer two points from this sub window –On the left side, notice the Attributes sub window that displays the various fields in the database.



3) In the Selected Attribute sub window, you can observe the following –

Selected attribute		Type: Nominal
Name: biscuits		Unique: 0 (0%)
Missing: 2022 (44%)	Distinct: 1	
No.	Label	Count
1	t	2605
		Weight
		2605

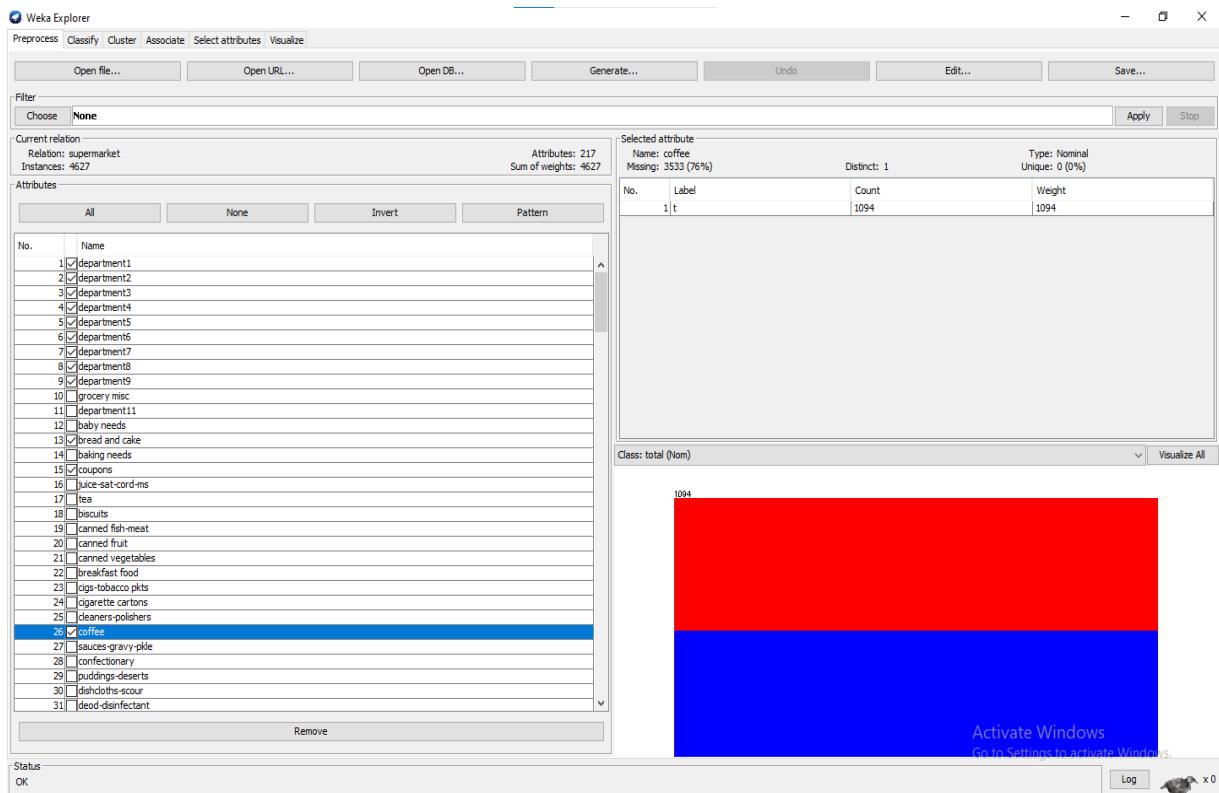
4) At the bottom of the window, you see the visual representation of the classvalues.

If you click on the Visualize All button, you will be able to see all features in one single window as shown here –



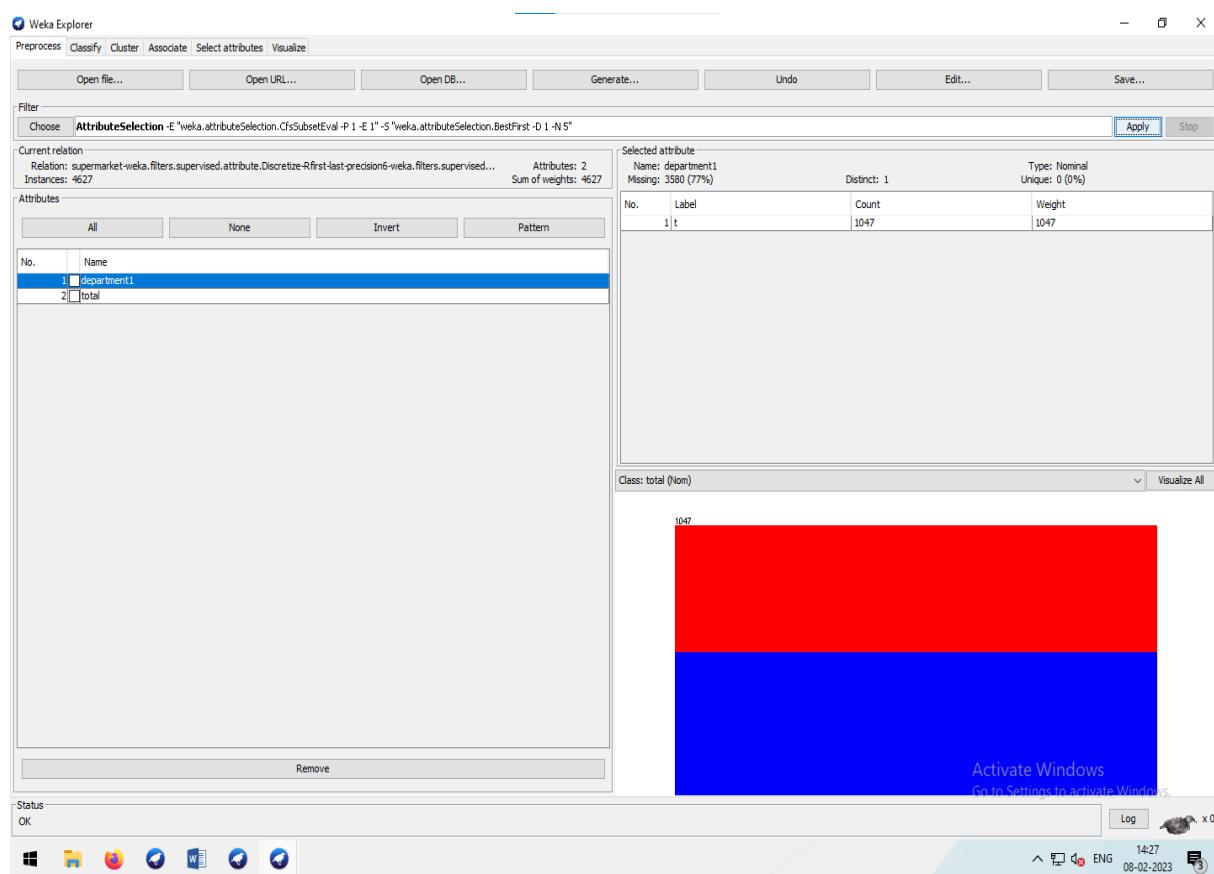
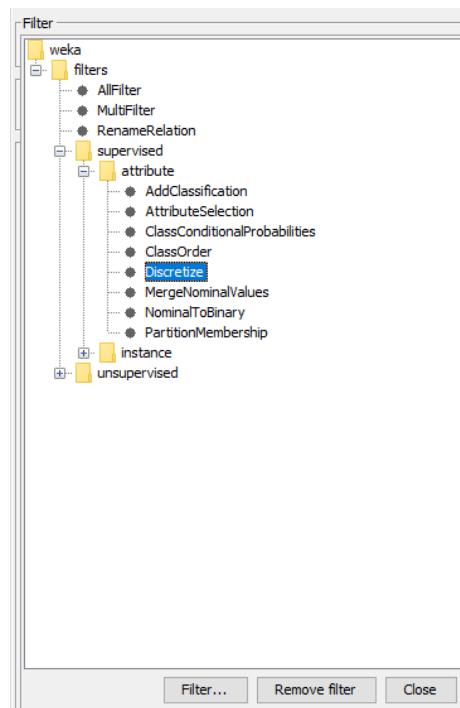
## 5) Removing Attributes

Many a time, the data that you want to use for model building comes with many irrelevant fields. For example, the customer database may contain his mobile number which is relevant in analyzing his credit rating.



## 6) Applying Filters

Some of the machine learning techniques such as association rule mining requires categorical data. We will convert these to nominal by applying a filter on our raw data. Click on the Choose button in the Filter sub window and select the following filter –

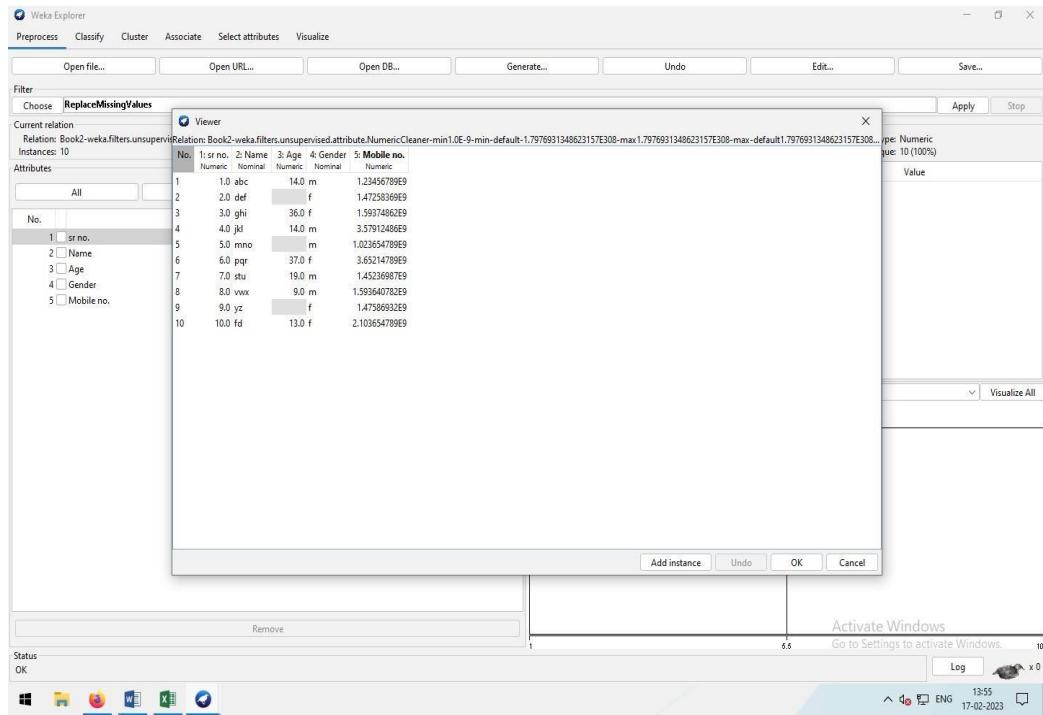


## Practical: 4

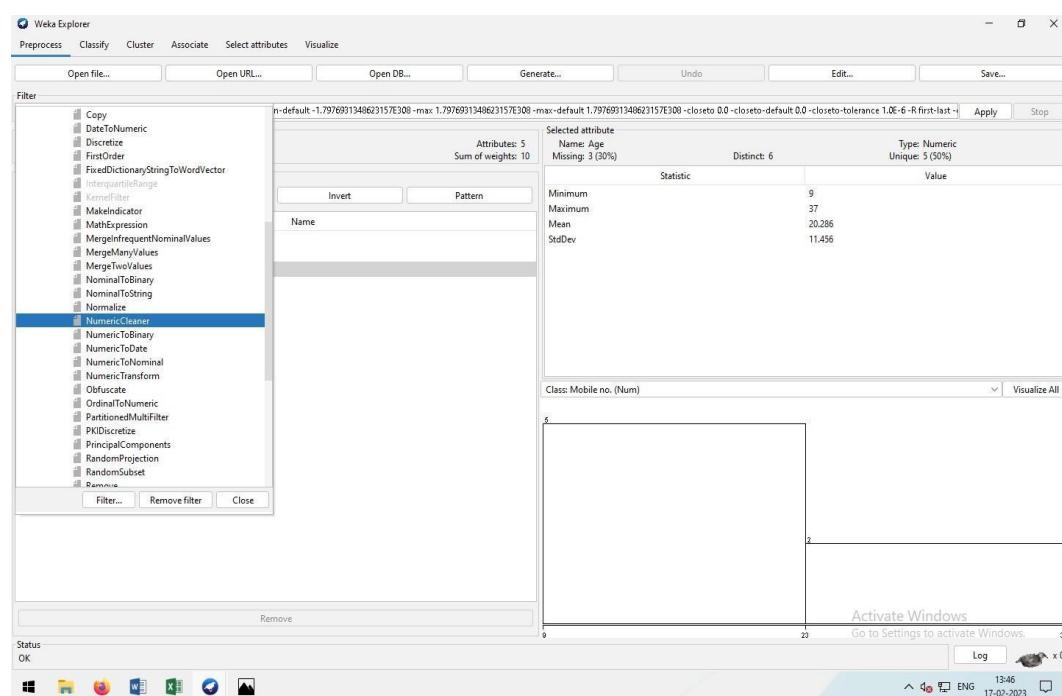
**Aim:** Pre-process a dataset based on handling missing values.

**Code:**

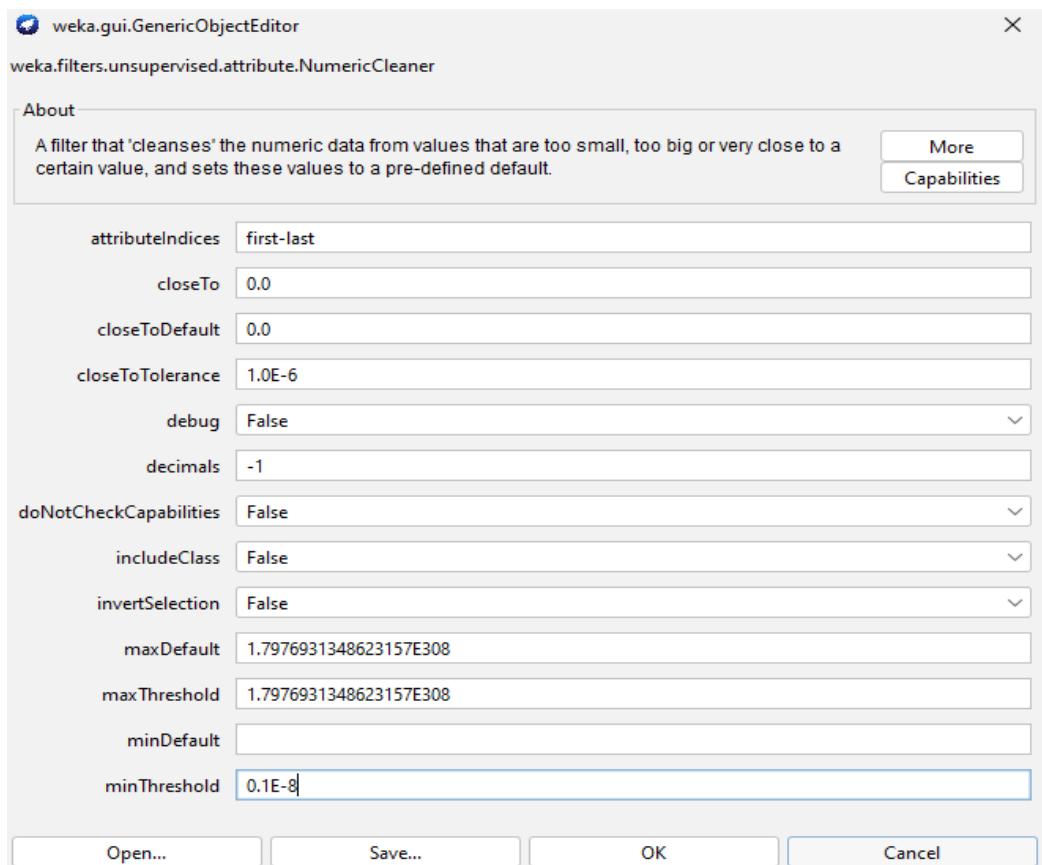
**STEP:1-** Load the weka onset of Book1 dataset.



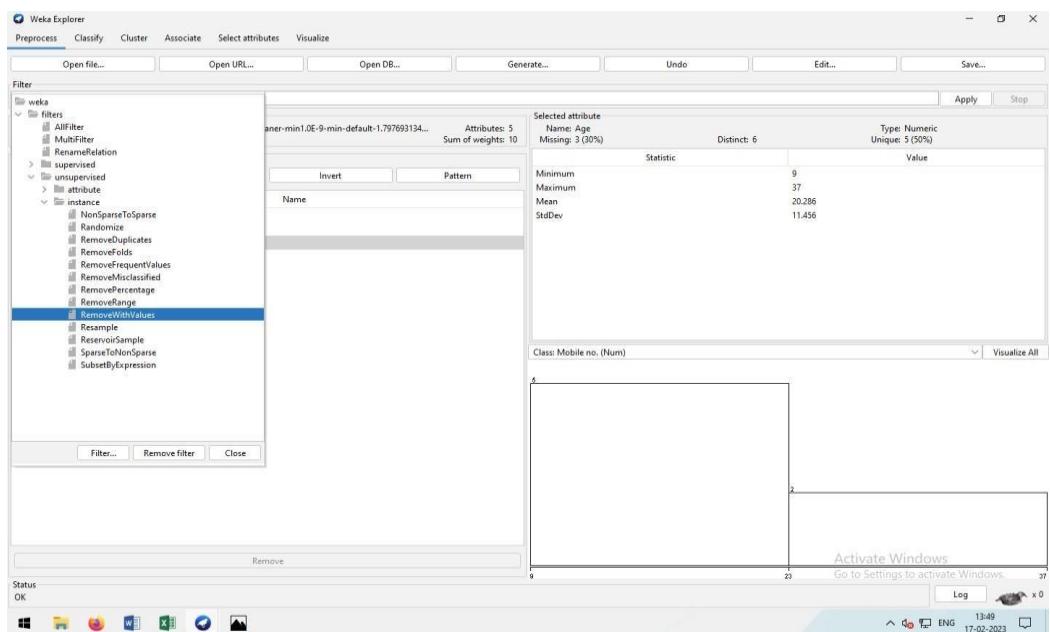
**STEP:2-** Click the “Choose” button for the Filter and select NumericalCleaner, it us under unsupervised.attribute.NumericalCleaner.



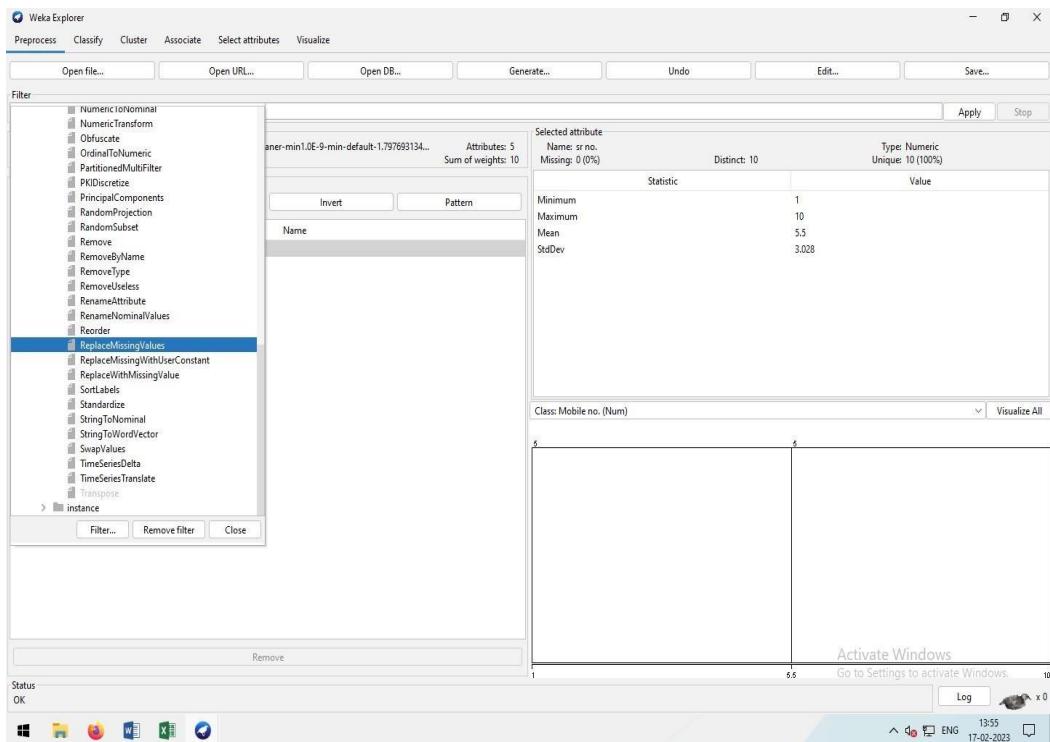
**STEP:3-** Set the *attributeIndices* to 6, the index of the mass attribute. Set *minThreshold* to 0.1E-8 (close to zero), which is the minimum value allowed for the attribute.



**STEP:4-** Click the “Choose” button for the Filter and select RemoveWithValues, it us under unsupervised.instance.RemoveWithValues.



**STEP:5-** Click the “Choose” button for the Filter and select ReplaceMissingValues, it us under unsupervised.attribute.ReplaceMissingValues.



**STEP:6-** Click the “Apply” button to apply the filter.

No.	1: sr no.	2: Name	3: Age	4: Gender	5: Mobile no.
	Numeric	Nominal	Numeric	Nominal	Numeric
1	1.0	abc	14.0	m	1.23456789E9
2	2.0	def	20.285...	f	1.47258369E9
3	3.0	ghi	36.0	f	1.59374862E9
4	4.0	jkl	14.0	m	3.57912486E9
5	5.0	mno	20.285...	m	1.023654789E9
6	6.0	pqr	37.0	f	3.65214789E9
7	7.0	stu	19.0	m	1.45236987E9
8	8.0	vwx	9.0	m	1.593640782E9
9	9.0	yz	20.285...	f	1.47586932E9
10	10.0	fd	13.0	f	2.103654789E9

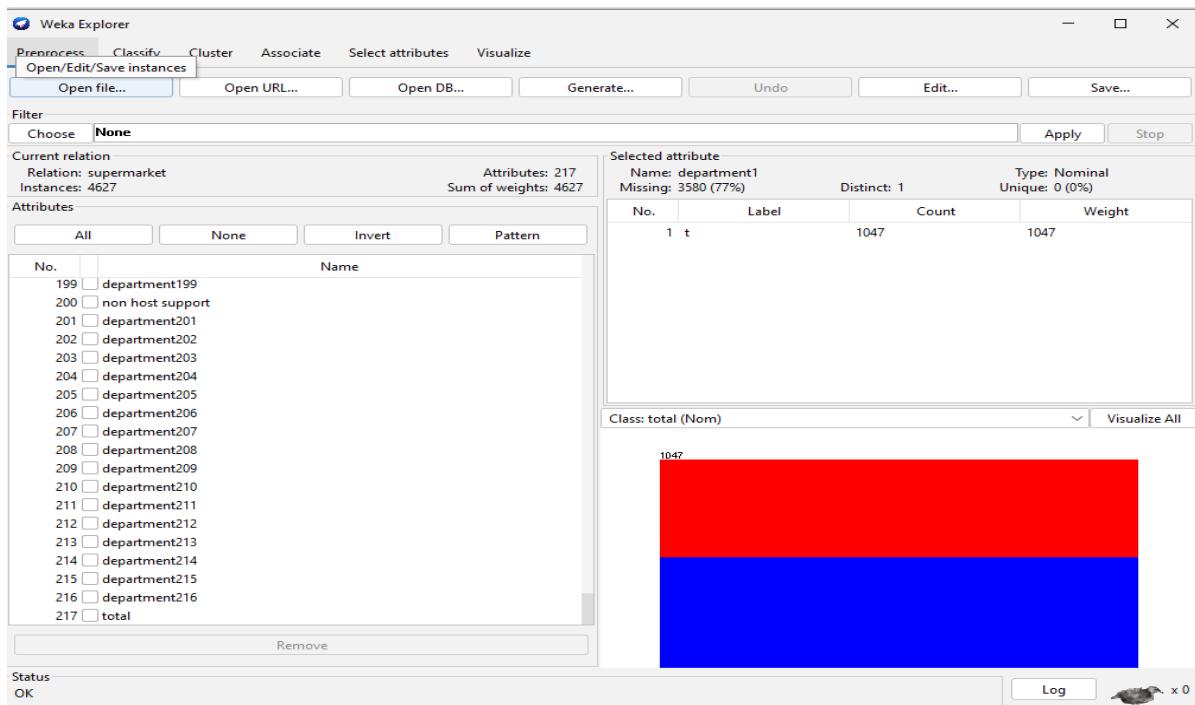
At the bottom of the table, there are buttons: 'Add instance', 'Undo', 'OK', and 'Cancel'.

## Practical: 5

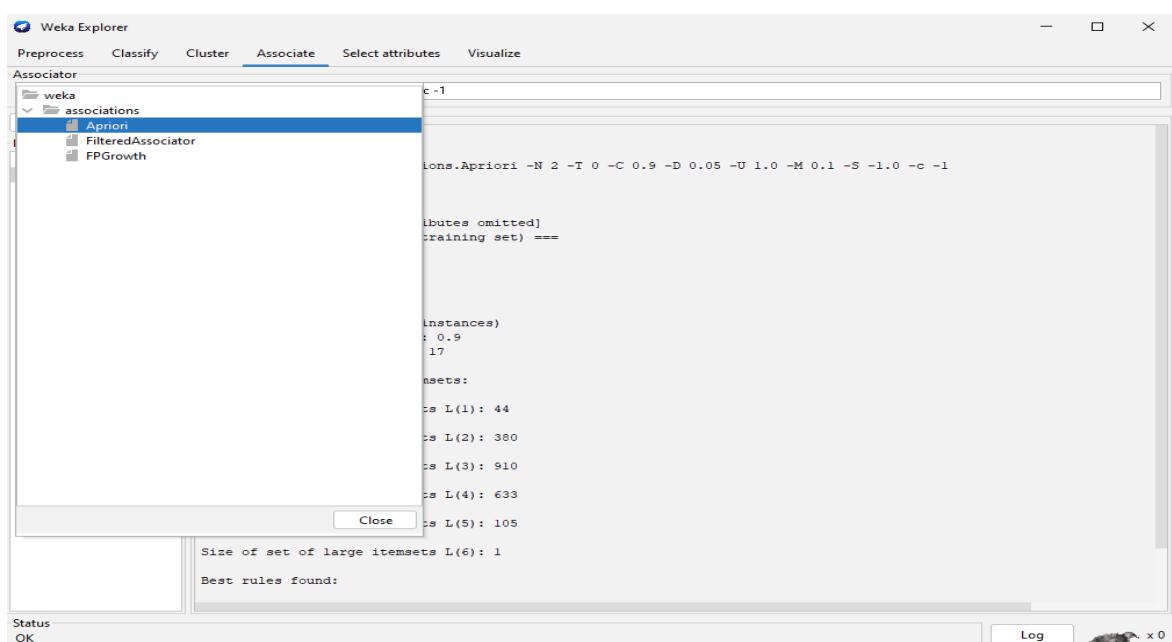
**Aim:** Generate association rules using the apriori algorithm.

**Code:**

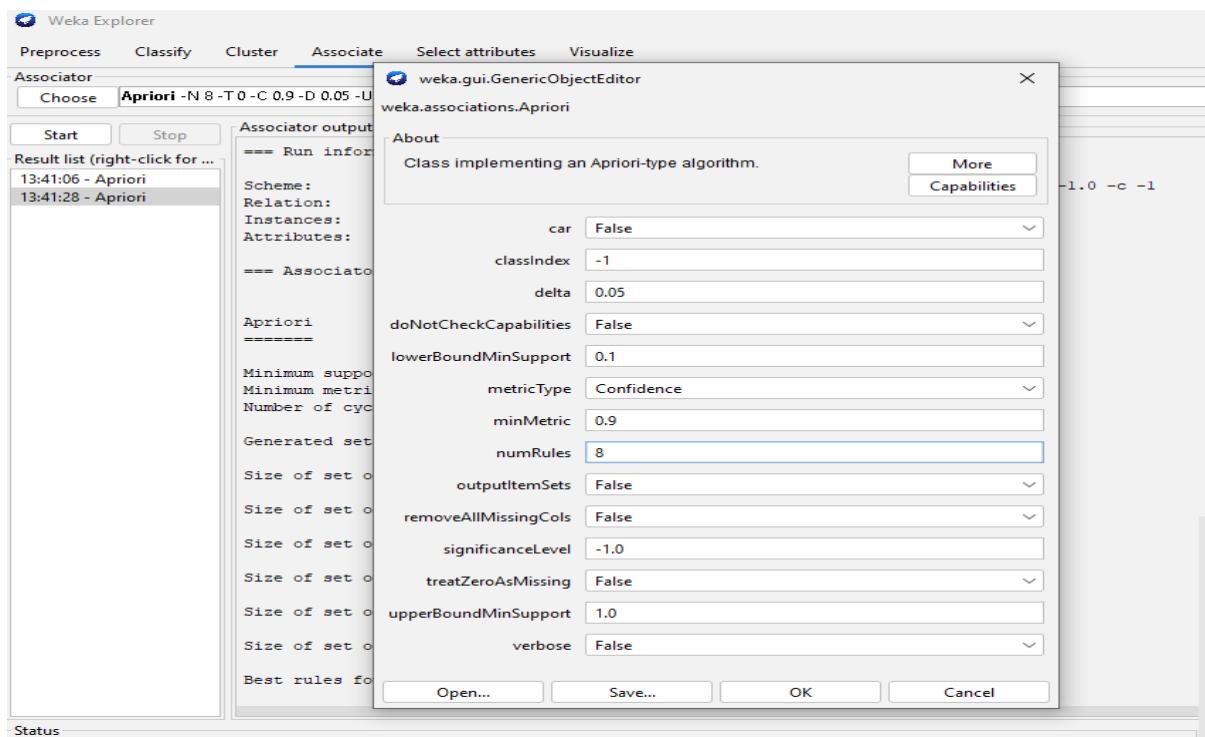
**Step-1:** In the WEKA open the Preprocess tab, click on the Open file button and select supermarket.arff dataset. After the data is loaded you will see the following screen



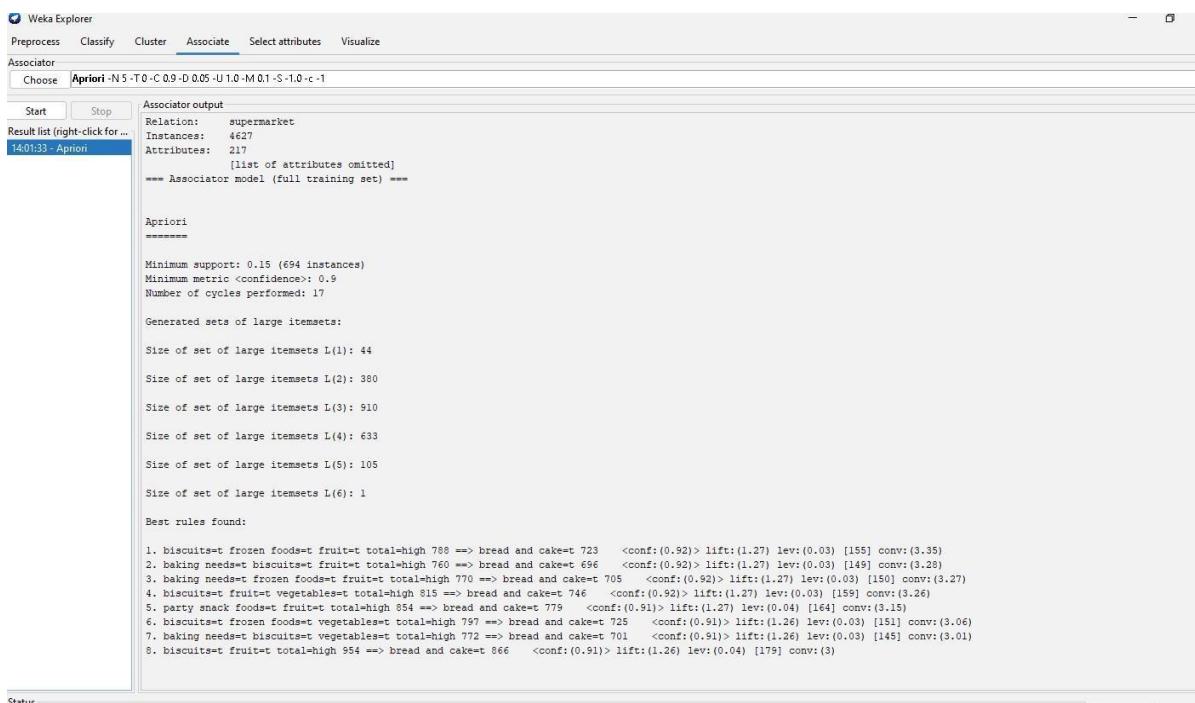
**Step-2:** Click on the Associate TAB and click on the Choose button. Select the Aprioriassociation as shown in below:



**Step-3:** set the parameters for the Apriori algorithm, click on its name, a window will pop up as shown below that allows you to set the parameters

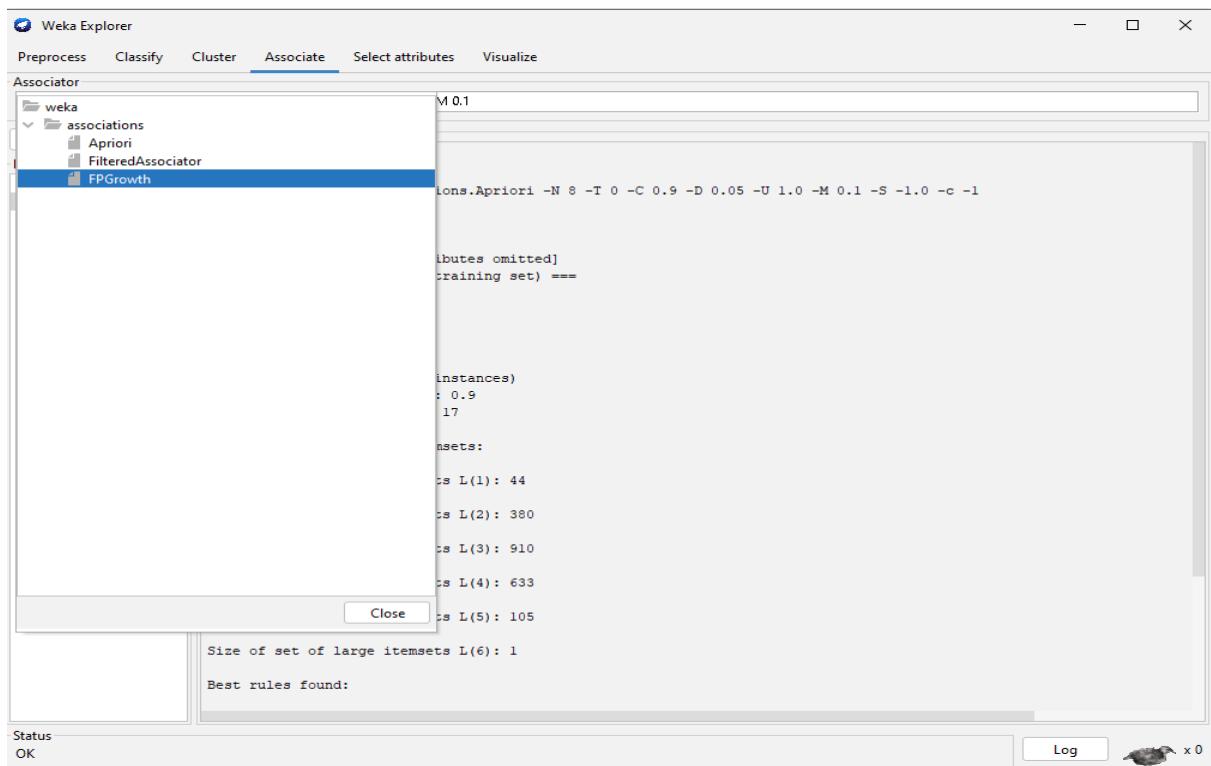


**Step-4:** set the parameters, click the Start button. After a while you will see the results as shown in the below.

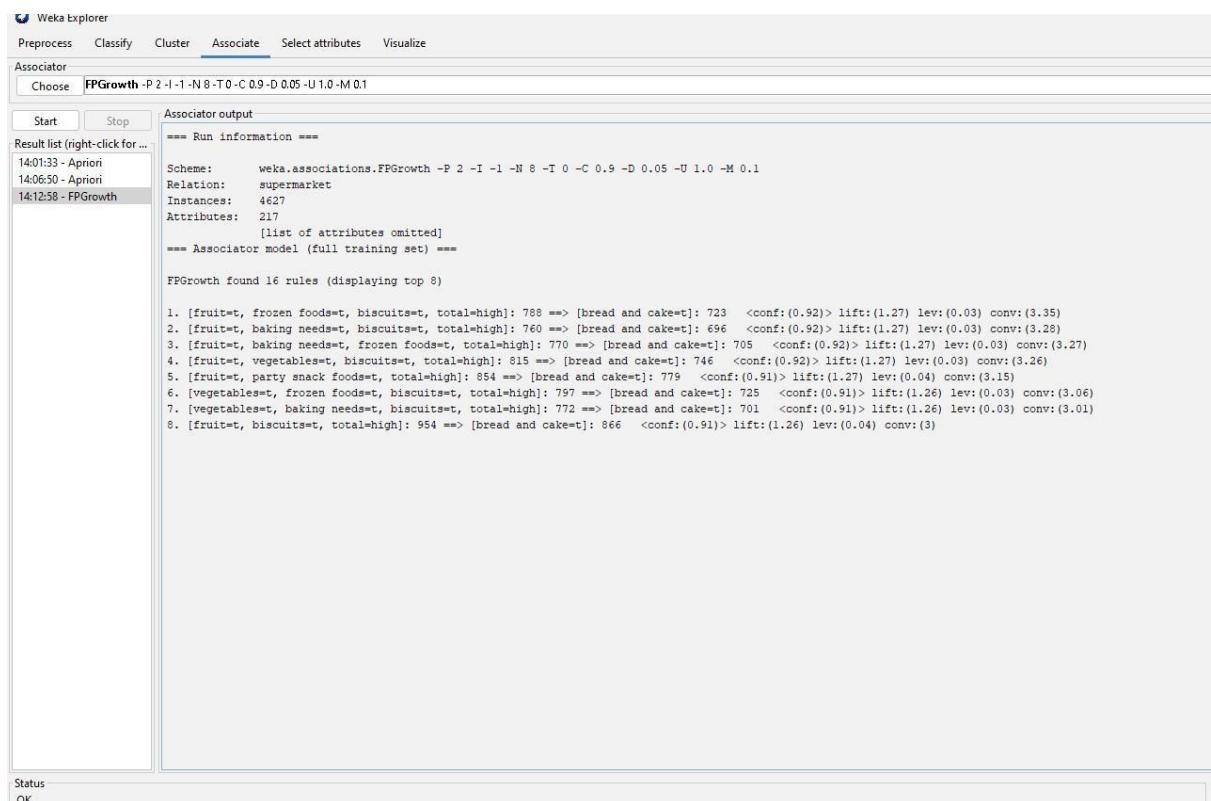


At the bottom, we will find the detected best rules of associations. This will help the supermarket in stocking their products in appropriate shelves.

**Step-5:** Click on the Associate TAB and click on the Choose button. Select the FP Growth association as shown in below:



**Step-6:** set the parameters, click the Start button. After a while you will see the results as shown in the below.



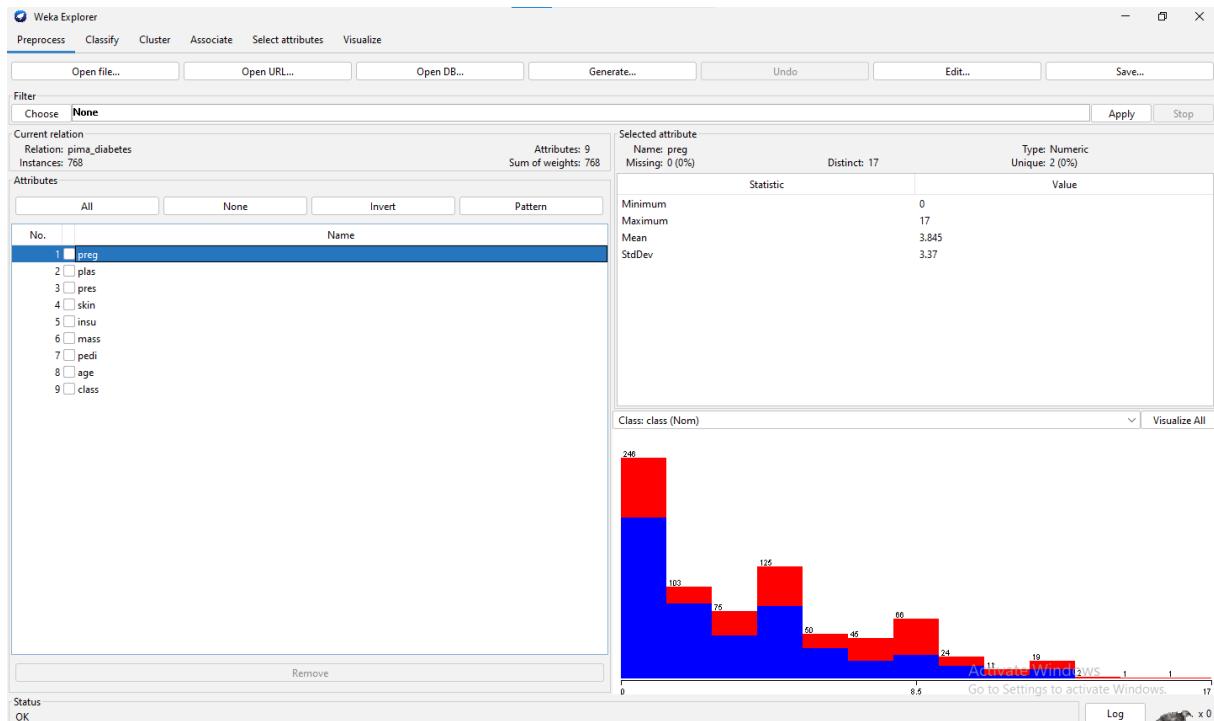
## Practical: 6

**Aim:** Implement classifier decision tree on any two dataset.

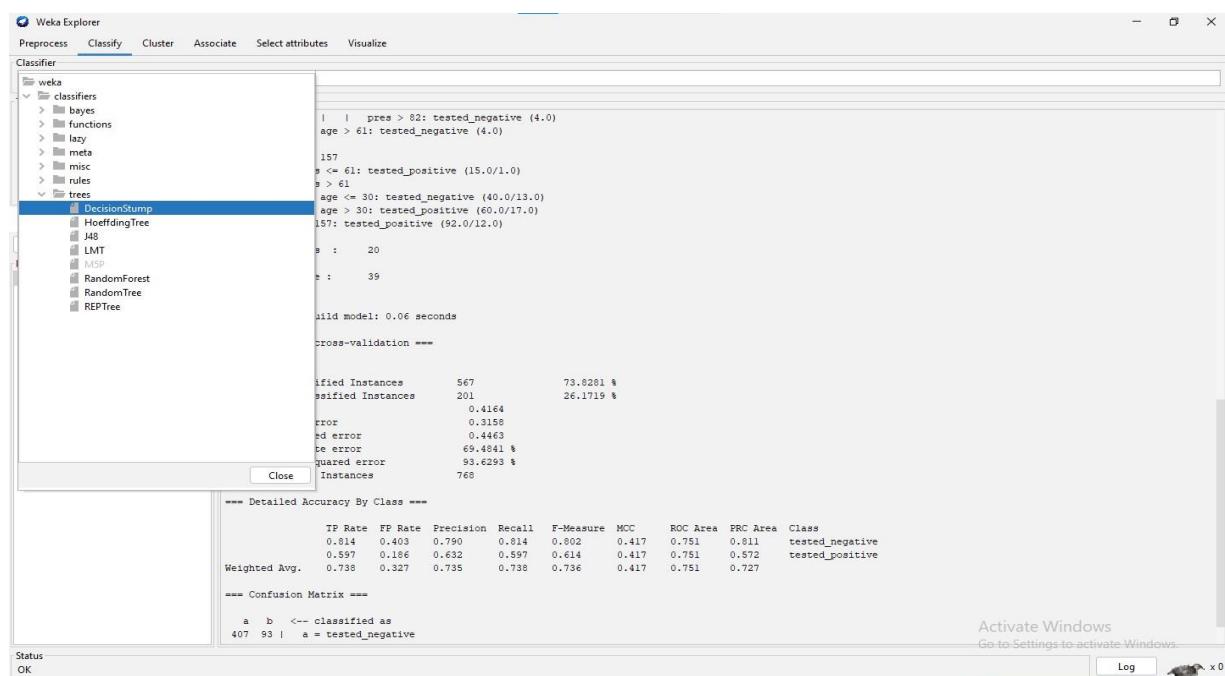
**Code:**

➤ Dataset 1

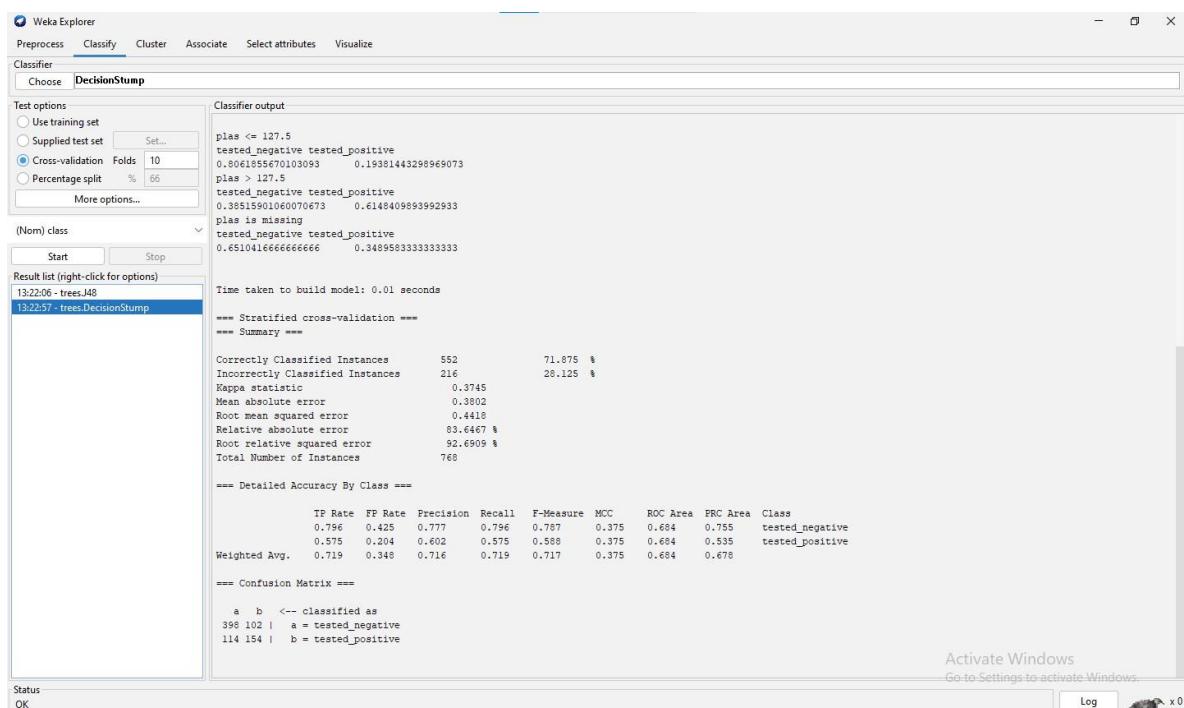
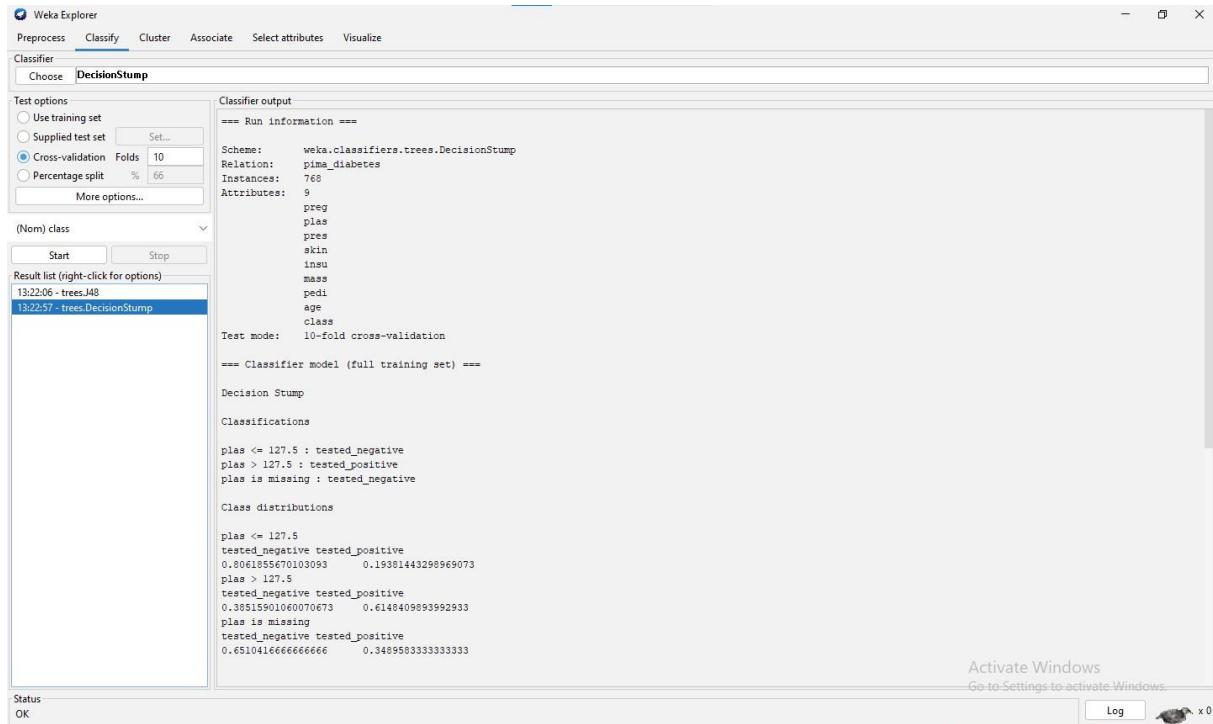
**Step 1)** Click on open file and load the data



**Step 2)** Click on Classify and choose the algorithm of Decision tree.



### Step 3) DecisionStump Algorithm



## Step 4) tree.J48 Algorithm

**Weka Explorer**

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose DecisionStump

Test options

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

(Nom) class Start Stop

Result list (right-click for options)

13:22:06 - trees.J48  
13:22:57 - trees.DecisionStump

```
==== Run information ====
Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2
Relation: pima_diabetes
Instances: 768
Attributes: 9
preg
plas
pres
skin
insu
mass
pedi
age
class

Test mode: 10-fold cross-validation

==== Classifier model (full training set) ====
J48 pruned tree
-----
plas <= 127
| mass <= 26.4: tested_negative (132.0/3.0)
| mass > 26.4
| | age <= 28: tested_negative (180.0/22.0)
| | age > 28
| | | plas <= 99: tested_negative (55.0/10.0)
| | | plas > 99
| | | | pedi <= 0.56: tested_negative (84.0/34.0)
| | | | pedi > 0.56
| | | | | preg <= 6
| | | | | | age <= 30: tested_positive (4.0)
| | | | | | age > 30
| | | | | | | age <= 34: tested_negative (7.0/1.0)
| | | | | | | age > 34
| | | | | | | | mass <= 33.1: tested_positive (6.0)
| | | | | | | | mass > 33.1: tested_negative (4.0/1.0)
| | | | | | | | preg > 6: tested_positive (13.0)

plas > 127
```

Activate Windows  
Go to Settings to activate Windows.

Status OK Log x 0

**Weka Explorer**

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose DecisionStump

Test options

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

(Nom) class Start Stop

Result list (right-click for options)

13:22:06 - trees.J48  
13:22:57 - trees.DecisionStump

```
plas > 127
| mass <= 29.9
| | plas <= 145: tested_negative (41.0/6.0)
| | plas > 145
| | | age <= 25: tested_negative (4.0)
| | | age > 25
| | | | age <= 61
| | | | | mass <= 27.1: tested_positive (12.0/1.0)
| | | | | mass > 27.1
| | | | | | pres <= 82
| | | | | | | pedi <= 0.396: tested_positive (8.0/1.0)
| | | | | | | pedi > 0.396: tested_negative (3.0)
| | | | | | | pres > 82: tested_negative (4.0)
| | | | | | | age > 61: tested_negative (4.0)
| mass > 29.9
| | plas <= 157
| | | pres <= 61: tested_positive (15.0/1.0)
| | | pres > 61
| | | | age <= 30: tested_negative (40.0/13.0)
| | | | age > 30: tested_positive (60.0/17.0)
| | plas > 157: tested_positive (92.0/12.0)

Number of Leaves : 20
Size of the tree : 39

Time taken to build model: 0.06 seconds

==== Stratified cross-validation ====
==== Summary ====
Correctly Classified Instances 567 73.8281 %
Incorrectly Classified Instances 201 26.1719 %
Kappa statistic 0.4164
Mean absolute error 0.3158
Root mean squared error 0.4463
Relative absolute error 69.4841 %
Root relative squared error 93.6293 %
Total Number of Instances 768
```

Activate Windows  
Go to Settings to activate Windows.

Status OK Log x 0

**Weka Explorer**

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose DecisionStump

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)

13:22:06 - trees.J48  
13:22:57 - trees.DecisionStump

Classifier output

```

|   |   plas <= 157
|   |   |   pres <= 61: tested_positive (15.0/1.0)
|   |   |   pres > 61
|   |   |   |   age <= 30: tested_negative (40.0/13.0)
|   |   |   |   age > 30: tested_positive (60.0/17.0)
|   |   plas > 157: tested_positive (92.0/12.0)

```

Number of Leaves : 20

Size of the tree : 39

Time taken to build model: 0.06 seconds

==== Stratified cross-validation ====  
==== Summary ====  
==== Detailed Accuracy By Class ====  
==== Confusion Matrix ====  
a b <- classified as  
407 93 | a = tested\_negative  
108 160 | b = tested\_positive

Status OK Log x 0

## Step 5) tree.RandomForest Algorithm

**Weka Explorer**

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier Choose RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1

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)

13:22:06 - trees.J48  
13:22:57 - trees.DecisionStump  
13:28:29 - trees.RandomForest

Classifier output

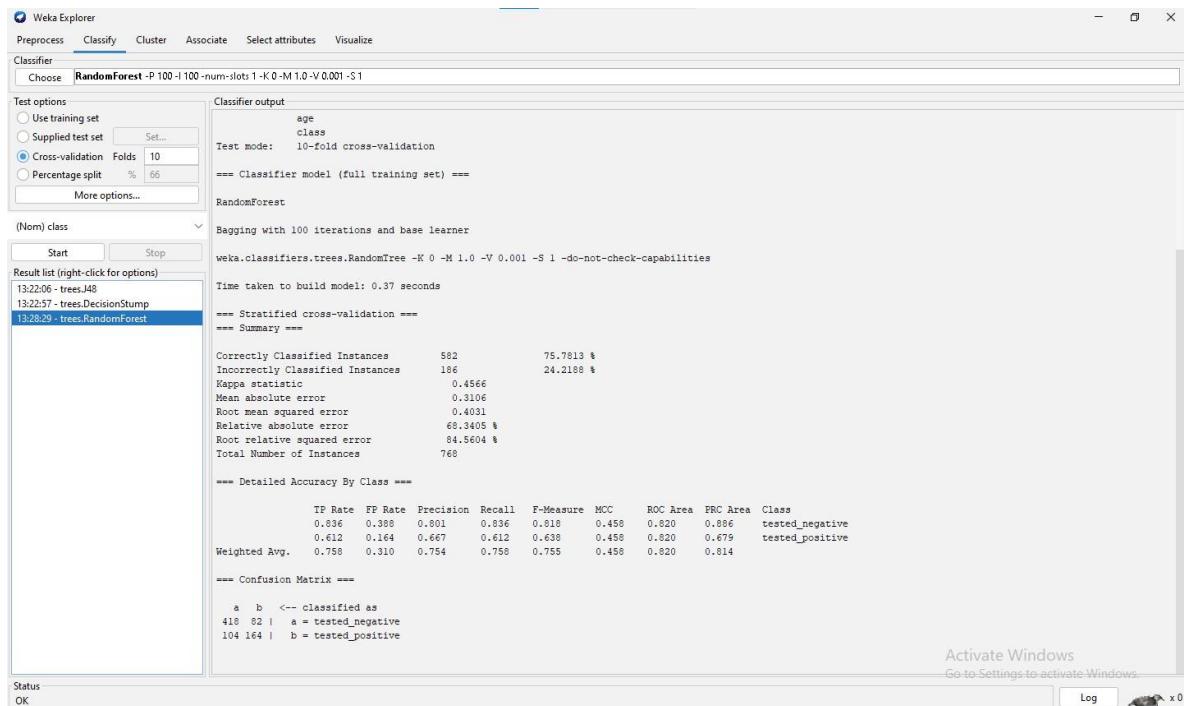
```

Scheme: weka.classifiers.trees.RandomForest -P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1
Relation: pima_diabetes
Instances: 768
Attributes: 9
preg
plas
pres
skin
insu
mass
pedi
age
class
Test mode: 10-fold cross-validation

```

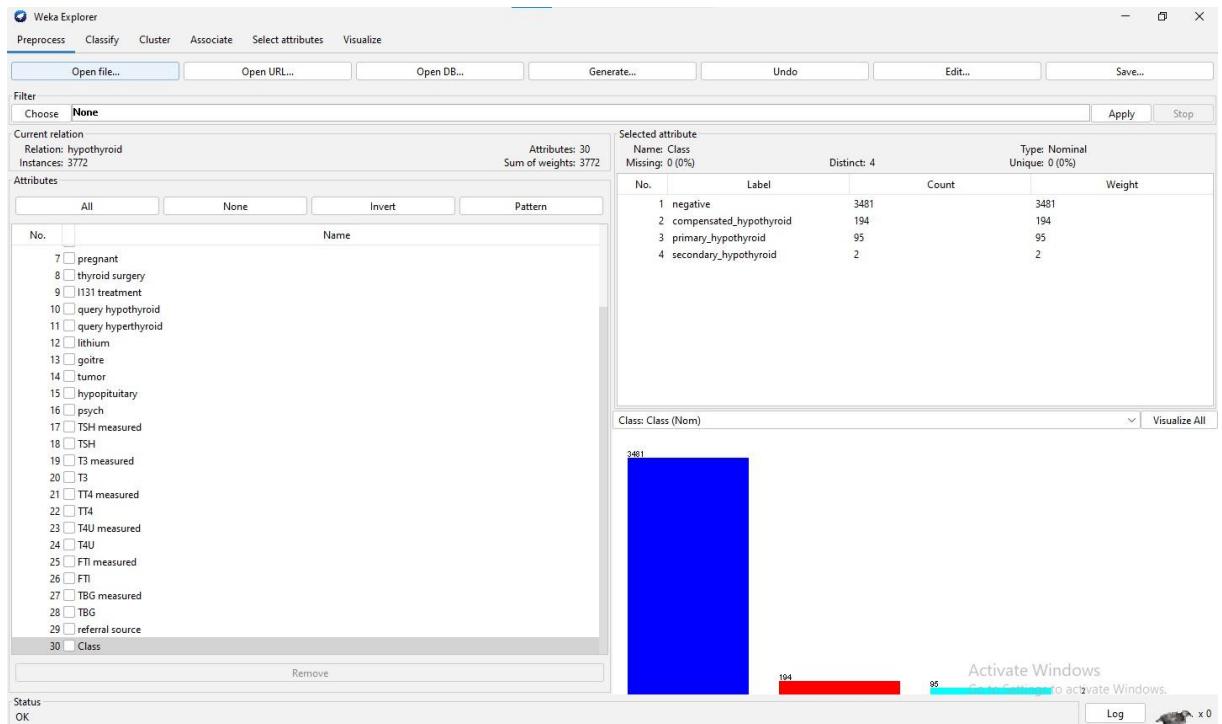
==== Classifier model (full training set) ====  
RandomForest  
Bagging with 100 iterations and base learner  
weka.classifiers.trees.RandomTree -K 0 -M 1.0 -V 0.001 -S 1 -do-not-check-capabilities  
Time taken to build model: 0.37 seconds  
==== Stratified cross-validation ====  
==== Summary ====  
==== Detailed Accuracy By Class ====  
Activate Windows  
Go to Settings to activate Windows.

Status OK Log x 0

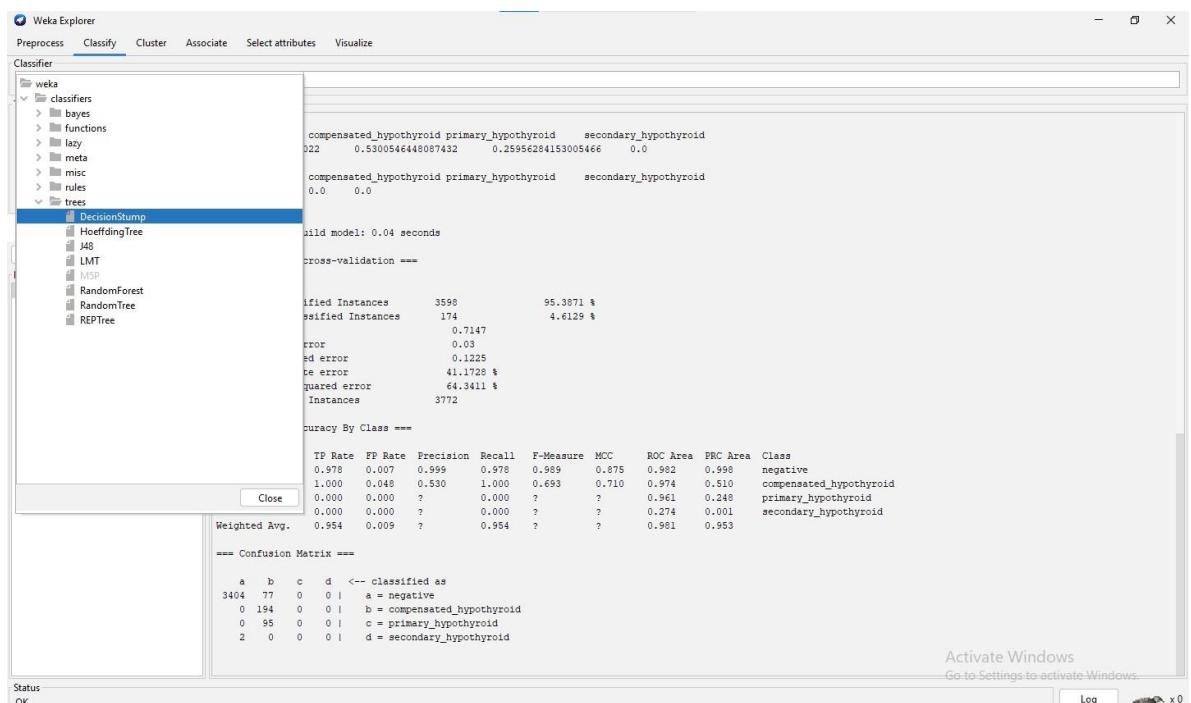


## ➤ Dataset 2

**Step 1)** Click on open file and load the data



**Step 2)** Click on Classify and choose the algorithm of Decision tree.



### Step 3) DecisionStump Algorithm

**Classifier output**

```

TSB = 6.05
negative compensated_hypothyroid primary_hypothyroid secondary_hypothyroid
0.210325136612022 0.5300546448007432 0.25956284153005466 0.0
TSB is missing
negative compensated_hypothyroid primary_hypothyroid secondary_hypothyroid
1.0 0.0 0.0 0.0

```

Time taken to build model: 0.04 seconds

==== Stratified cross-validation ====  
==== Summary ====  
Correctly Classified Instances 3586 95.3871 %  
Incorrectly Classified Instances 174 4.6129 %  
Kappa statistic 0.7147  
Mean absolute error 0.03  
Root mean squared error 0.1225  
Relative absolute error 41.1728 %  
Root relative squared error 64.3411 %  
Total Number of Instances 3772

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

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.978	0.007	0.999	0.978	0.989	0.975	0.982	0.998	0.998	negative
1.000	0.048	0.530	1.000	0.693	0.710	0.974	0.510	0.510	compensated_hypothyroid
0.000	0.000	?	0.000	?	?	0.961	0.248	0.248	primary_hypothyroid
0.000	0.000	?	0.000	?	?	0.274	0.001	0.001	secondary_hypothyroid
Weighted Avg.	0.954	0.009	?	0.954	?	?	0.981	0.953	

  
==== Confusion Matrix ====  

a	b	c	d	--- classified as
3404	77	0	0	a = negative
0	194	0	0	b = compensated_hypothyroid
0	95	0	0	c = primary_hypothyroid
2	0	0	0	d = secondary_hypothyroid

### Step 4) tree.J48 Algorithm

**Classifier output**

```

| | | TSB measured = f: negative (30.75)
| | | on thyroxine = t: negative (56.17)
Number of Leaves : 15
Size of the tree : 29

```

Time taken to build model: 0.07 seconds

==== Stratified cross-validation ====  
==== Summary ====  
Correctly Classified Instances 3756 99.5758 %  
Incorrectly Classified Instances 16 0.4242 %  
Kappa statistic 0.9707  
Mean absolute error 0.003  
Root mean squared error 0.0414  
Relative absolute error 4.1612 %  
Root relative squared error 21.7445 %  
Total Number of Instances 3772

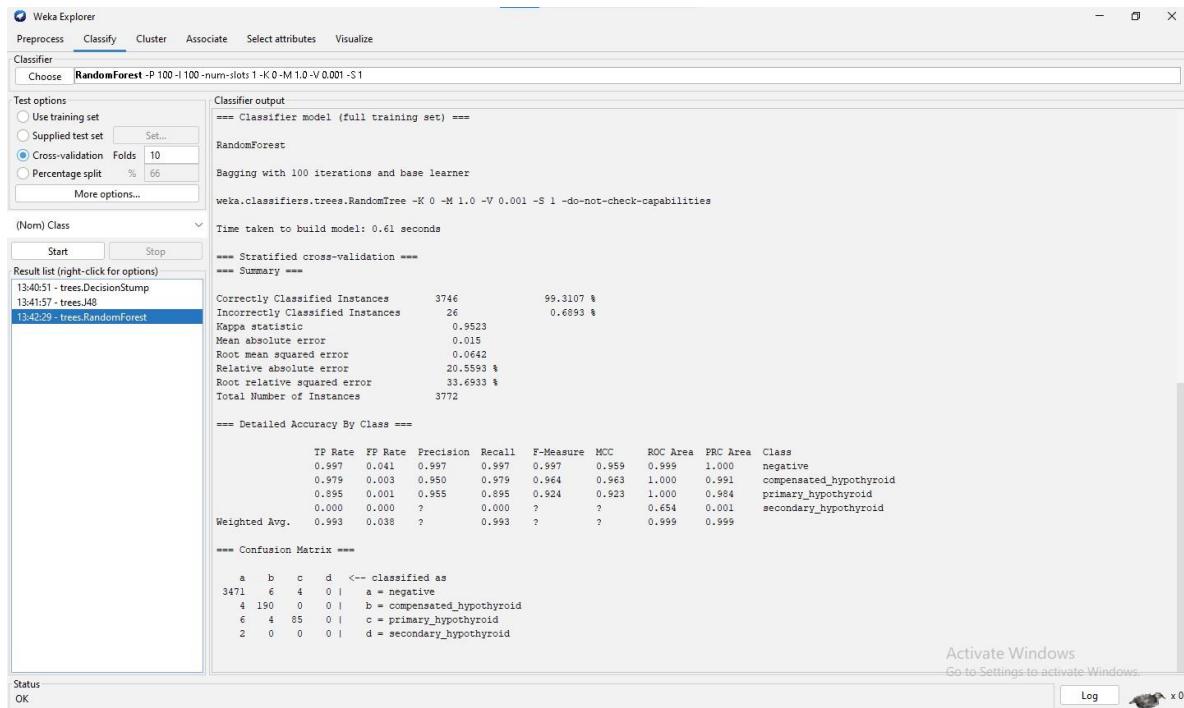
==== Detailed Accuracy By Class ====  

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.999	0.021	0.998	0.999	0.998	0.979	0.993	0.999	0.999	negative
0.985	0.002	0.970	0.985	0.977	0.976	0.999	0.964	0.964	compensated_hypothyroid
0.937	0.001	0.957	0.937	0.947	0.946	1.000	0.988	0.988	primary_hypothyroid
0.000	0.000	?	0.000	?	?	0.197	0.000	0.000	secondary_hypothyroid
Weighted Avg.	0.996	0.019	?	0.996	?	?	0.993	0.996	

  
==== Confusion Matrix ====  

a	b	c	d	--- classified as
3476	3	2	0	a = negative
1	191	2	0	b = compensated_hypothyroid
3	3	89	0	c = primary_hypothyroid
2	0	0	0	d = secondary_hypothyroid

## Step 5) tree.RandomForest Algorithm



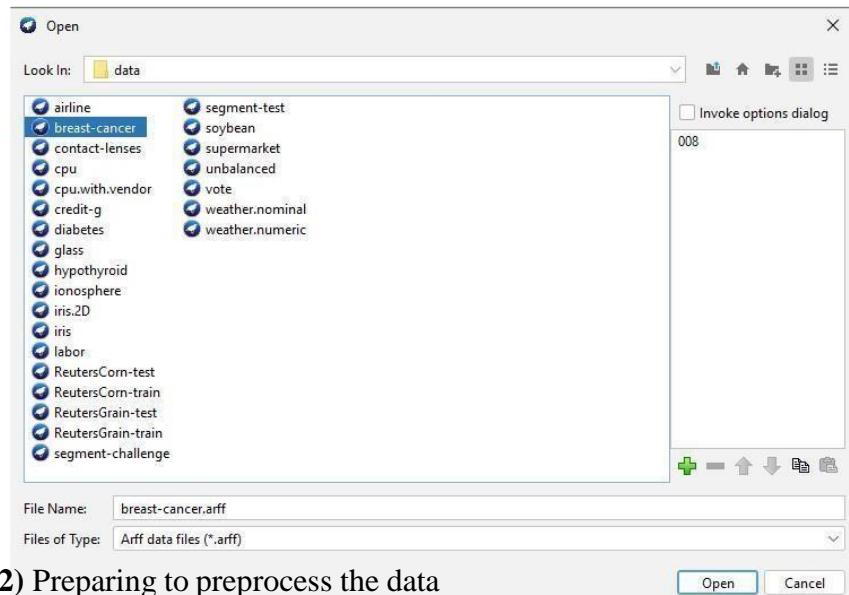
## Practical: 7

**Aim:** Implement classification algorithm (Naïve Bayes) on any two datasets.

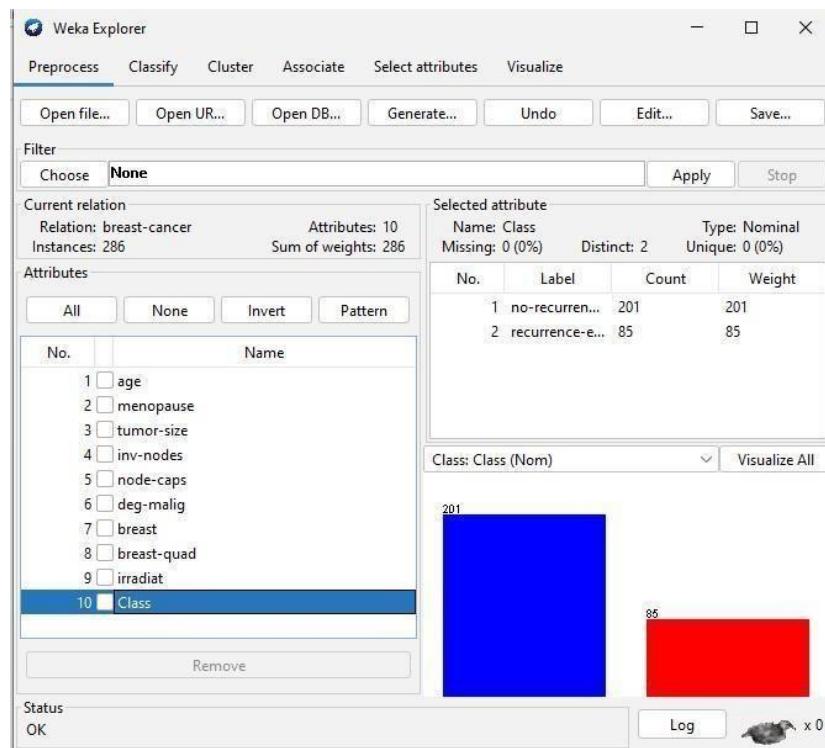
**Code:**

- Dataset: Breast Cancer

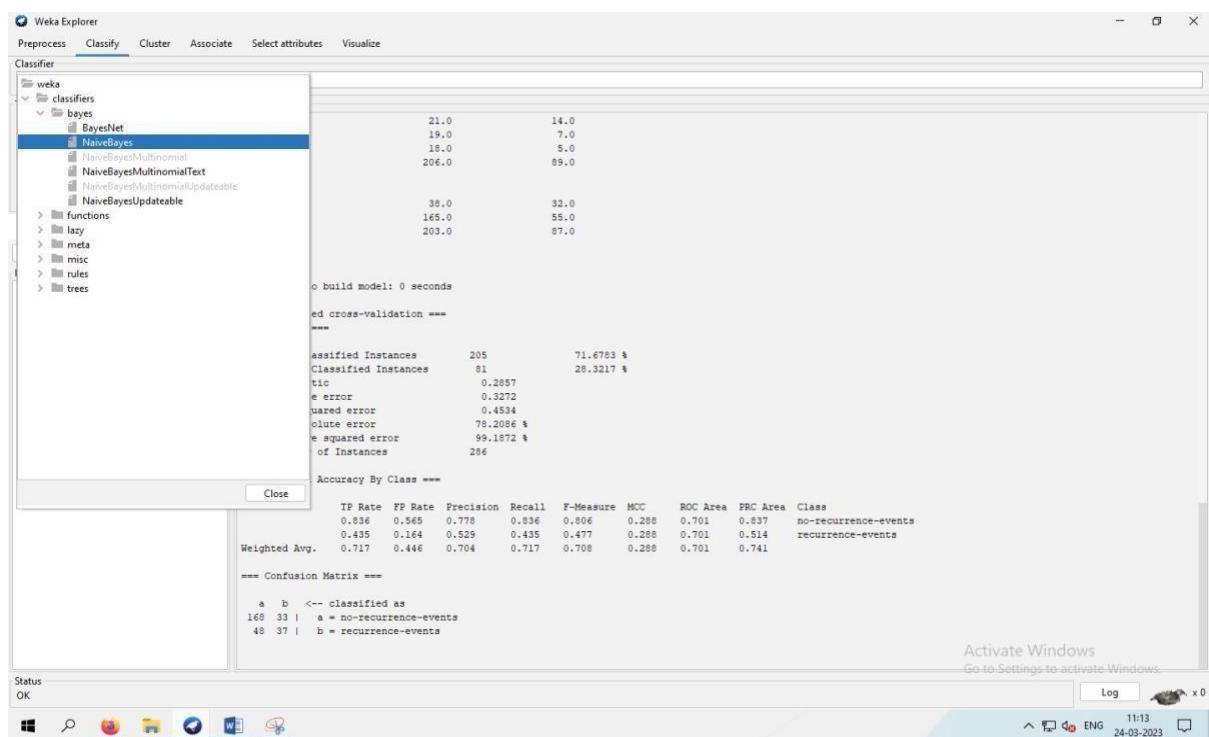
**Step 1)** Select the dataset



**Step 2)** Preparing to preprocess the data



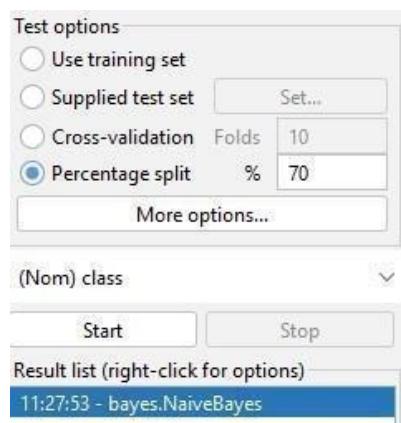
**Step 3)** Now go to the classify tab. Choose the Naïve Bayes classifier. Be sure that the Playattribute is selected as a class selector, and then press the Start button to build a model.

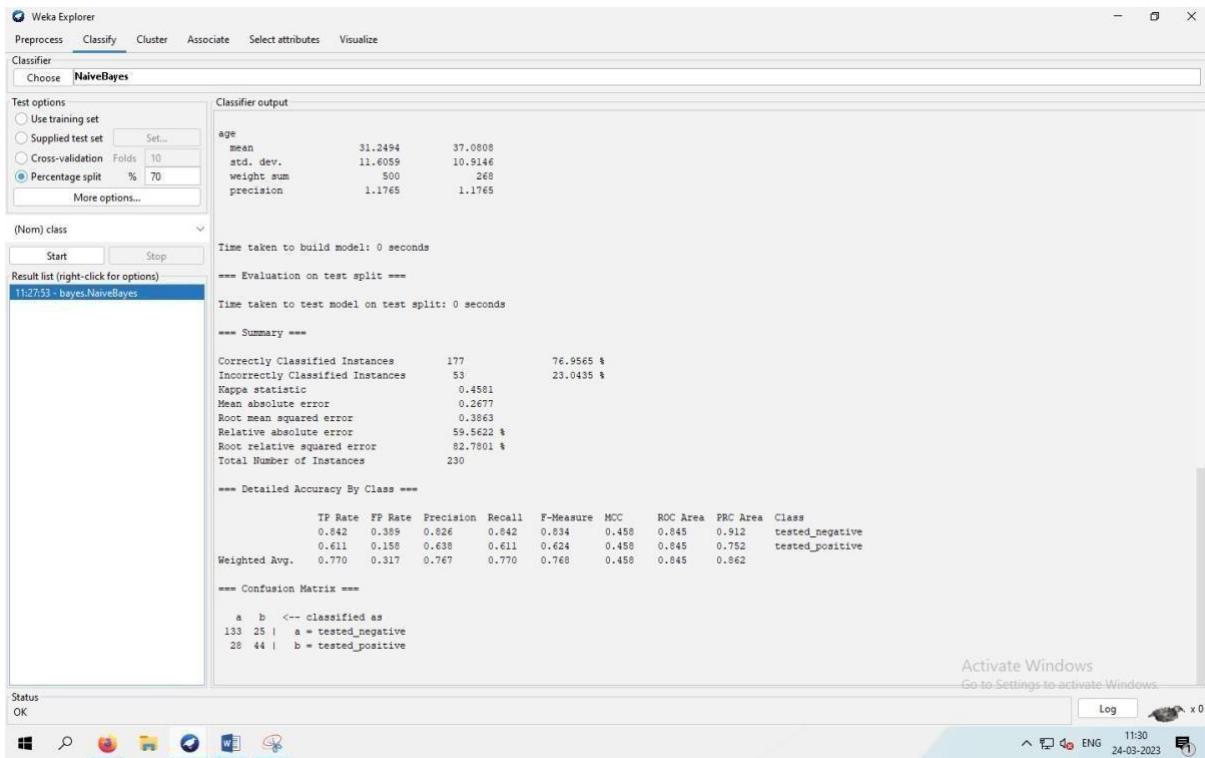


**Step 4)** Model outputs some information on how accurate it classifies and otherparameters.

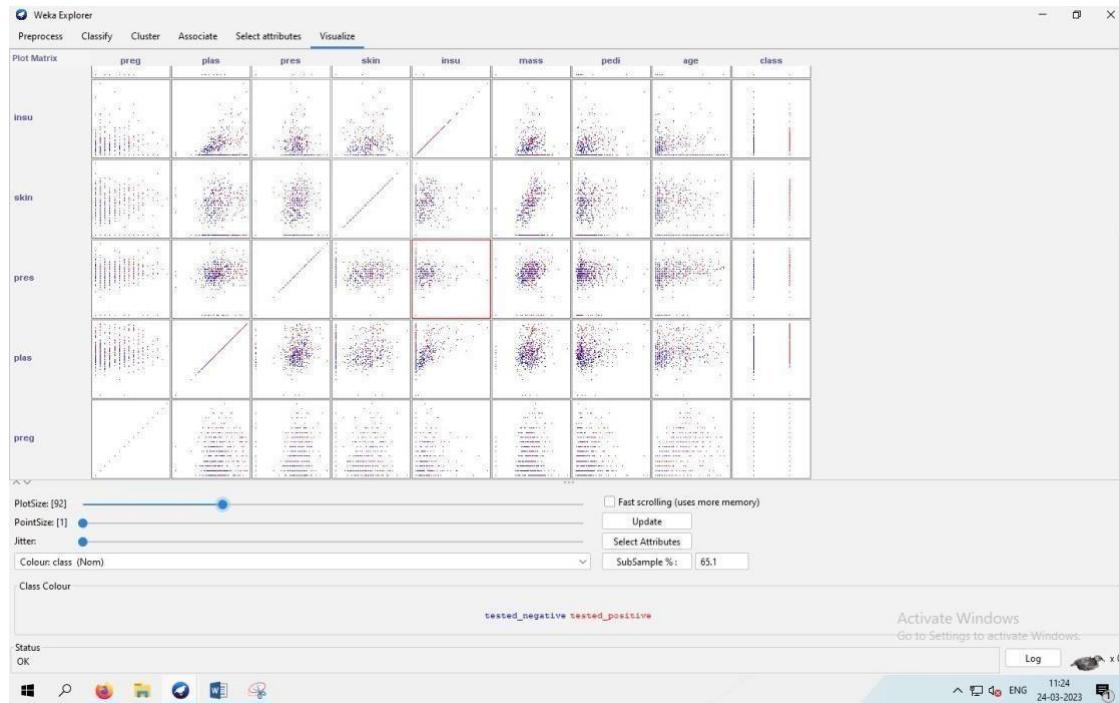
Correctly Classified Instances	205	71.6783 %
Incorrectly Classified Instances	81	28.3217 %

**Step 5)** Testing option to percentage split



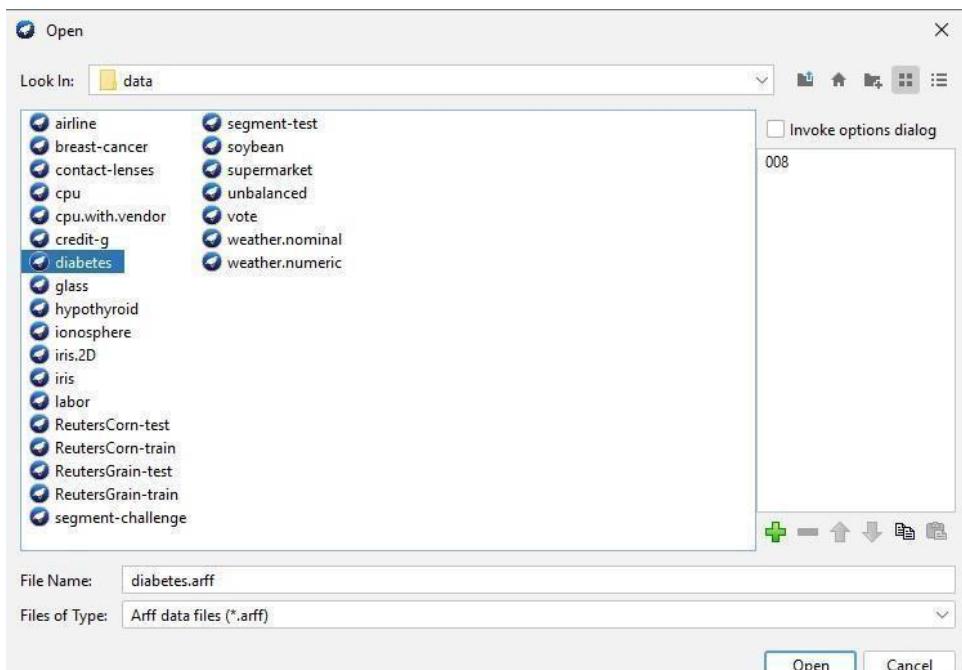


## Step 6) Visualize

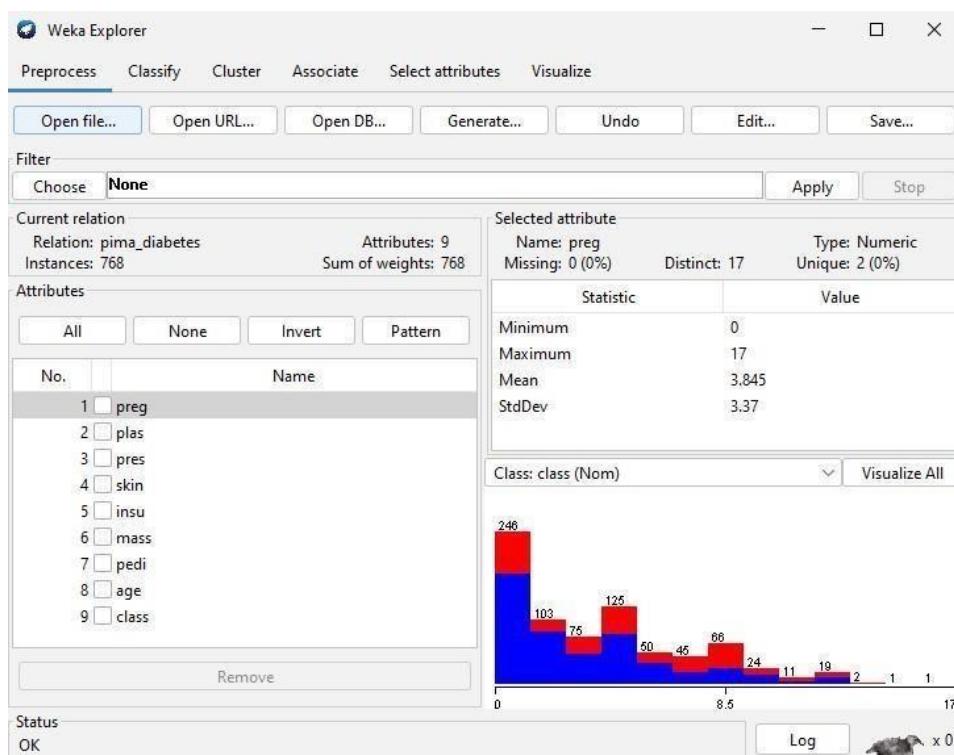


➤ Dataset: Diabetes

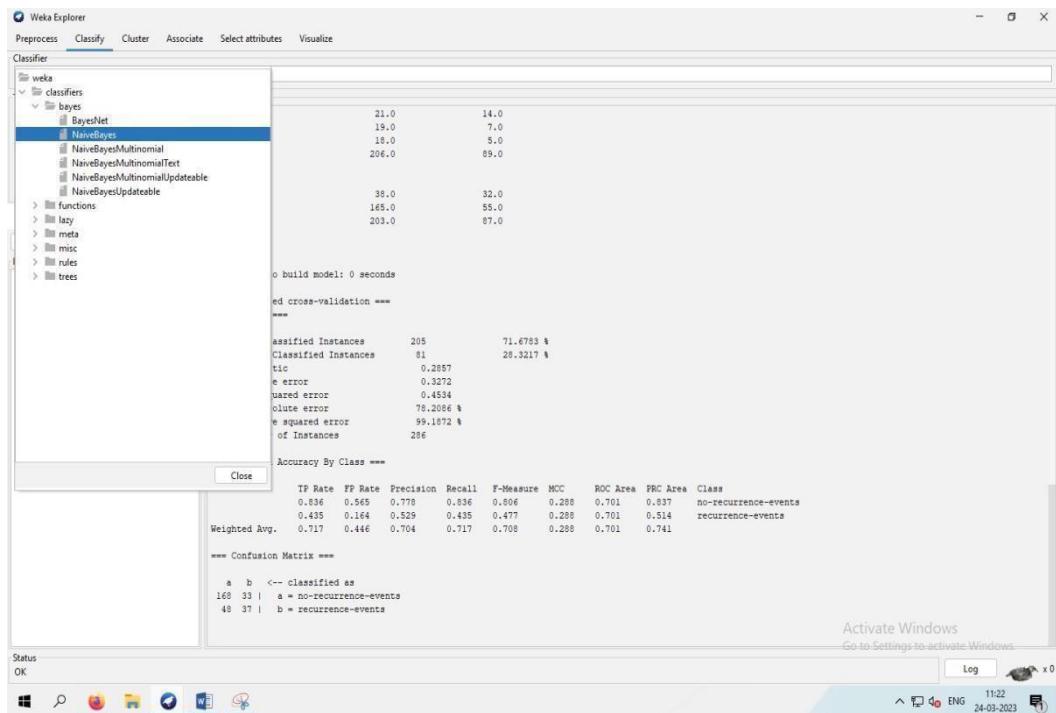
**Step 1)** Select the dataset



**Step 2)** Preparing to preprocess the data



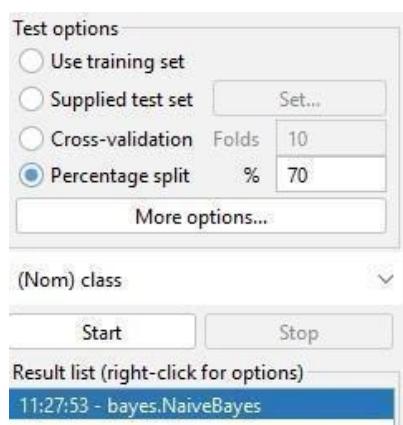
**Step 3)** Now go to the classify tab. Choose the Naïve Bayes classifier. Be sure that the Play attribute is selected as a class selector, and then press the Start button to build a model.



**Step 4)** Model outputs some information on how accurate it classifies and other parameters.

Correctly Classified Instances	205	71.6783 %
Incorrectly Classified Instances	81	28.3217 %

**Step 5)** Testing option to percentage split



**Weka Explorer**

Preprocess Classify Cluster Associate Select attributes Visualize

Classifier  
Choose **NaiveBayes**

Test options  
 Use training set  
 Supplied test set Set...  
 Cross-validation Folds 10  
 Percentage split % 70  
[More options...](#)

(Nom) class

Result list (right-click for options)  
 11:27:53 - bayes.NaiveBayes  
**11:31:23 - bayes.NaiveBayes**

Classifier output

```

age
mean           31.2494      37.0088
std. dev.       11.6055      10.9146
weight sum       500          268
precision        1.1765      1.1765

Time taken to build model: 0 seconds
*** Evaluation on test split ***
Time taken to test model on test split: 0 seconds
*** Summary ***
Correctly Classified Instances   177      76.5565 %
Incorrectly Classified Instances 53       23.0435 %
Kappa statistic                 0.4891
Mean absolute error             0.3077
Root mean squared error         0.3863
Relative absolute error         59.5622 %
Root relative squared error    82.7001 %
Total Number of Instances      230

*** Detailed Accuracy By Class ***
               TP Rate FP Rate Precision Recall  F-Measure MCC ROC Area PRC Area Class
tested_negative 0.842   0.389   0.826   0.842   0.834   0.458   0.845   0.912   tested_negative
tested_positive 0.611   0.158   0.638   0.611   0.624   0.458   0.845   0.752   tested_positive
Weighted Avg.    0.770   0.317   0.767   0.770   0.768   0.458   0.845   0.862

*** Confusion Matrix ***
a   b   <-- Classified as
133  25 |  a = tested_negative
28  44 |  b = tested_positive

```

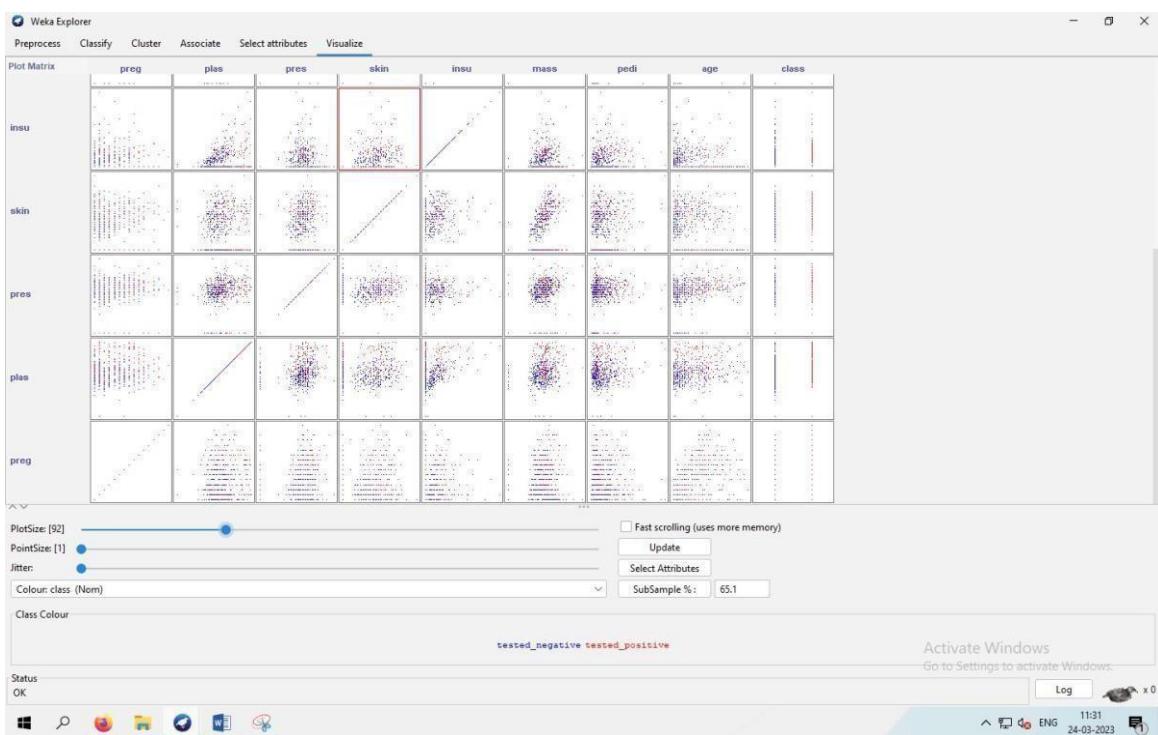
Status OK

Activate Windows  
Go to Settings to activate Windows.

Log x 0

11:31 ENG 24-03-2023

## Step 6) Visualize



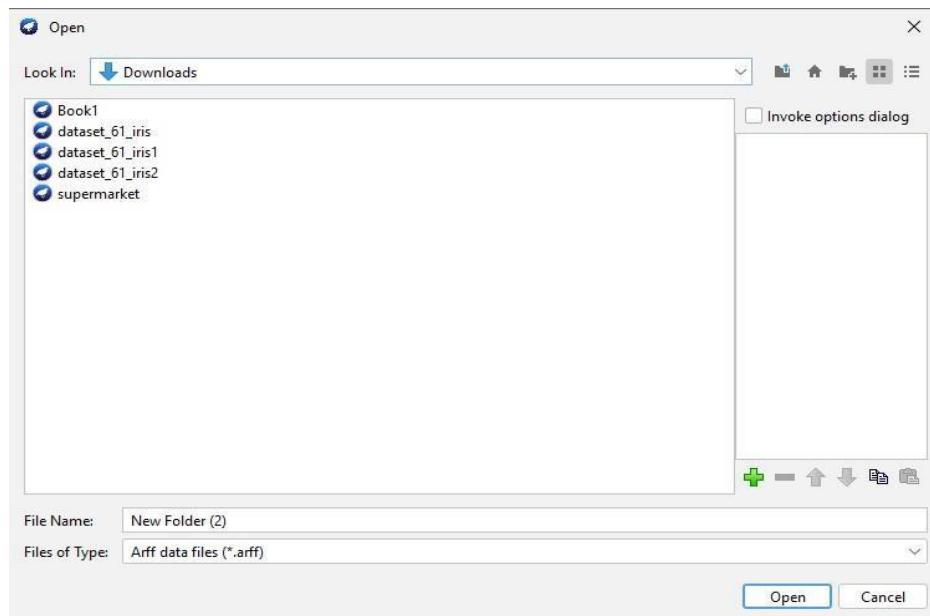
## Practical: 8

**Aim:** Applying clustering algorithm using K-means.

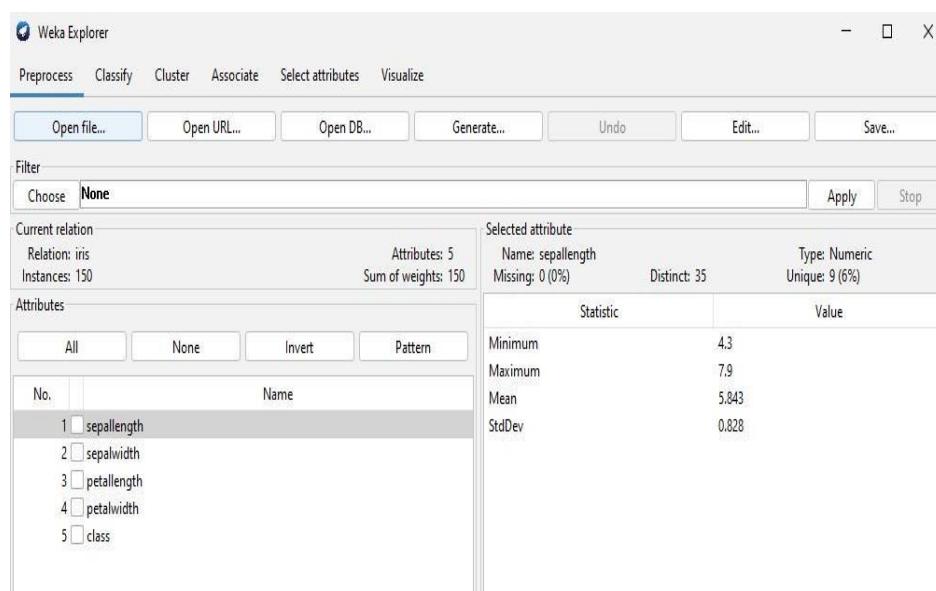
**Code:**

**Method: 1**

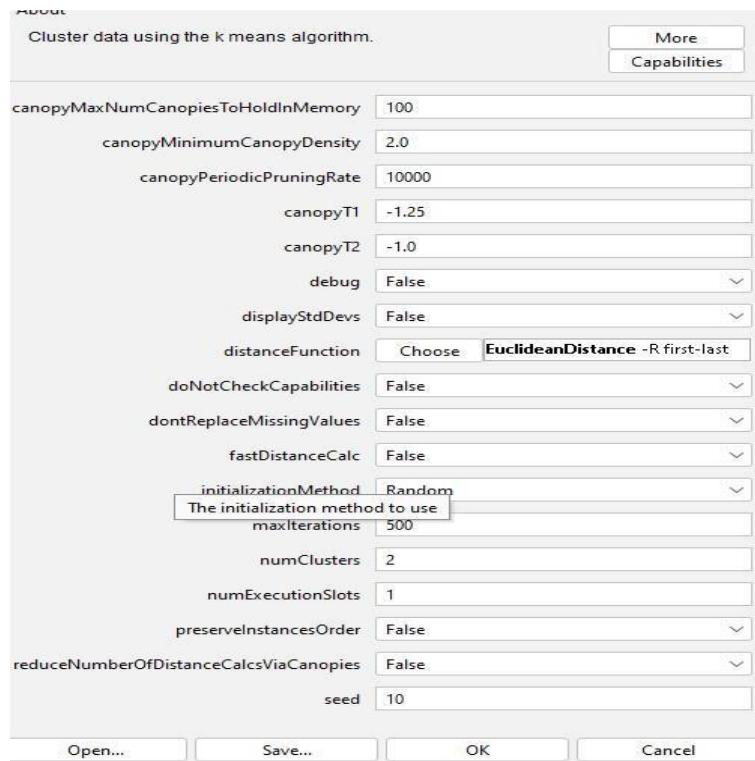
**Step 1):** In the pre-processing interface, open the Weka Explorer and load the required dataset, and we are taking the iris.arff dataset



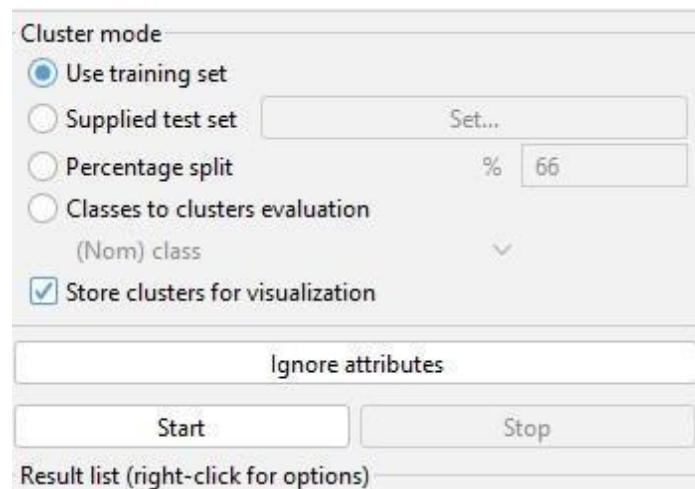
**Step 2):** Find the ‘cluster’ tab in the explorer and press the choose button to execute clustering.



**Step 3):** Then, to the right of the choose icon, press the text button to bring up the popup window shown in the screenshots. We enter three for the number of clusters in this window and leave the seed value alone.



**Step 4):** One of the choices has been chosen. We must ensure that they are in the ‘cluster mode’ panel before running the clustering algorithm. The choice to use a training set is selected, and then the ‘start’ button is pressed.



**Step 5:** The centroid of each cluster is shown in the result window, along with statistics on the number and percent of instances allocated to each cluster. Each cluster centroid is represented by a mean vector. This cluster can be used to describe a cluster.

```

mean           2.7519   3.418   3.1103  2.8088
std. dev.      0.3103   0.3772   0.2952  0.2361

petallength
mean          4.2267   1.464   5.8559  5.0993
std. dev.      0.445    0.1718   0.4626  0.2462

petalwidth
mean          1.3134   0.244   2.1495  1.8254
std. dev.      0.1864   0.1061   0.232   0.2152

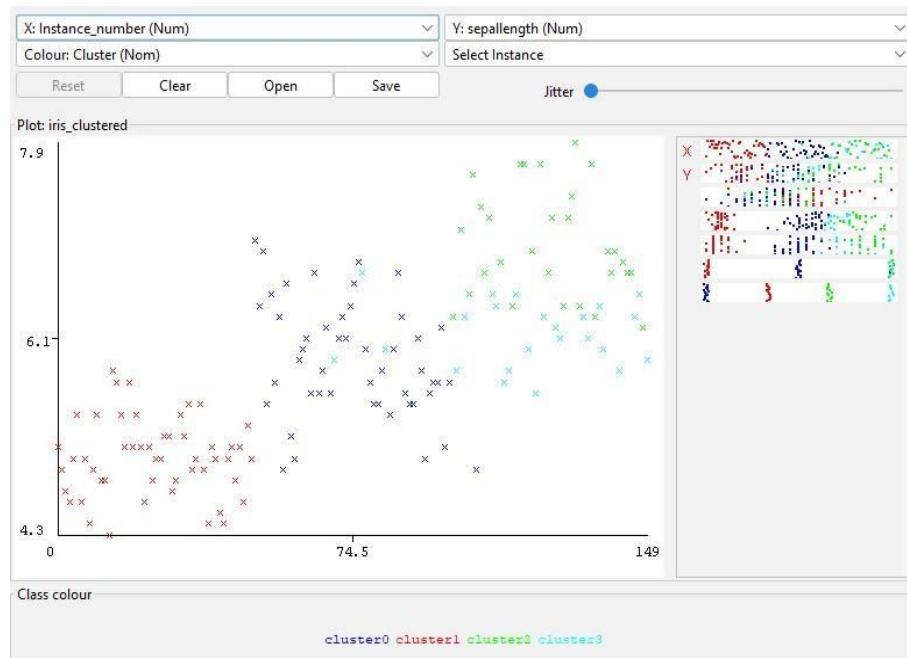
class
Iris-setosa     1       51      1       1
Iris-versicolor 48.1125  1       1.0182  3.8693
Iris-virginica  2.0983   1       31.0375 19.8641
[total]         51.2108  53     33.0557  24.7335

Time taken to build model (full training data) : 0.48 seconds
==== Model and evaluation on training set ====
Clustered Instances
0      48 ( 32%)
1      50 ( 33%)
2      29 ( 19%)
3      23 ( 15%)

Log likelihood: -2.03504

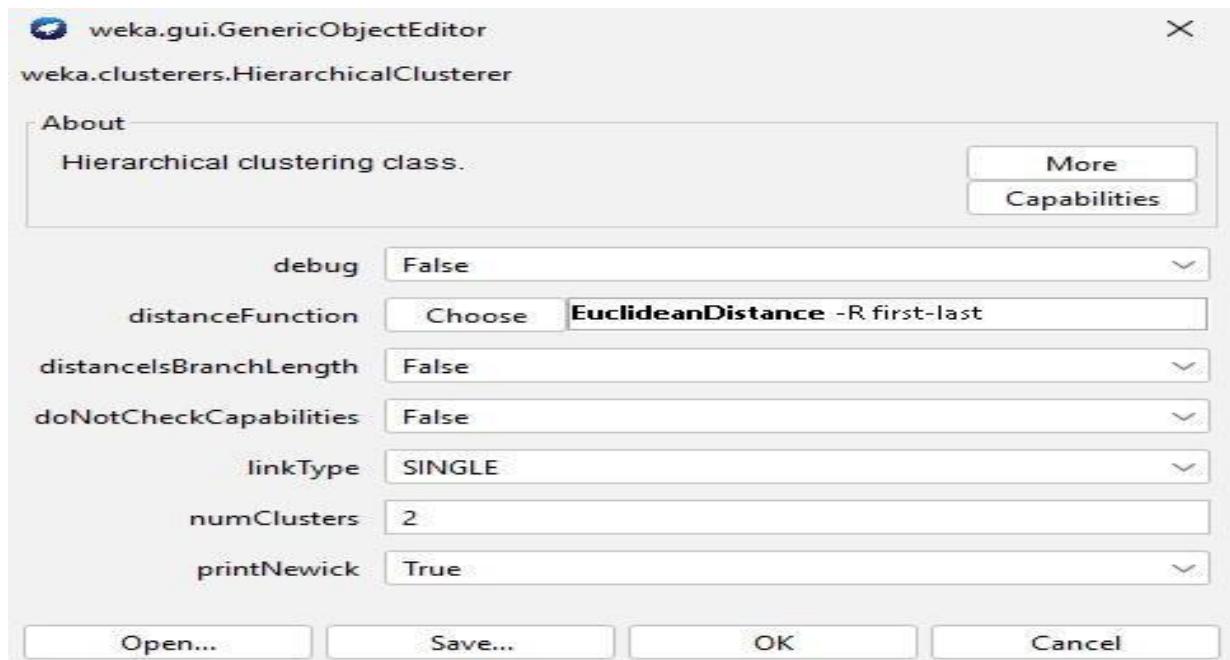
```

**Step 6:** Another way to grasp the characteristics of each cluster is to visualize them. To do so, right-click the result set on the result. Selecting to visualize cluster assignments from the list column.

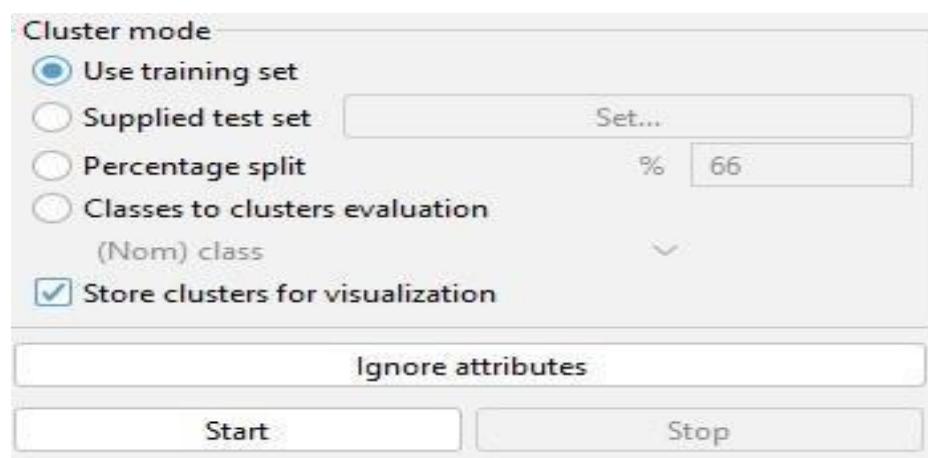


## Method: 2(Hierarchical clustering)

**Step 1):** Then, to the right of the choose icon, press the text button to bring up the popup window shown in the screenshots. We enter three for the number of clusters in this window and leave the seed value alone.



**Step 2):** One of the choices has been chosen. We must ensure that they are in the ‘cluster mode’ panel before running the clustering algorithm. The choice to use a training set is selected, and then the ‘start’ button is pressed.



**Step 3):** The centroid of each cluster is shown in the result window, along with statistics on the number and percent of instances allocated to each cluster. Each cluster centroid is represented by a mean vector. This cluster can be used to describe a cluster.

```
Clusterer output
Scheme: weka.clusterers.HierarchicalClusterer -N 2 -L SINGLE -P -A "weka.core.EuclideanDistance"
Relation: iris
Instances: 150
Attributes: 5
    sepallength
    sepalwidth
    petallength
    petalwidth
    class
Test mode: evaluate on training data

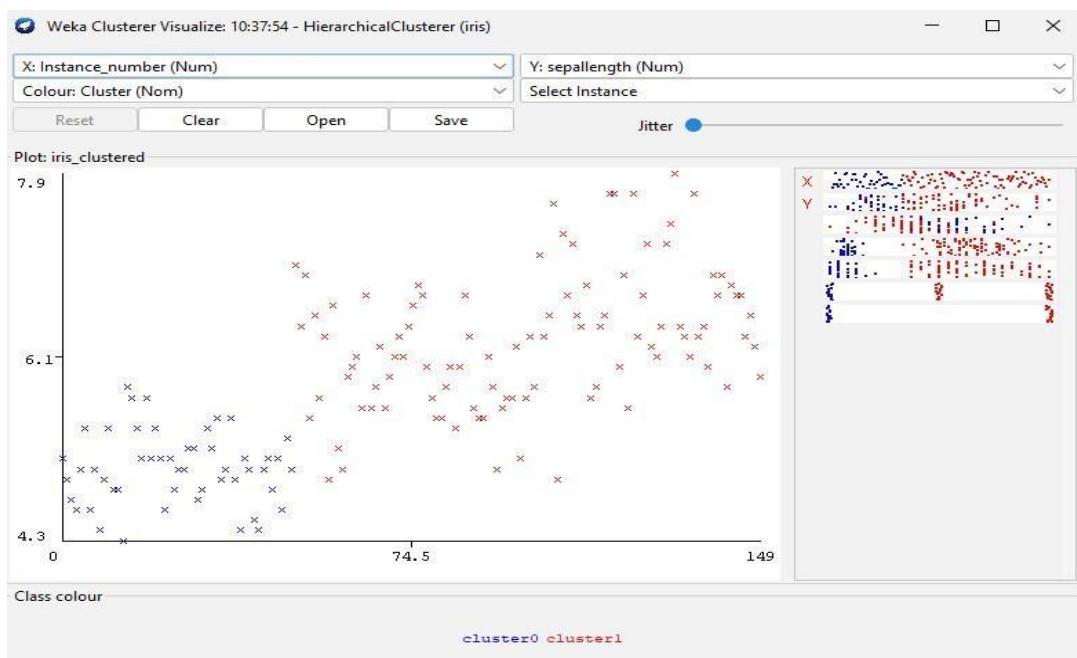
==== Clustering model (full training set) ====

Cluster 0
((0.0:0.03254,0.0:0.03254):0.00913,(0.0:0.03254,0.0:0.03254):0.00913):0.00913

Cluster 1
((1.0:0.07344,((1.0:0.06508,1.0:0.06508):0.00066,(1.0:0.05008,1.0:0.05008):0.00066):0.00066):0.00066

Time taken to build model (full training data) : 0.04 seconds
==== Model and evaluation on training set ====
Clustered Instances
0      50 ( 33%)
1     100 ( 67%)
```

**Step 4):** Another way to grasp the characteristics of each cluster is to visualize them. To do so, right-click the result set on the result. Selecting to visualize cluster assignments from the list column.



## Practical: 9

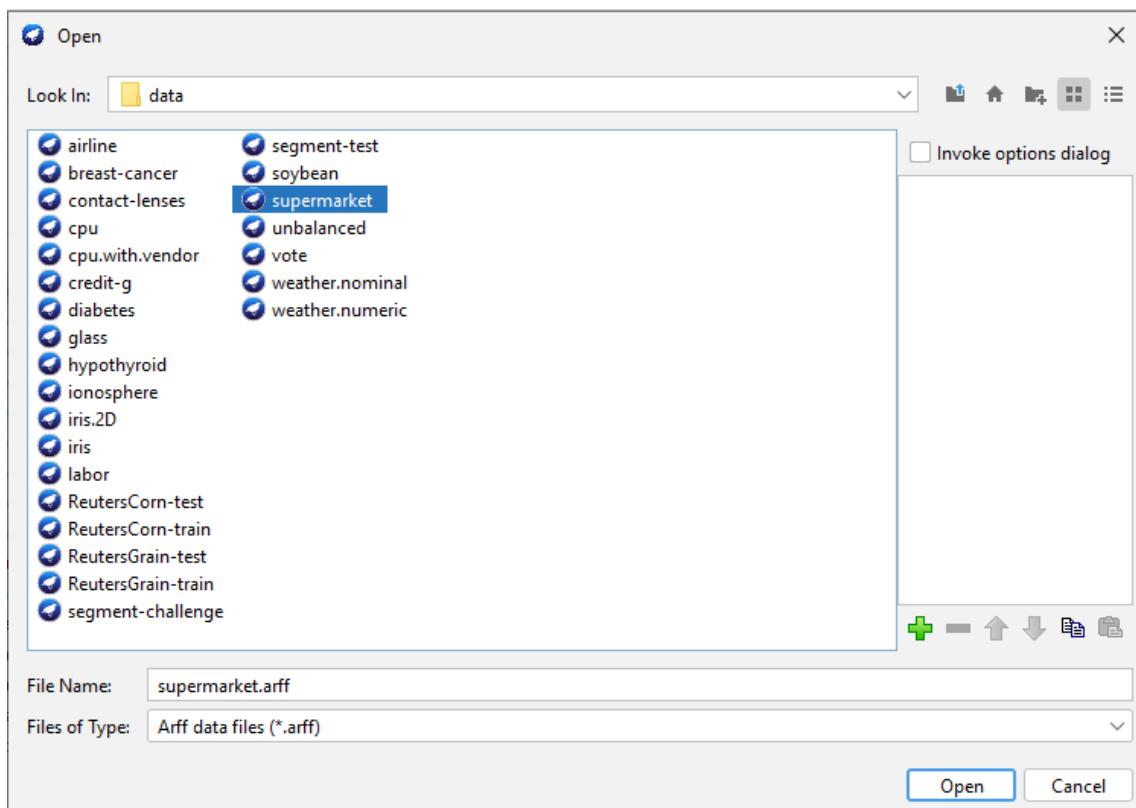
**Aim:** Generate Association rule, build classification model and clustering model for your dataset.

### Code:

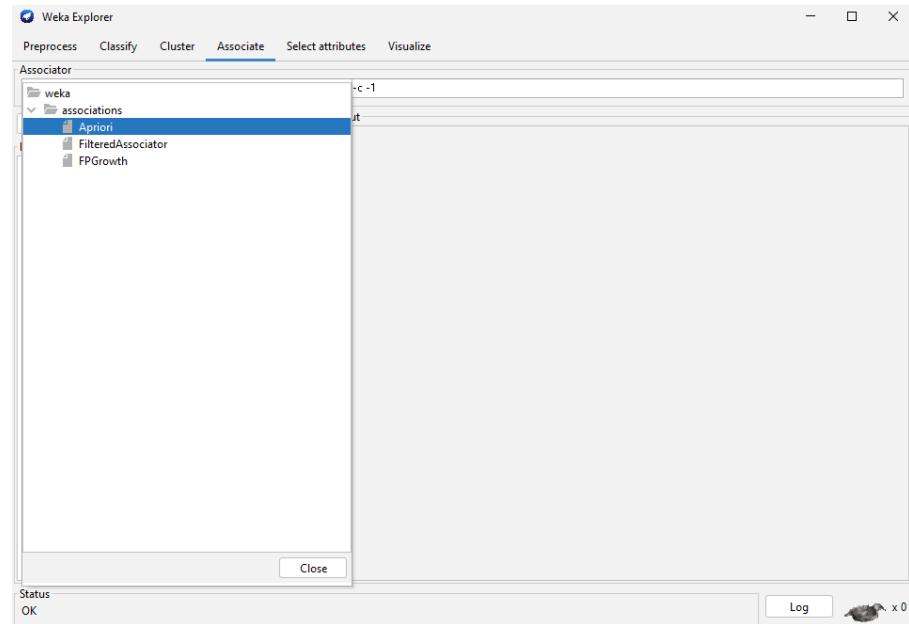
#### ➤ ASSOCIATION RULE

##### **Algorithm 1: Apriori Algorithm**

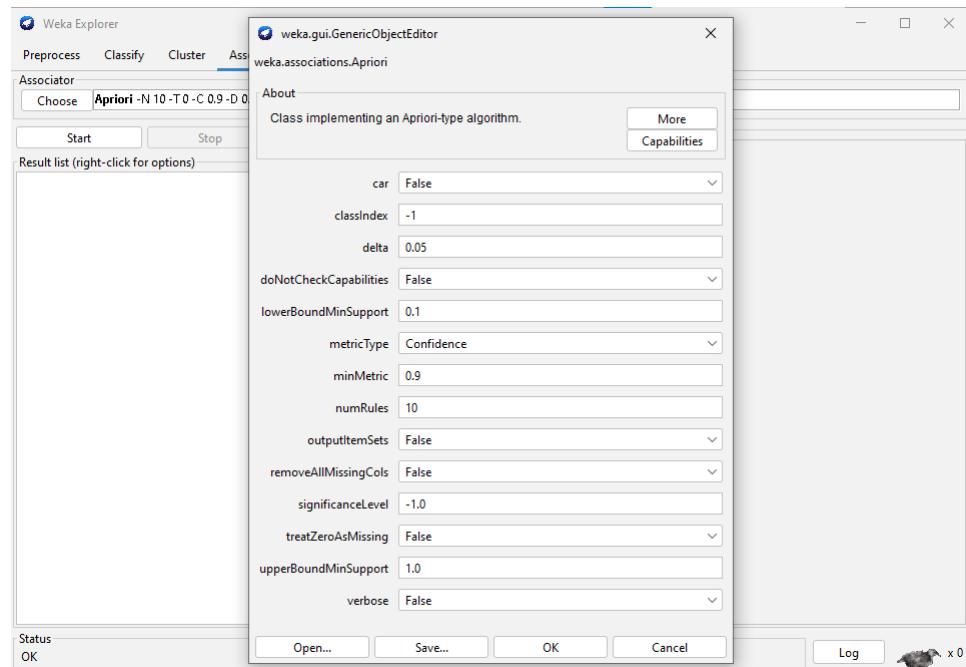
**Step 1)** In the WEKA explorer, open the **Preprocess** tab, click on the **Open file ...** button and select **supermarket.arff** database from the installation folder. After the data is loaded you will see the following screen –



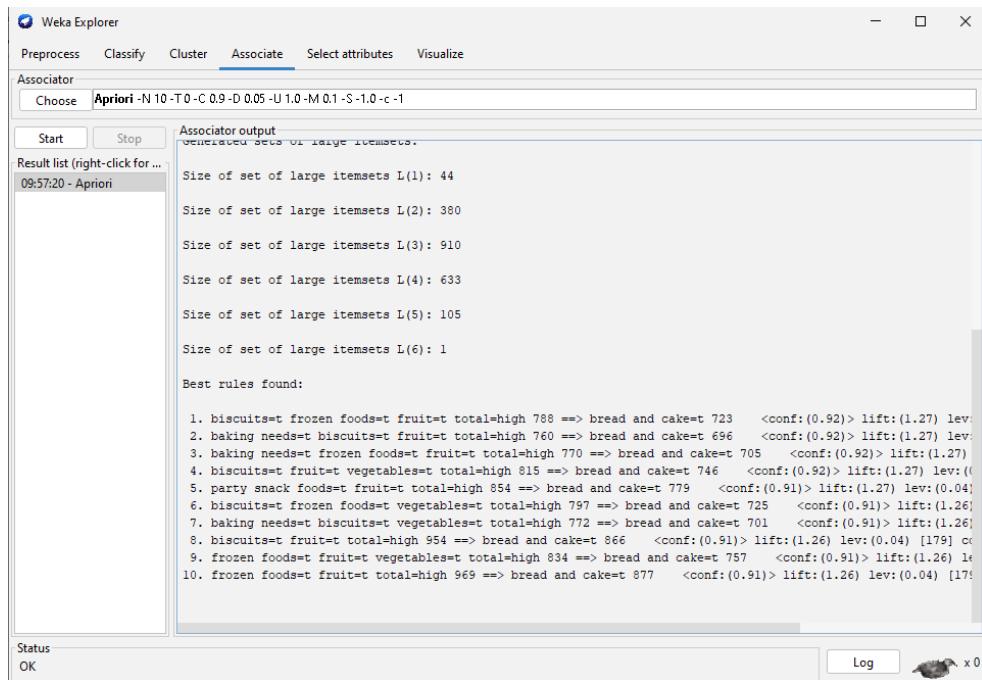
**Step 2)** Click on the **Associate TAB** and click on the **Choose** button. Select the **Apriori** association as shown in the screenshot –



**Step 3)** To set the parameters for the Apriori algorithm, click on its name, a window will pop up as shown below that allows you to set the parameters –



**Step 4)** After you set the parameters, click the **Start** button. After a while you will see the results as shown in the screenshot below –



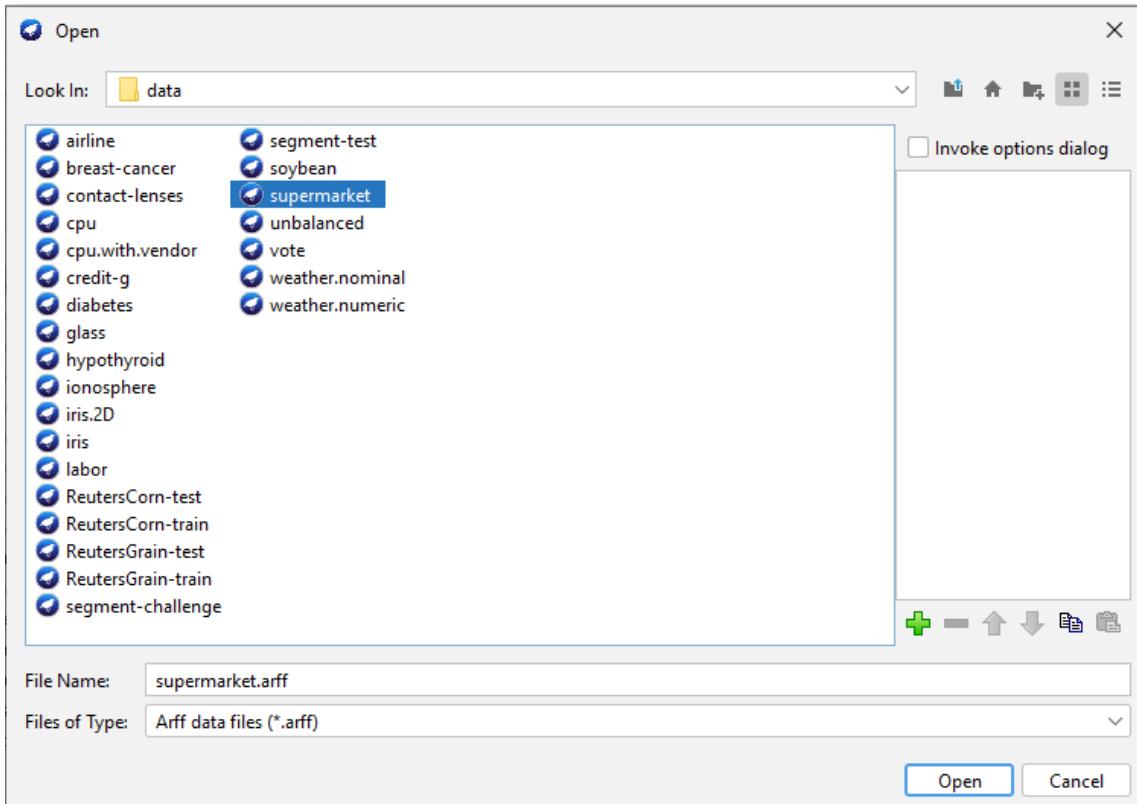
**Step 5)** At the bottom, you will find the detected best rules of associations. This will help the supermarket in stocking their products in appropriate shelves.

Best rules found:

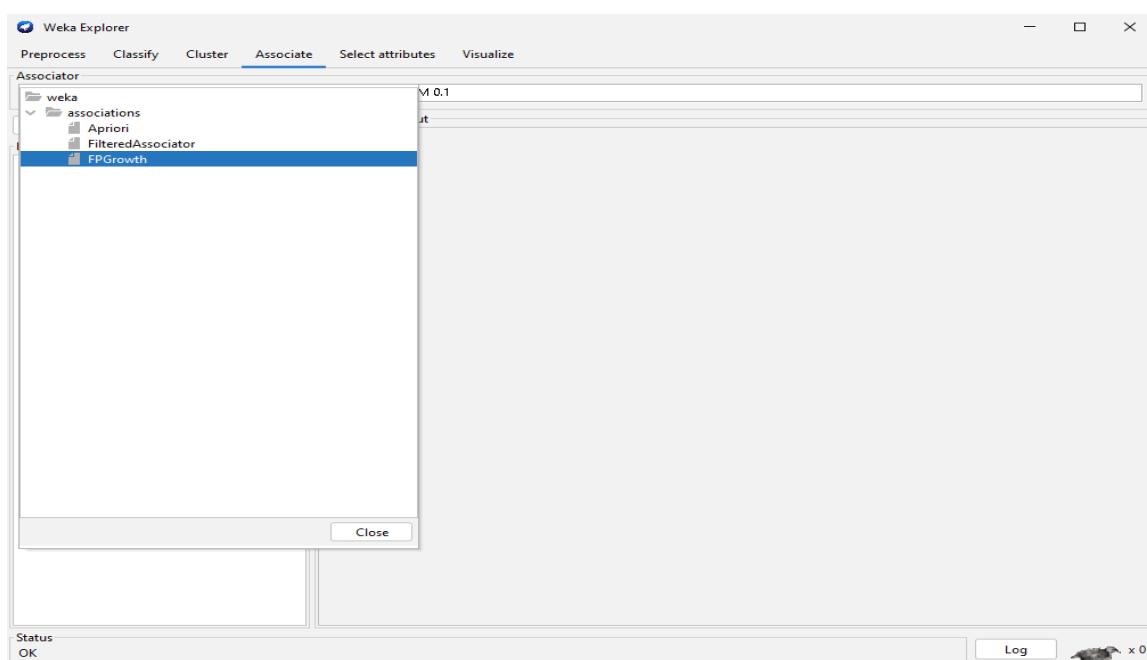
1. biscuits=t frozen foods=t fruit=t total=high 788 => bread and cake=t 723 <conf:(0.92)> lift:(1.27) lev:(0.03) [155] conv:(3.35)
2. baking needs=t biscuits=t fruit=t total=high 760 => bread and cake=t 696 <conf:(0.92)> lift:(1.27) lev:(0.03) [149] conv:(3.28)
3. baking needs=t frozen foods=t fruit=t total=high 770 => bread and cake=t 705 <conf:(0.92)> lift:(1.27) lev:(0.03) [150] conv:(3.27)
4. biscuits=t fruit=t vegetables=t total=high 815 => bread and cake=t 746 <conf:(0.92)> lift:(1.27) lev:(0.03) [159] conv:(3.26)
5. party snack foods=t fruit=t total=high 854 => bread and cake=t 779 <conf:(0.91)> lift:(1.27) lev:(0.04) [179] conv:(3.15)
6. biscuits=t frozen foods=t vegetables=t total=high 797 => bread and cake=t 725 <conf:(0.91)> lift:(1.26) lev:(0.03) [151] conv:(3.06)
7. baking needs=t biscuits=t vegetables=t total=high 772 => bread and cake=t 701 <conf:(0.91)> lift:(1.26) lev:(0.03) [145] conv:(3.01)
8. biscuits=t fruit=t total=high 954 => bread and cake=t 866 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(3)
9. frozen foods=t fruit=t vegetables=t total=high 834 => bread and cake=t 757 <conf:(0.91)> lift:(1.26) lev:(0.03) [156] conv:(3)
10. frozen foods=t fruit=t total=high 969 => bread and cake=t 877 <conf:(0.91)> lift:(1.26) lev:(0.04) [179] conv:(2.92)

## Algorithm 2: FP-Growth Algorithm

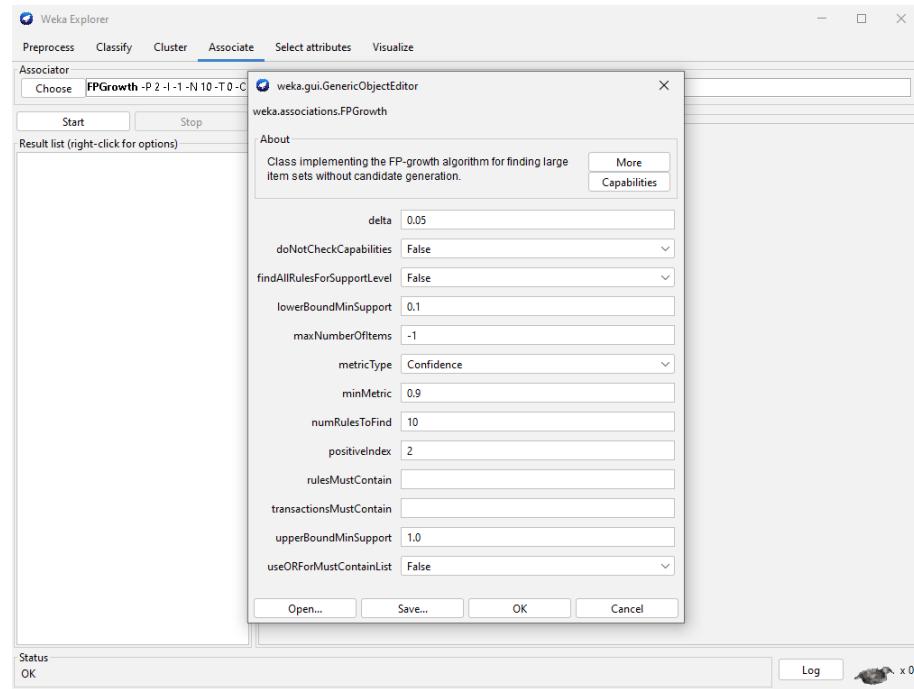
**Step 1)** In the WEKA explorer, open the **Preprocess** tab, click on the **Open file ...** button and select **supermarket.arff** database from the installation folder. After the data is loaded you will see the following screen –



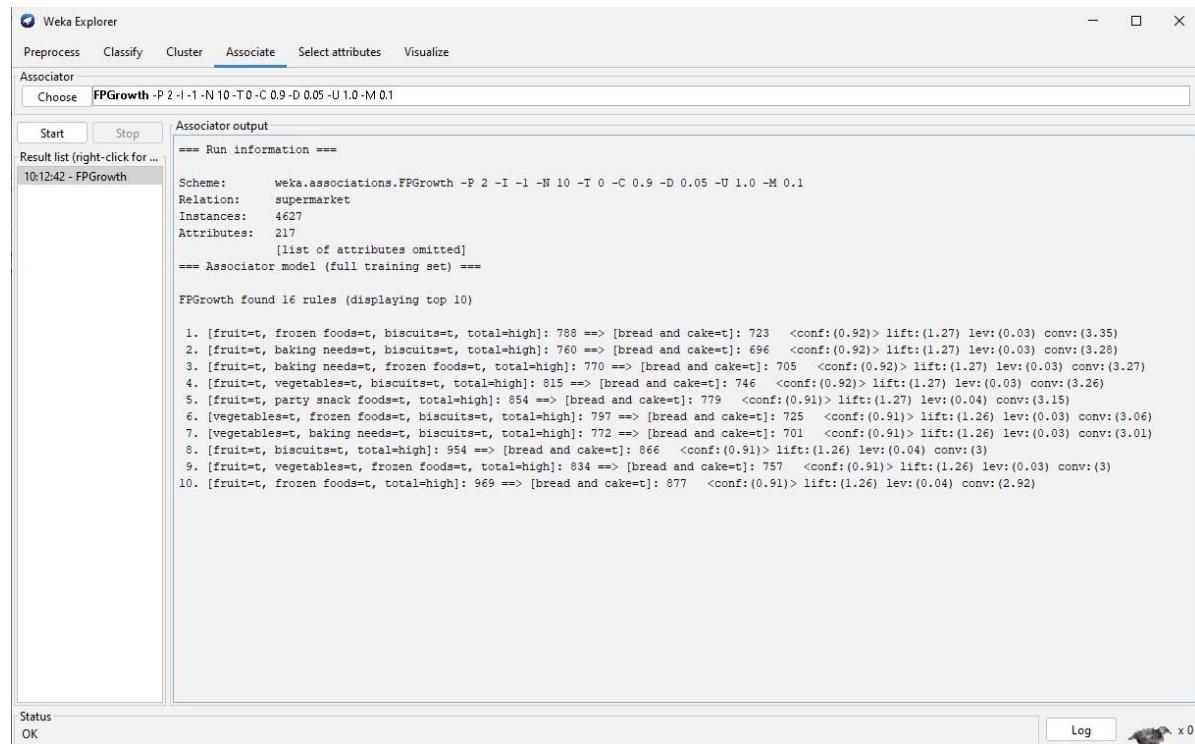
**Step 2)** Click on the **Associate TAB** and click on the **Choose** button. Select the **FPGrowth** association as shown in the screenshot –



**Step 3)** To set the parameters for the FP-Growth algorithm, click on its name, a window will pop up as shown below that allows you to set the parameters –



**Step 4)** After you set the parameters, click the **Start** button. After a while you will see the results as shown in the screenshot below –



**Step 5)** At the bottom, you will find the detected best rules of associations. This will help the supermarket in stocking their products in appropriate shelves.

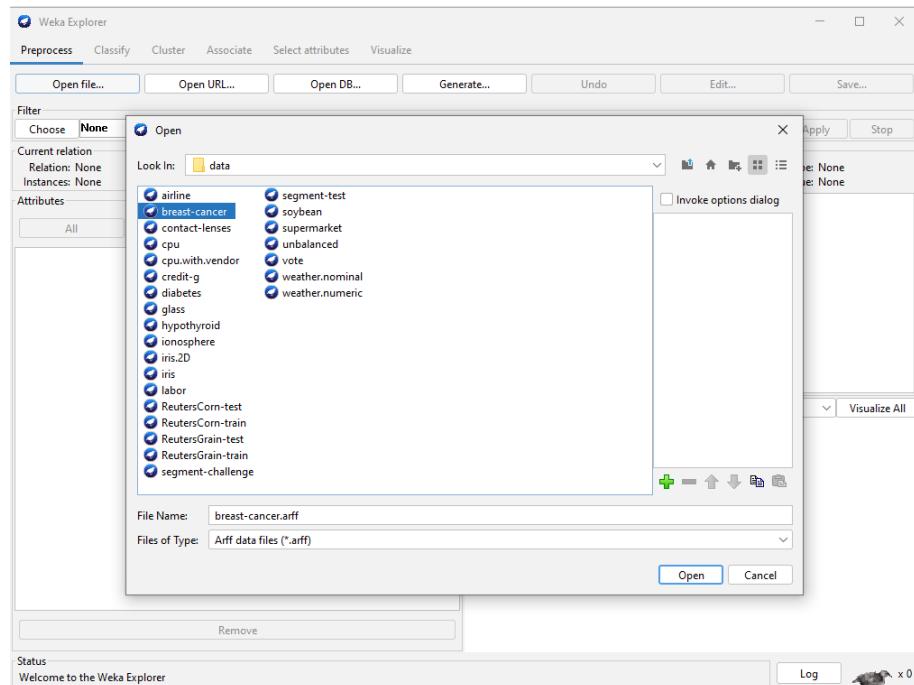
```
FPGrowth found 16 rules (displaying top 10)

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)
```

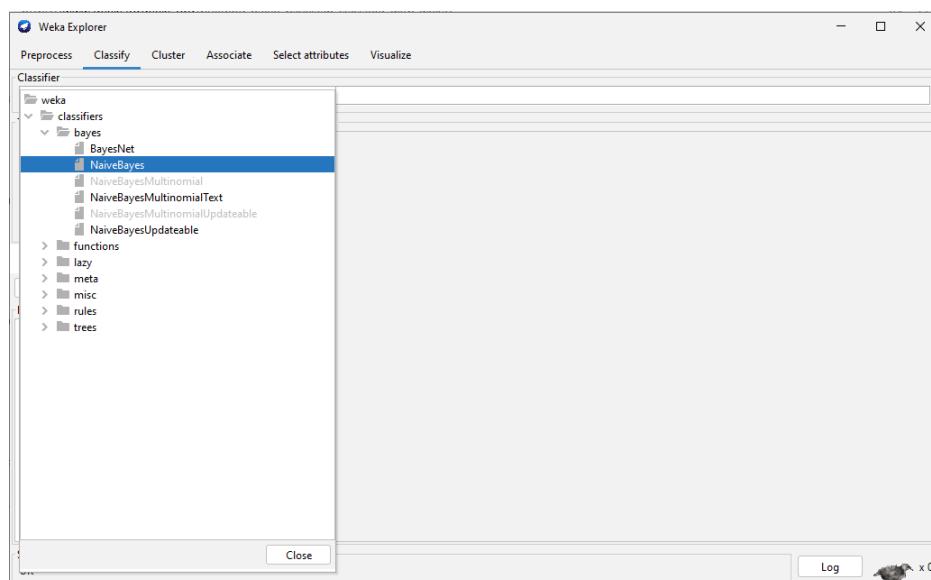
## ➤ CLASSIFICATION MODEL

### Algorithm 1 : NaiveByes

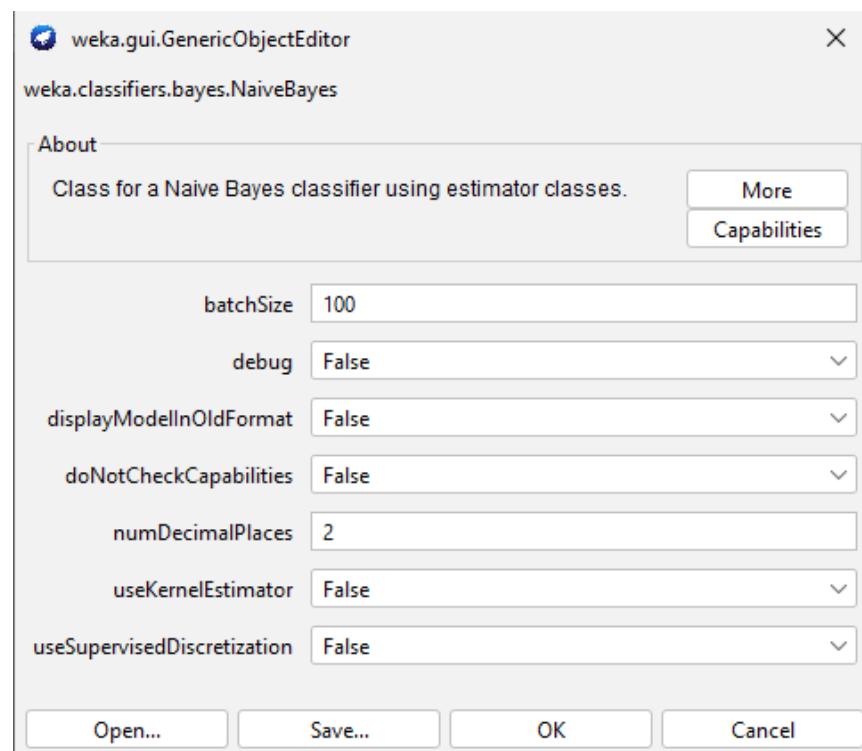
**Step 1)** Initially, we have to load the required dataset in the weka tool using choose file option. Here we are selecting the weather-nominal dataset to execute.



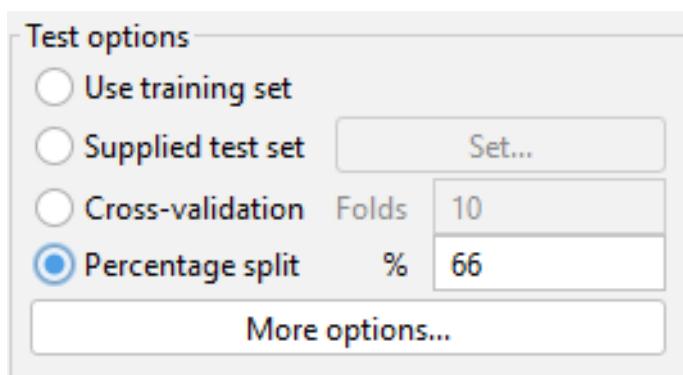
**Step 2)** Now we have to go to the classify tab on the top left side and click on the choose button and select the Naive Bayesian algorithm in it.



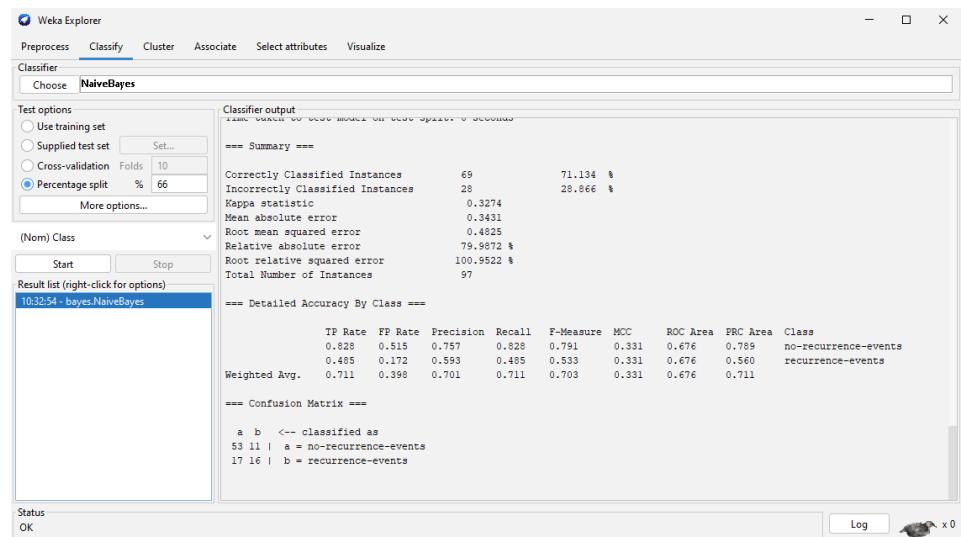
**Step 3)** Now to change the parameters click on the right side at the choose button, and we are accepting the default values in this example.



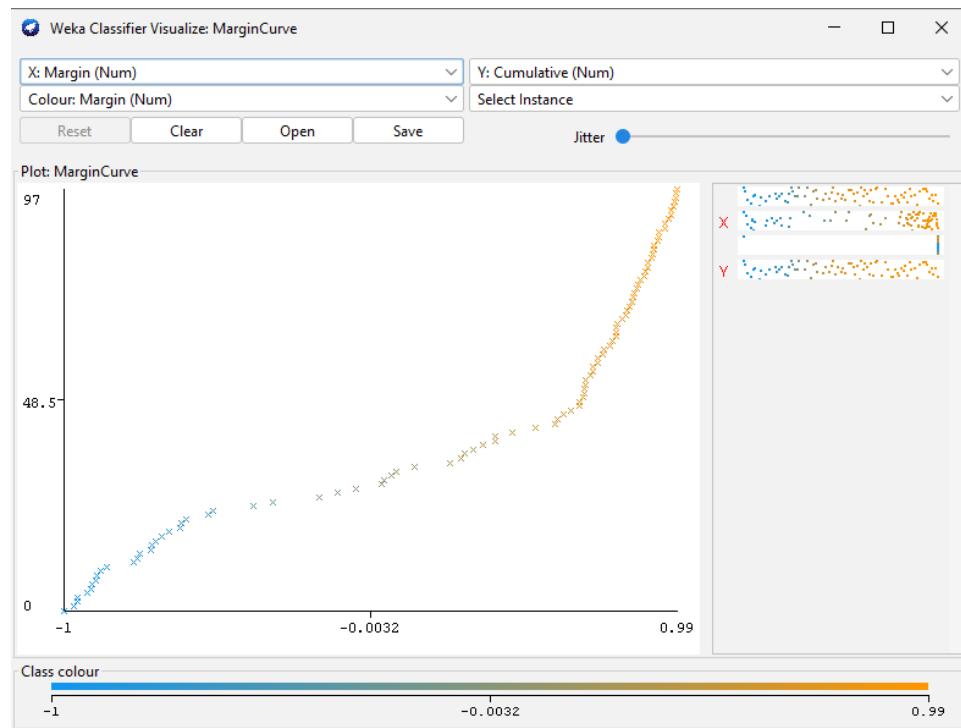
**Step 4)** We choose the Percentage split as our measurement method from the “Test” choices in the main panel. Since we don’t have a separate test data collection, we’ll use the percentage split of 66 percent to get a good idea of the model’s accuracy. Our dataset contains 14 examples, with h9 being used for training and 5 being used for testing.



**Step 5)** To generate the model, we now click “start.” When the model is done, the evaluation statistic will appear in the right panel.

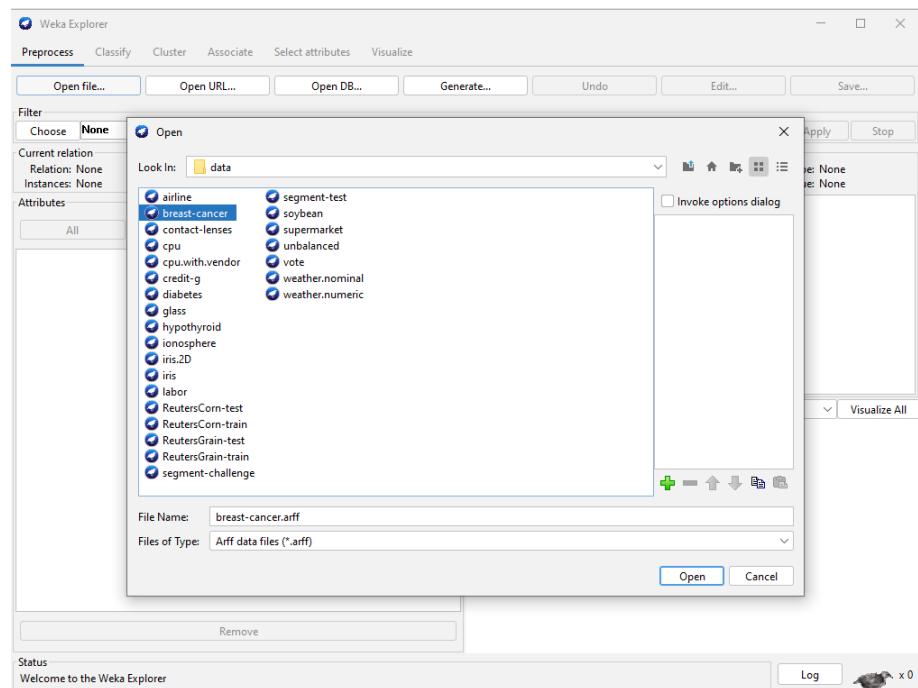


**Step 6)** Visualize the graph

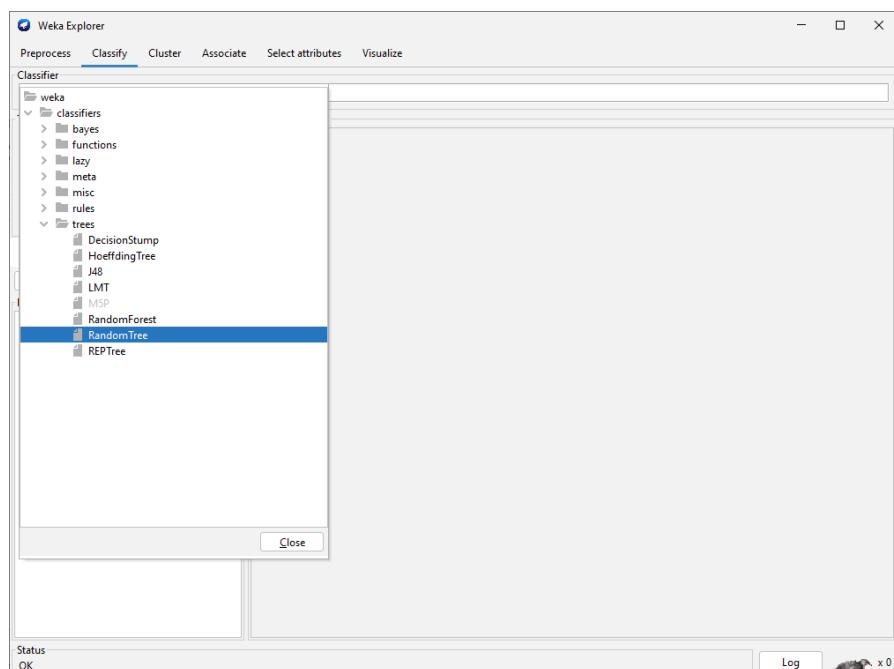


## Algorithm 2: Random Tree

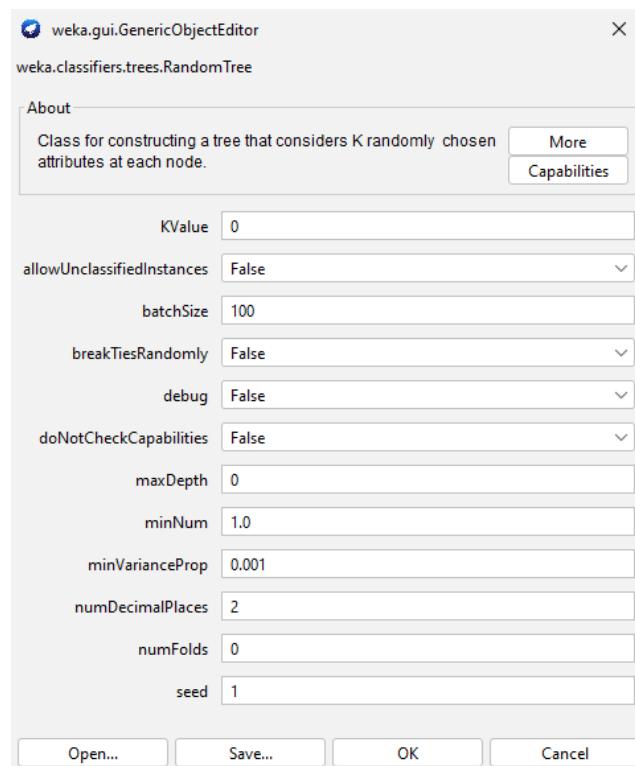
**Step 1)** Initially, we have to load the required dataset in the weka tool using choose file option. Here we are selecting the weather-nominal dataset to execute.



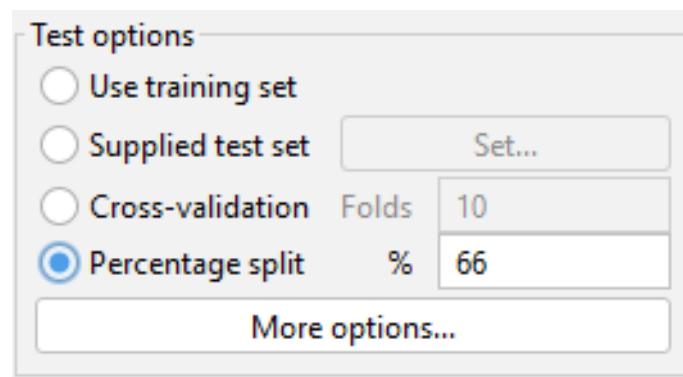
**Step 2)** Now we have to go to the classify tab on the top left side and click on the choose button and select the Random Tree algorithm in it.



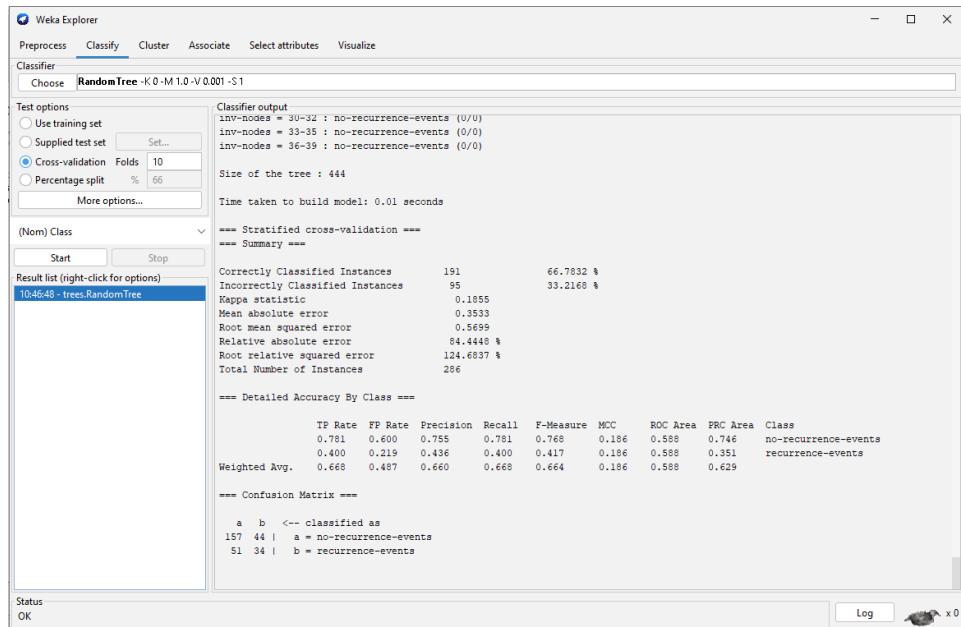
**Step 3)** Now to change the parameters click on the right side at the choose button, and we are accepting the default values in this example.



**Step 4)** We choose the Percentage split as our measurement method from the “Test” choices in the main panel. Since we don’t have a separate test data collection, we’ll use the percentage split of 66 percent to get a good idea of the model’s accuracy. Our dataset contains 14 examples, with h9 being used for training and 5 being used for testing.

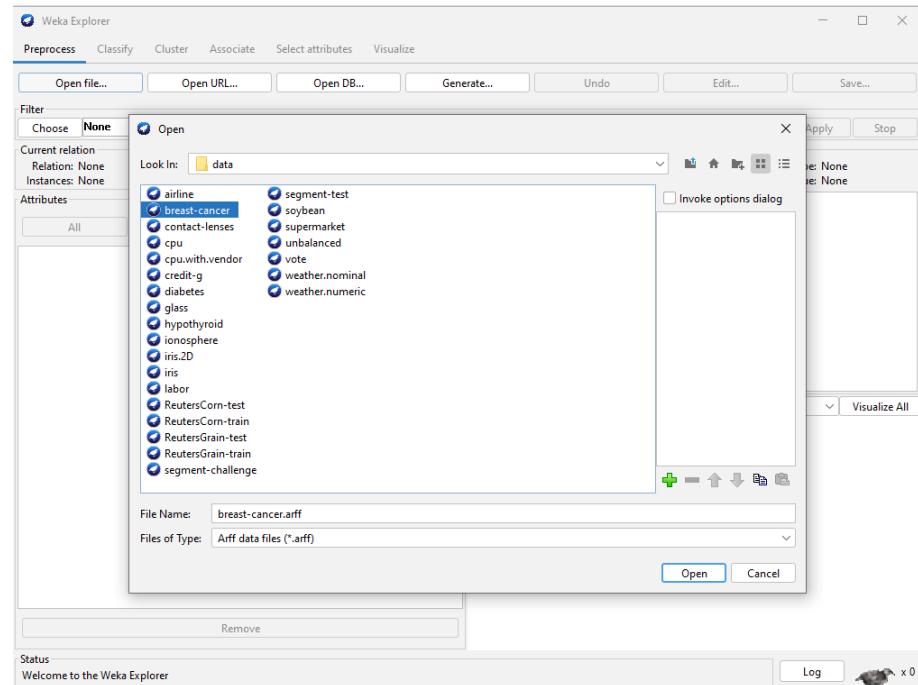


**Step 5)** To generate the model, we now click “start.” When the model is done, the evaluation statistic will appear in the right panel.

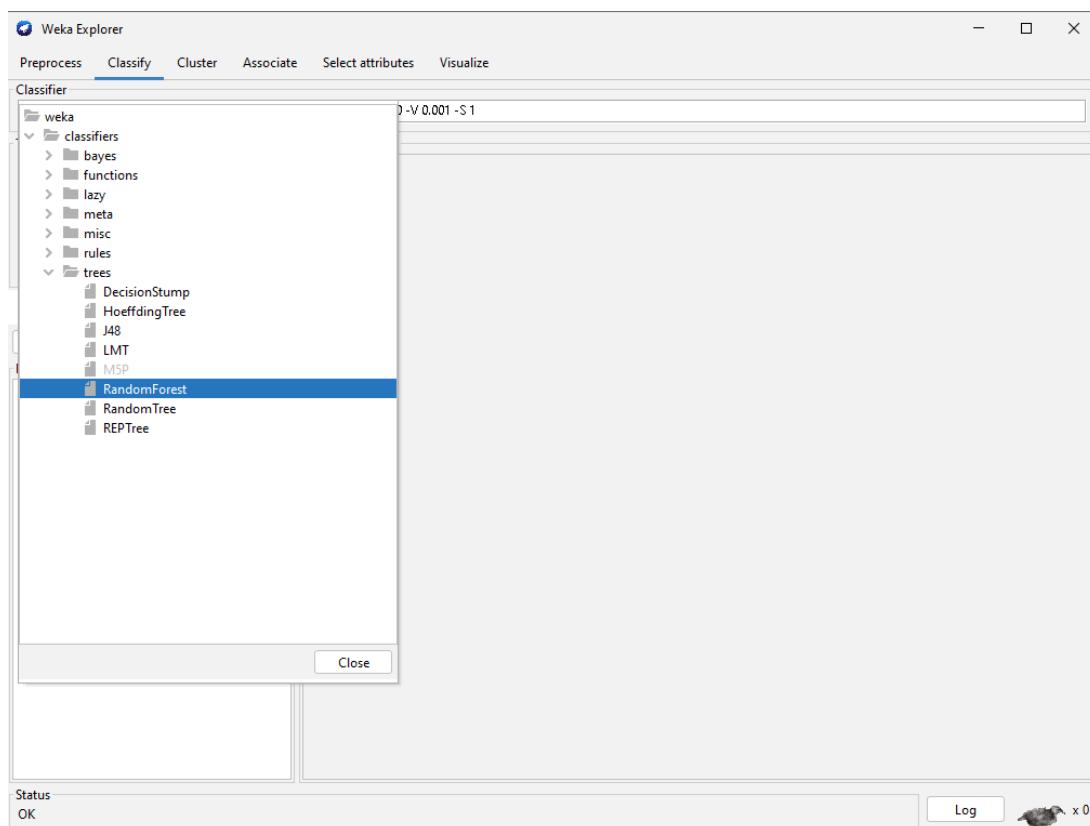


### Algorithm 3: Random Forest

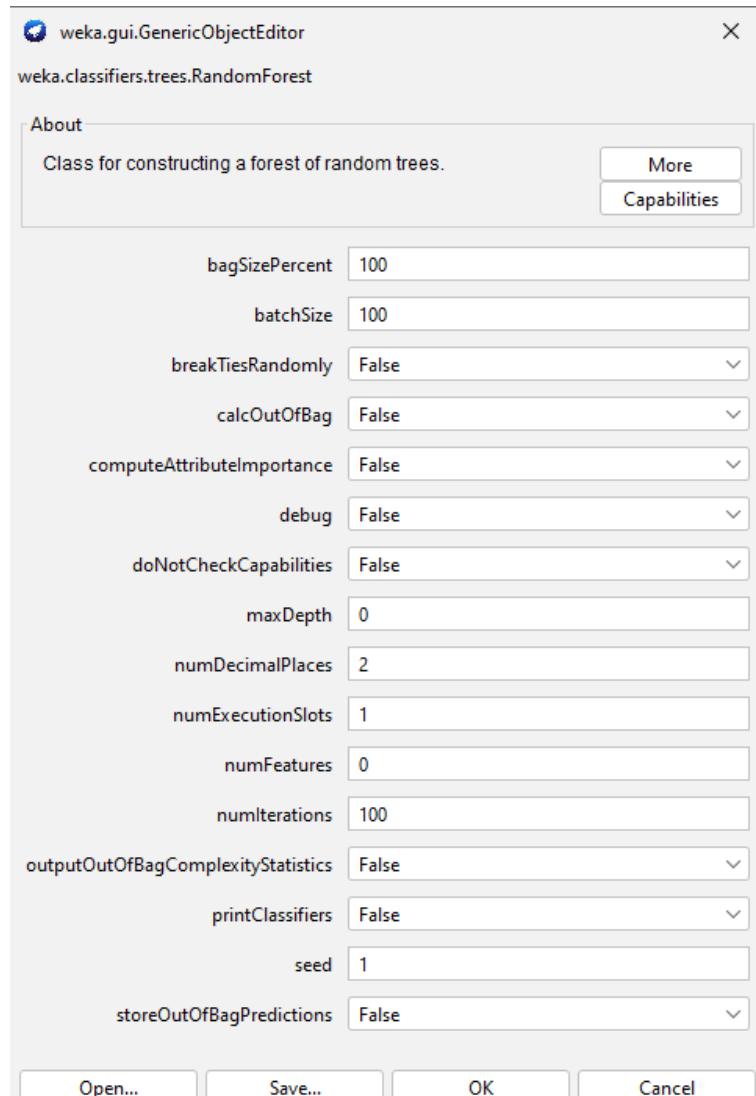
**Step 1)** Initially, we have to load the required dataset in the weka tool using choose file option. Here we are selecting the weather-nominal dataset to execute.



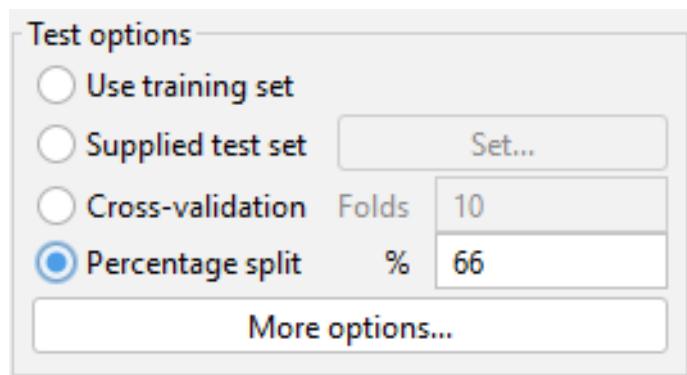
**Step 2)** Now we have to go to the classify tab on the top left side and click on the choose button and select the Random Tree algorithm in it.



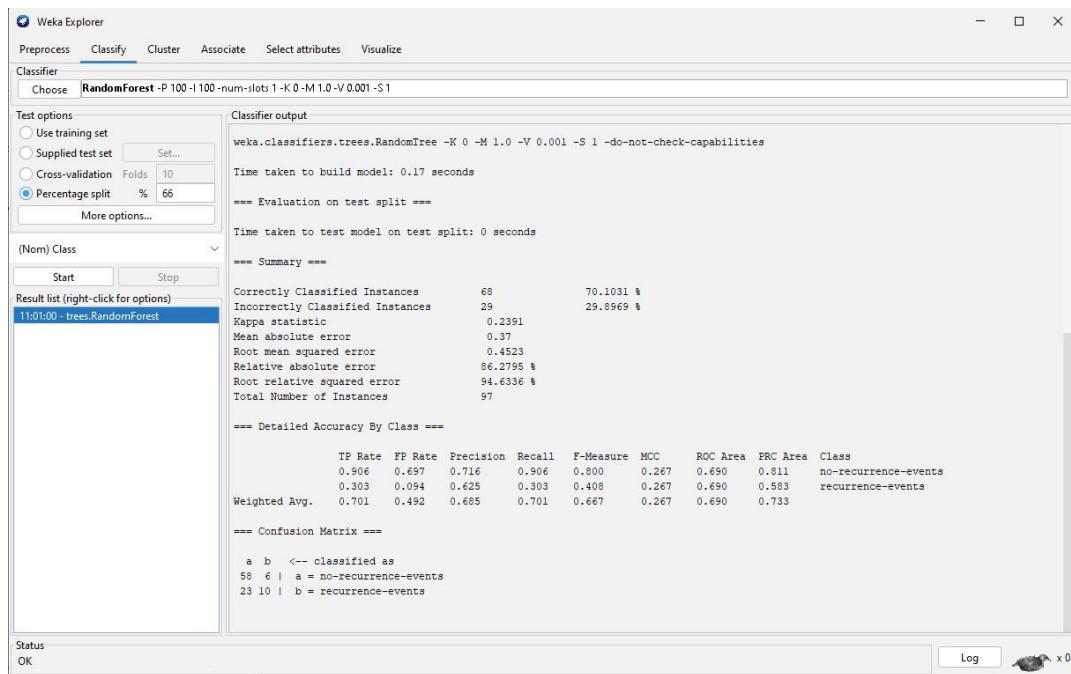
**Step 3)** Now to change the parameters click on the right side at the choose button, and we are accepting the default values in this example.



**Step 4)** We choose the Percentage split as our measurement method from the “Test” choices in the main panel. Since we don’t have a separate test data collection, we’ll use the percentage split of 66 percent to get a good idea of the model’s accuracy. Our dataset contains 14 examples, with h9 being used for training and 5 being used for testing.



**Step 5)** To generate the model, we now click “start.” When the model is done, the evaluation statistic will appear in the right panel.



## **Practical: 10**

**Aim:** Calculating Information gains measures.

**Code:**

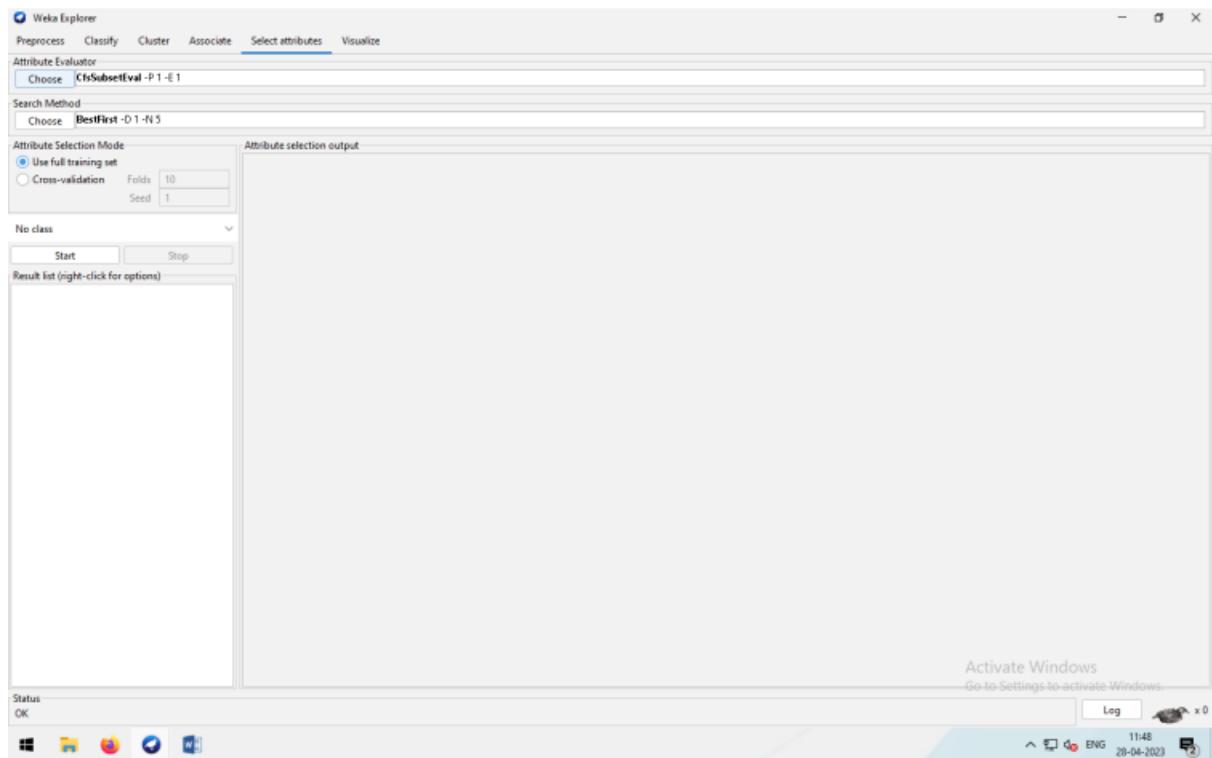
- Information gain (IG) measures how much “information” a feature gives us about the class. – Features that perfectly partition should give maximal information. – Unrelated features should give no information.
- It measures the reduction in entropy.
- CfsSubsetEval aims to identify a subset of attributes that are highly correlated with the target while not being strongly correlated with one another. It searches through the space of possible attribute subsets for the “best” one using theBestFirst search method by default, although other methods can be chosen.
- To use the wrapper method rather than a filter method, such as CfsSubsetEval, first select WrapperSubsetEval and then configure it by choosing a learning algorithm to apply and setting the number of crossvalidation folds to use when evaluating it on each attribute subset.  
Steps:

- Open WEKA Tool.
- Click on WEKA Explorer.
- Click on Preprocessing tab button.
- Click on open file button.
- Select and Click on data option button.
- Choose a data set and open file.
- Click on select attribute tab and Choose attribute evaluator, search method algorithm

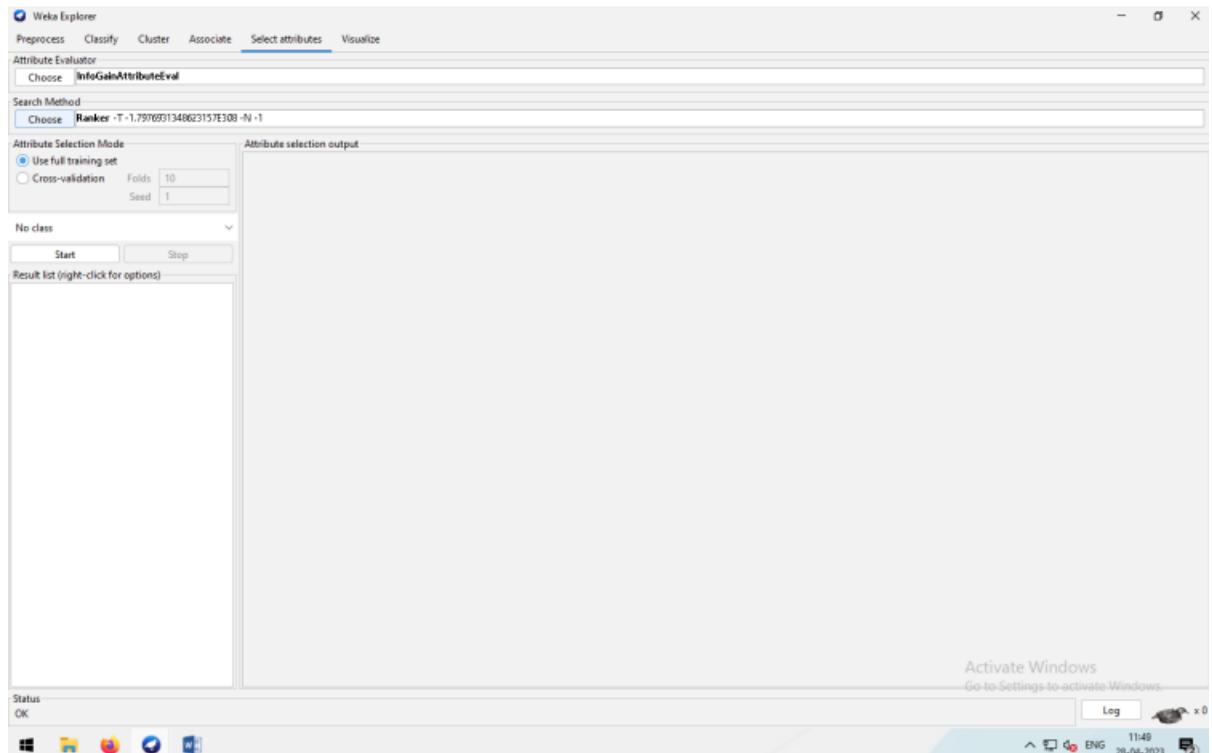
**Step 1)** In the WEKA explorer, open the **Preprocess** tab, click on the **Open file ...** button and select **supermarket.arff** database from the installation folder. After the data is loaded you will see the following screen –

No.	Label	Count	Weight
1	april	26	26
2	may	75	75
3	june	93	93
4	july	118	118
5	august	131	131
6	september	149	149
7	october	90	90

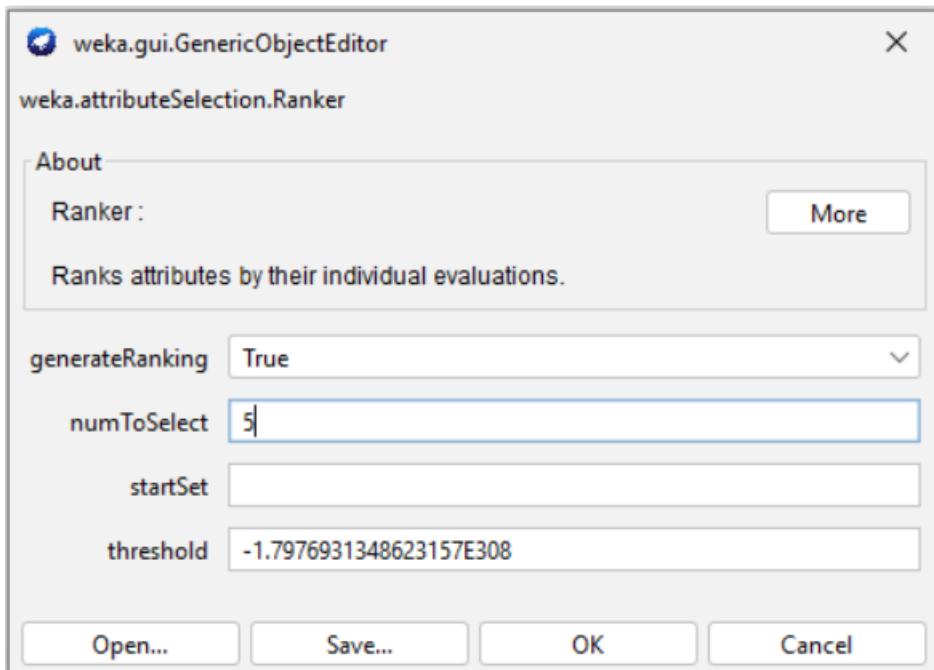
## Step 2) Click on Select Attributes tab



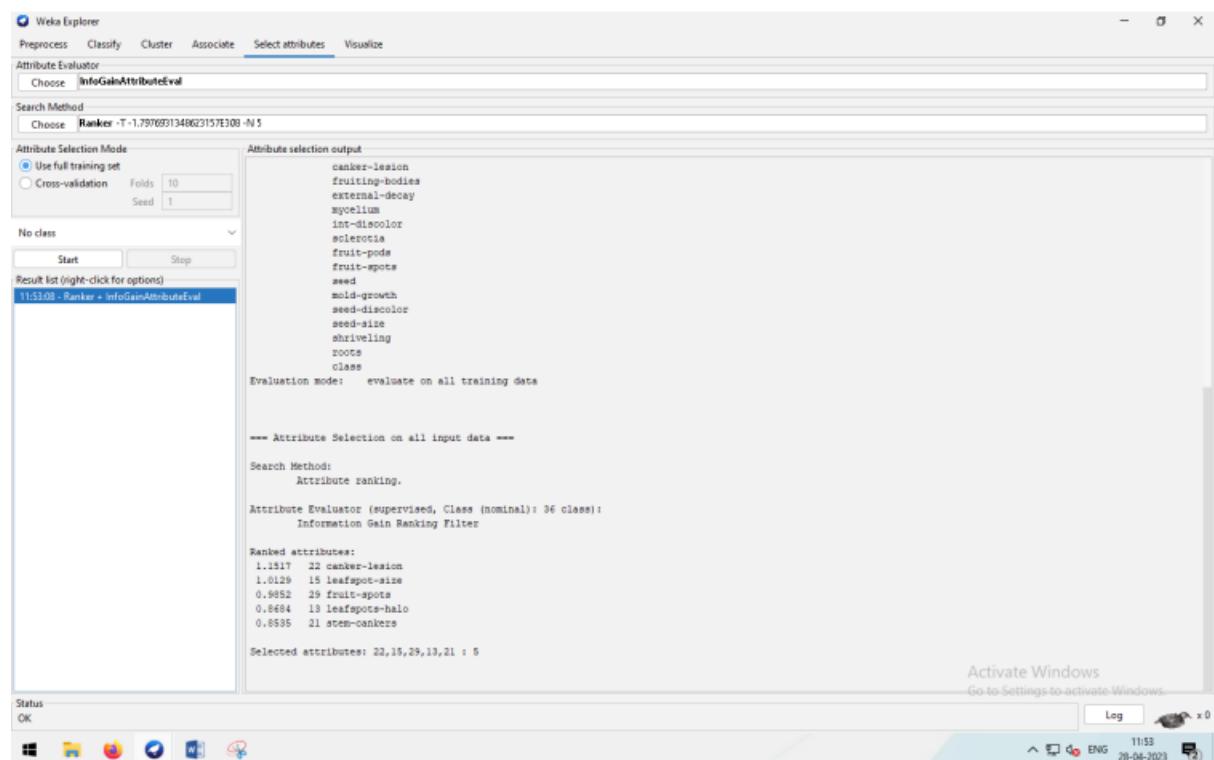
## Step 3) Click on Choose [Attribute Evaluator -> InfoGainAttributeEval], [Search Method -> Ranker]



**Step 4)** Back to our task, where we need to find 5 and 10 best attributes, so first I will show how to get the best 5 attributes, then we just need to repeat step 5 to do for 10 best attributes. Click on Ranker, type 5 in numToSelect. Then click okay.



**Step 5)** - For Attribute Selection Mode, by default it will tick Use for training set, so don't change anything. Click Start. Then it will display the result at the right hand side.



From the result, we get 5 best attributes based on Information Gain and Ranker search method.

```

Ranked attributes:
1.1517 22 canker-lesion
1.0129 15 leafspot-size
0.9852 29 fruit-spots
0.8684 13 leafspots-halo
0.8535 21 stem-cankers

```

Output:

```

Attribute selection output
canker-lesion
fruiting-bodies
external-decay
mycelium
int-discolor
sclerotia
fruit-pods
fruit-spots
seed
mold-growth
seed-discolor
seed-size
shriveling
roots
class
Evaluation mode: evaluate on all training data

==== Attribute Selection on all input data ====
Search Method:
    Attribute ranking.

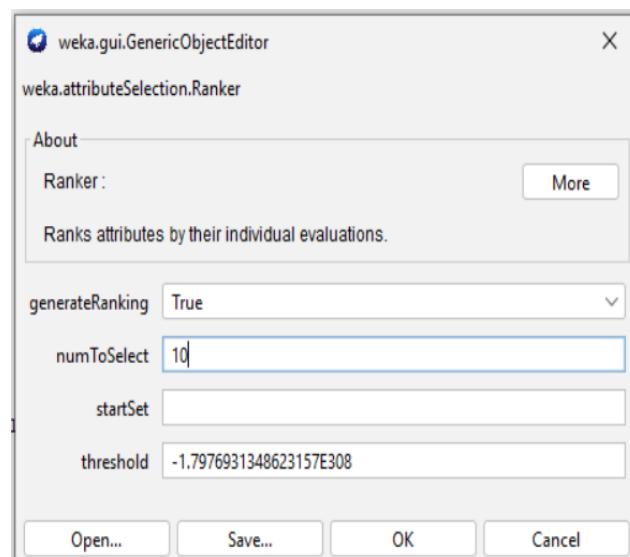
Attribute Evaluator (supervised, Class (nominal): 36 class):
    Information Gain Ranking Filter

Ranked attributes:
1.1517 22 canker-lesion
1.0129 15 leafspot-size
0.9852 29 fruit-spots
0.8684 13 leafspots-halo
0.8535 21 stem-cankers

Selected attributes: 22,15,29,13,21 : 5

```

## Step 6)



## Step 7)

```
Ranked attributes:  
1.1517 22 canker-lesion  
1.0129 15 leafspot-size  
0.9852 29 fruit-spots  
0.8684 13 leafspots-halo  
0.8535 21 stem-cankers  
0.8504 14 leafspots-marg  
0.8437 28 fruit-pods  
0.6918 19 stem  
0.6715 1 date  
0.6265 11 plant-growth
```

## Step 8)

```
Attribute selection output  
sclerotia  
fruit-pods  
fruit-spots  
seed  
mold-growth  
seed-discolor  
seed-size  
shriveling  
roots  
class  
Evaluation mode: evaluate on all training data  
  
==== Attribute Selection on all input data ====  
Search Method:  
Attribute ranking.  
Attribute Evaluator (supervised, Class (nominal): 36 class):  
Information Gain Ranking Filter  
Ranked attributes:  
1.1517 22 canker-lesion  
1.0129 15 leafspot-size  
0.9852 29 fruit-spots  
0.8684 13 leafspots-halo  
0.8535 21 stem-cankers  
0.8504 14 leafspots-marg  
0.8437 28 fruit-pods  
0.6918 19 stem  
0.6715 1 date  
0.6265 11 plant-growth  
Selected attributes: 22,15,29,13,21,21,14,28,19,1,11 : 10
```