# Experiment 7 - Perform data Pre-processing task and Demonstrate performing Classification, Clustering, Association algorithm on data sets using data mining tool

| | | | | | |
|---|---|---|---|---|---|
| **Name:** | **Shaikh Shadab** | | **Rollno** | **:** | **17DCO74** |
| **Class :** | **TE.CO** | | **Batch** | **:** | **B3** |

## #Theory

➢ Weka Tool

Weka is a collection of machine learning algorithms for data mining tasks. It contains tools for data preparation, classification, regression, clustering, association rules mining, and visualization. Found only on the islands of New Zealand, the Weka is a flightless bird with an inquisitive nature. The name is pronounced like this, and the bird sounds like this. Weka is open source software issued under the GNU General Public License. It is possible to apply Weka to process big data and perform deep learning!

➢ Testing methods

- Use training set - Means we will test our knowledge on the same data we learned. Not very accepted because we can just make build our code to memorize the training instances (which will be in the test).
- Supplied test set - it is an external file that we can use as training set. It can be used when we want/need to test the algorithm's knowledge against a specific test set.
- Percentage split: Splits the data and separates x% of the data for learning and the rest of it for testing. It is useful when your algorithm is slow.
- Cross-validation (CV): Works like many percentage splits. we fold the data in 10 folds (for example) and repeat 10 (because it is 10-folds) the following process: Use 9 folds for learning and leave 1-fold out for testing. Every time leaving a different fold for testing. This is the most used testing method in papers. They say "anything over 5 folds is acceptable", but no one has any good explanation for that.

➢ Confusion Matrix

A confusion matrix is a table that is often used to describe the performance of a classification model (or "classifier") on a set of test data for which the true values are known. It allows the visualization of the performance of an algorithm. It allows easy identification of confusion between classes e.g. one class is commonly mislabelled as the other. Most performance measures are computed from the confusion matrix.

```
=== Confusion Matrix ===

  a  b  c   <-- classified as
 50  0  0 |  a = Iris-setosa
  0 48  2 |  b = Iris-versicolor
  0  3 47 |  c = Iris-virginica
```

For example, from the above table we can say that the actual class was a i.e Iris-setosa and predicted class a count was 50, b is 0 and c is 0.similaryly the actual class was c i.e Iris-virginica and predicted class a count was 0 b is 3 and c is 47 (using CV 2 fold)

➤ True positive, True negative, False positive, False negative
  • When the actual class was true and we predict it also as true is TP
  • When the actual class was false and we predict it also as false is TN
  • When the actual class was false and we predict it as true is FP (Type I error)
  • When the actual class was true and we predict it as false is FN (Type II error)

In simpler terms, A true positive is an outcome where the model correctly predicts the positive class. Similarly, a true negative is an outcome where the model correctly predicts the negative class.

A false positive is an outcome where the model incorrectly predicts the positive class. And a false negative is an outcome where the model incorrectly predicts the negative class.

Real life example of TP (correctly identified)
Umpire gives a Batsman NOT OUT when he is NOT OUT.
Real life example of TN (correctly rejected)
Umpire gives a Batsman OUT when he is OUT.
Real life example of FP (Type I error) (incorrectly identified)
Umpire gives a Batsman NOT OUT when he is OUT.
Real life example of FN (Type II error) (incorrectly rejected)
Umpire gives a Batsman OUT when he is NOT OUT.

➤ Sensitivity, Specificity, significance
  True positive rate (sensitivity) - measures the proportion of actual positives that are correctly identified as such
  True negative rate (Specificity) - measures the proportion of actual negatives that are correctly identified as such
  False positive rate (significance) - proportion of all negatives that still yield positive test outcomes, i.e., the conditional probability of a positive test result given an event that was not present.

➤ Precision and Recall
  Precision tries to answer such type of queries:
  What proportion of positive identifications was actually correct?

$$\text{Precision} = \frac{TP}{TP + FP}$$

  Recall tries to answer such type of queries:
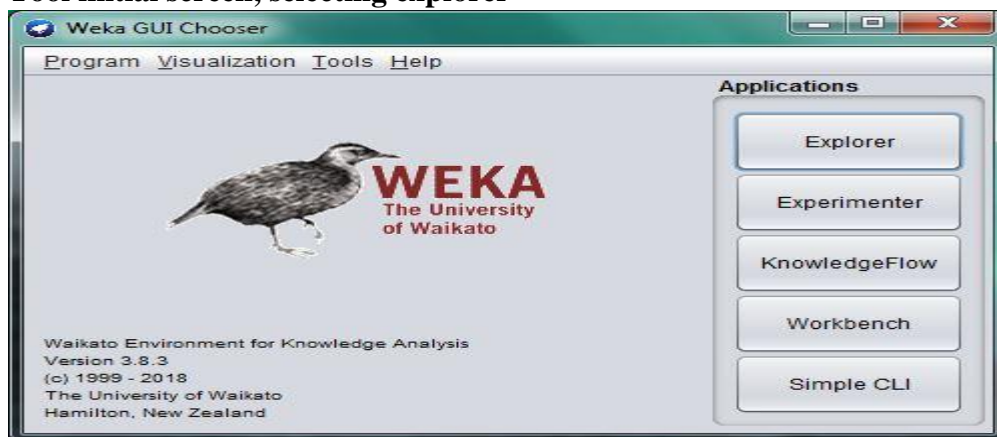  What proportion of actual positives was identified correctly?

$$\text{Recall} = \frac{TP}{TP + FN}$$

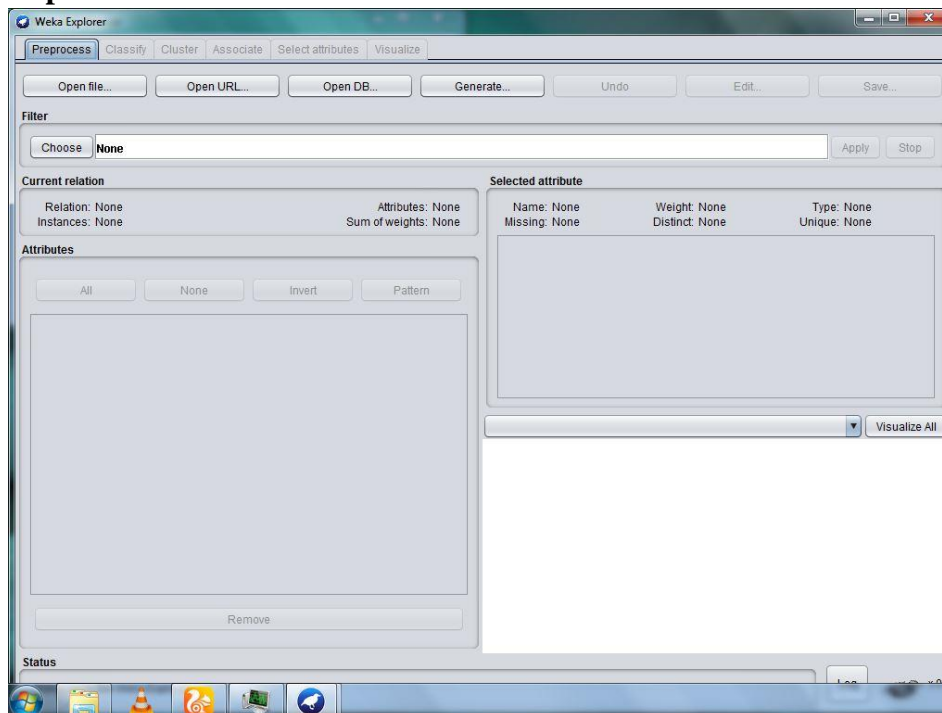# #Using Weka-tool for data pre-processing task

1. **Tool loading screen**
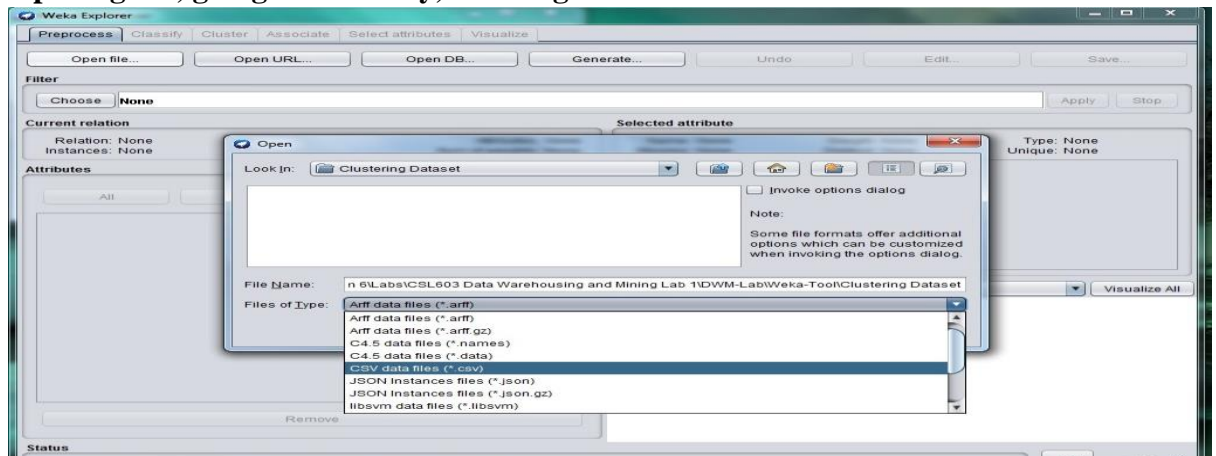


2. **Tool initial screen, selecting explorer**



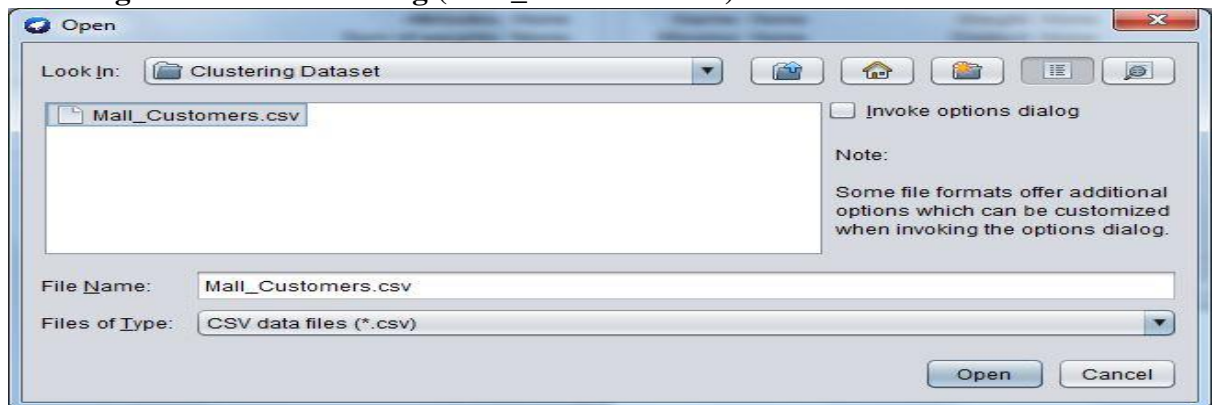3. **Explorer screen**

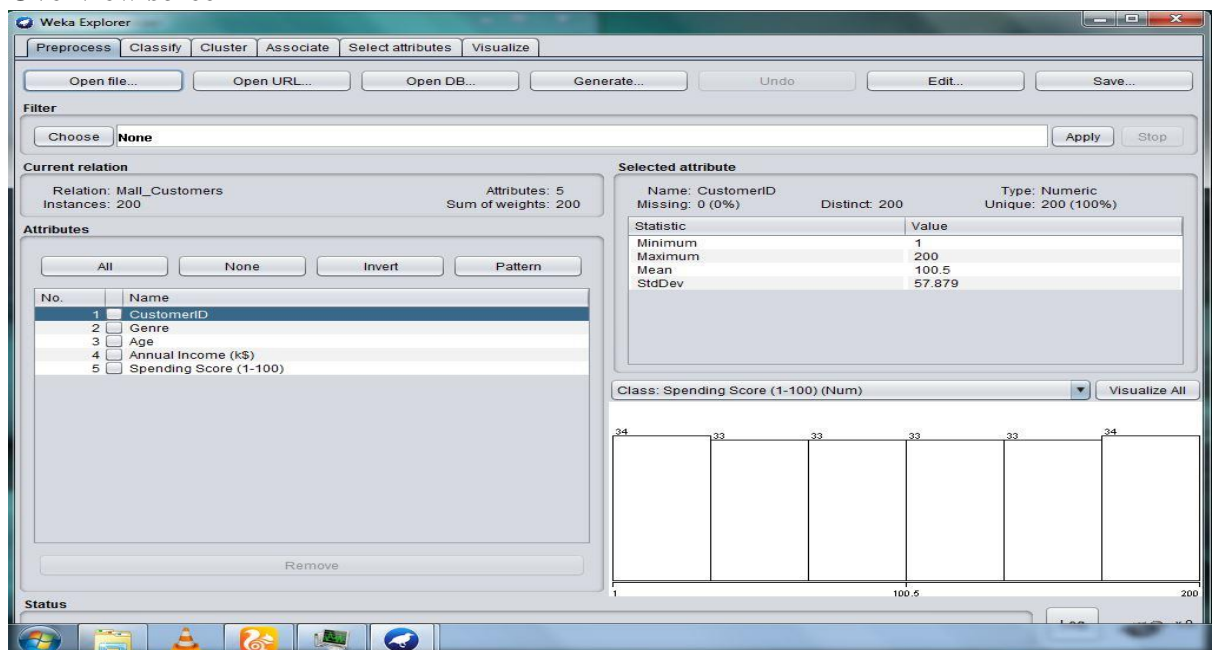# #Performing Clustering

1.  **Opening file, going to directory, selecting .csv format**
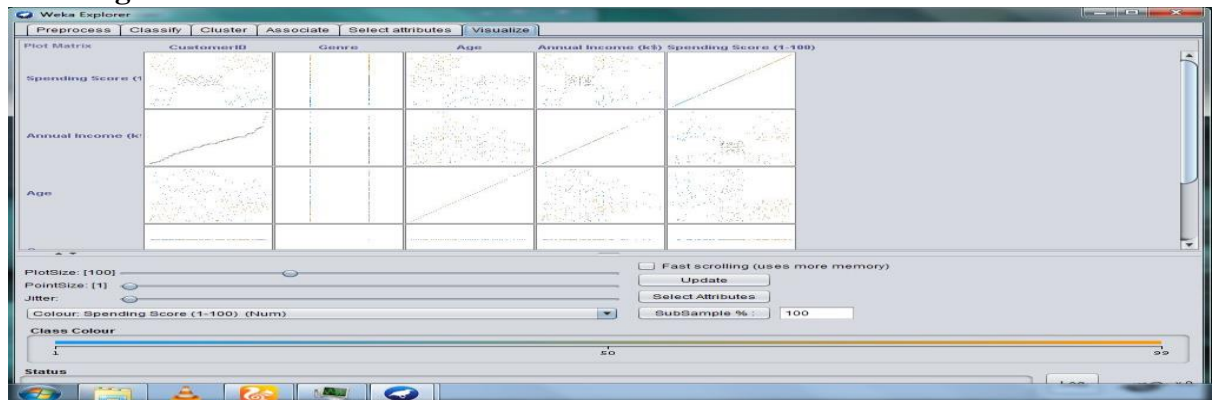
    

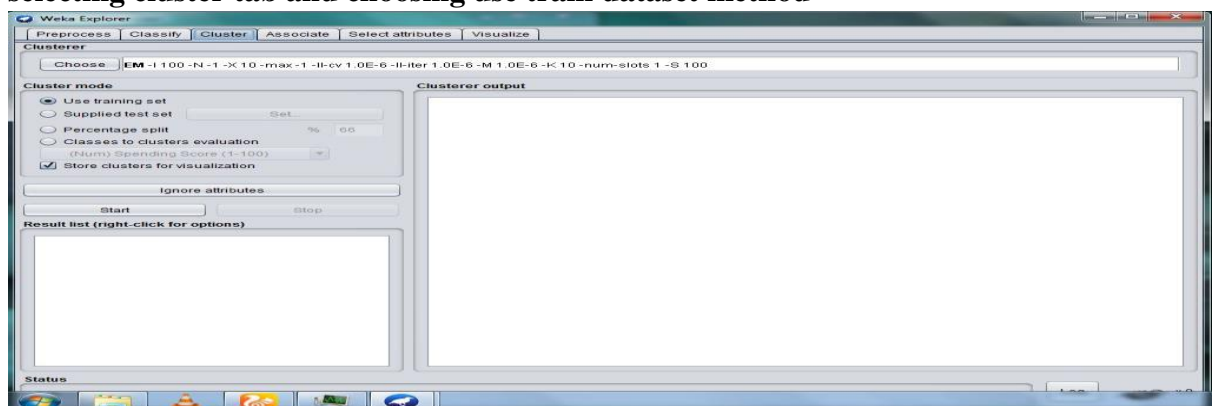2.  **Selecting dataset for clustering (Mall_Customers.csv)**
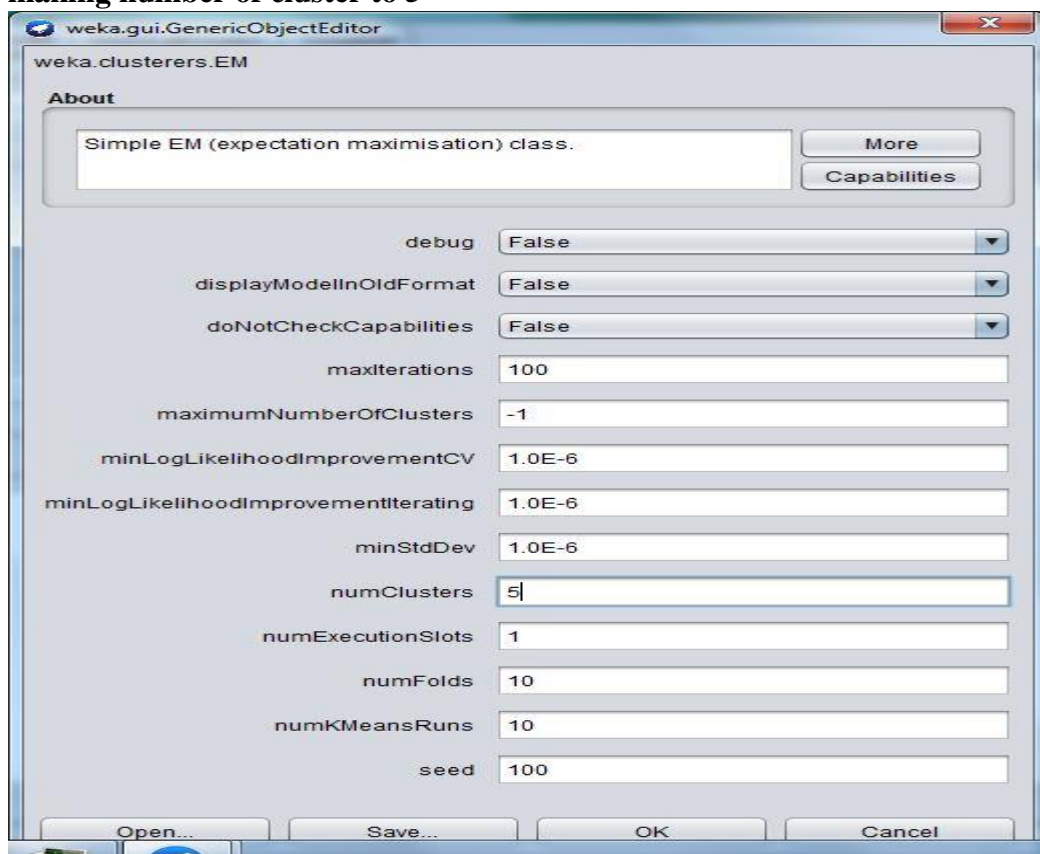
    

3.  **Overview screen**

    

**4. Selecting visualise tab**



**5. selecting cluster tab and choosing use train dataset method**



**6. making number of cluster to 5**

## 7. Removing unwanted attribute (CustomerID, Genre,Age)



## 8. right clicking on result buffer and clicking on visualise



## 9. Changing colour for better representation

**10. selecting annual income on x and spending on y and visualising**



# #Performing Classification

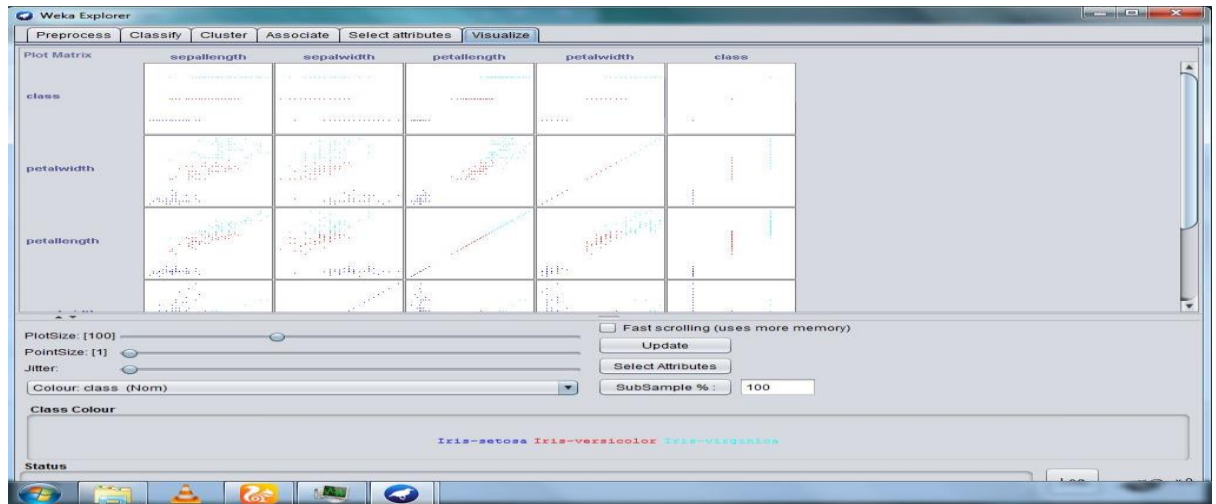1.  **opening file, going to directory, selecting .arff file format**



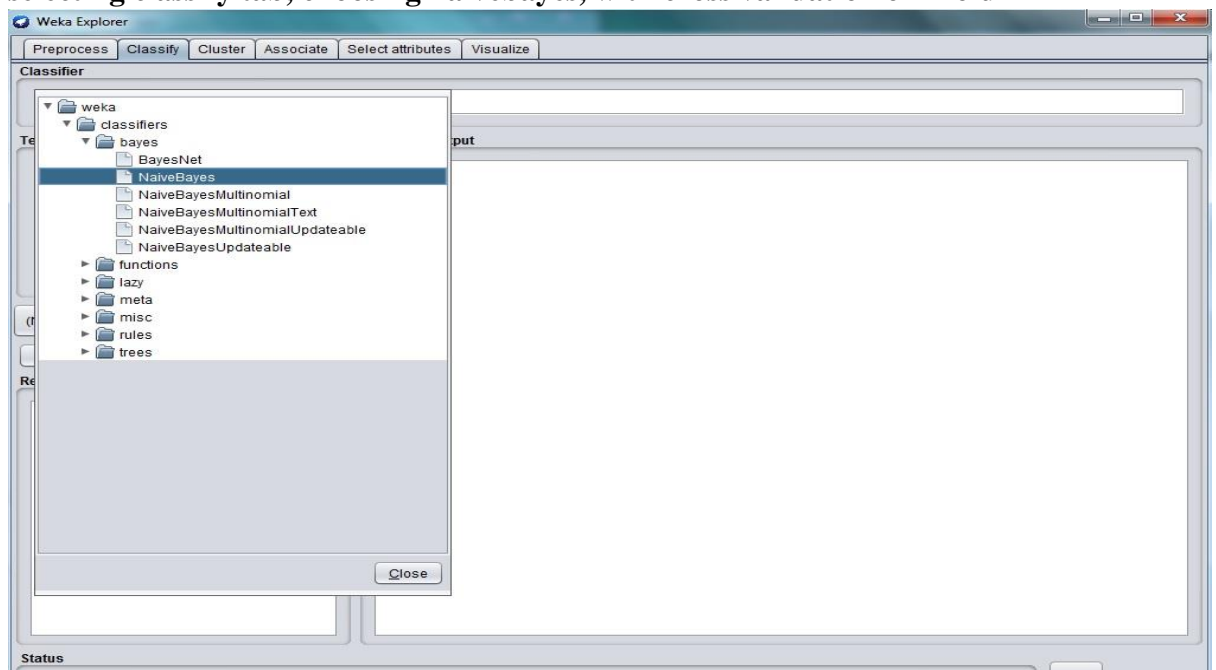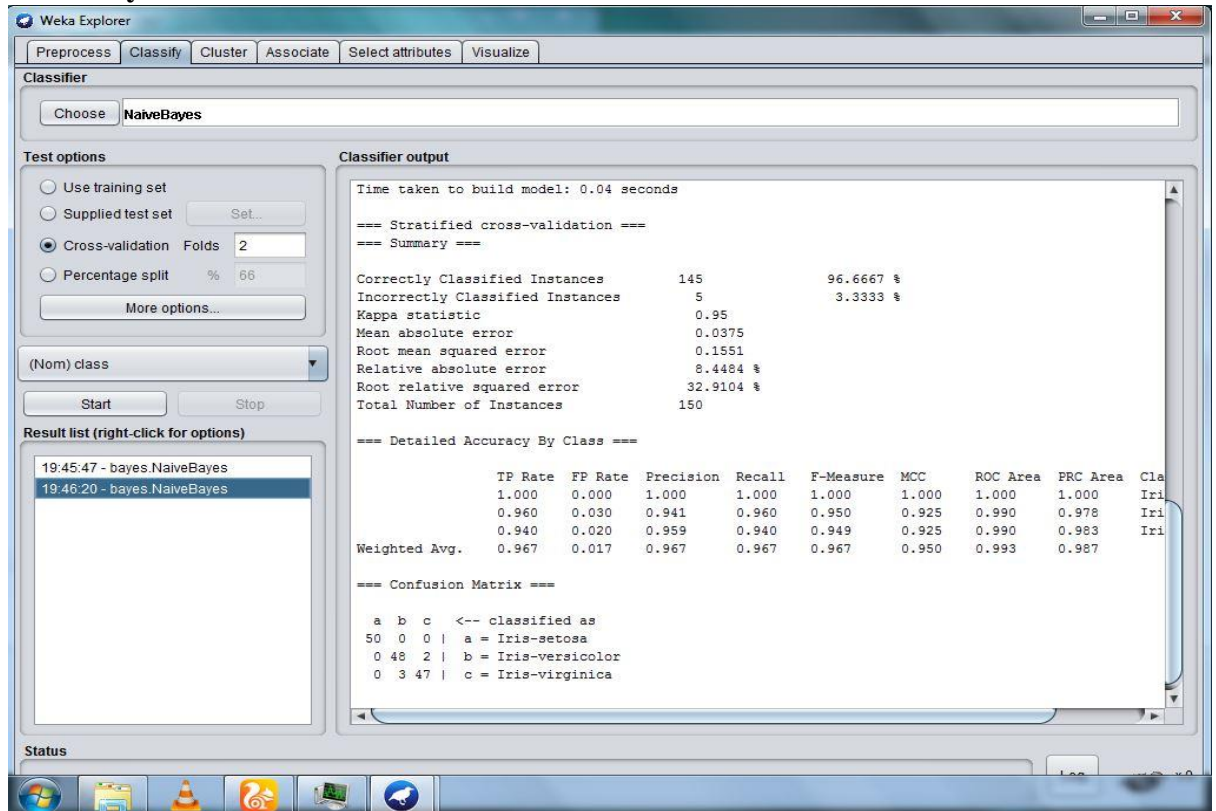2.  **selecting dataset for classification (iris.arff)**

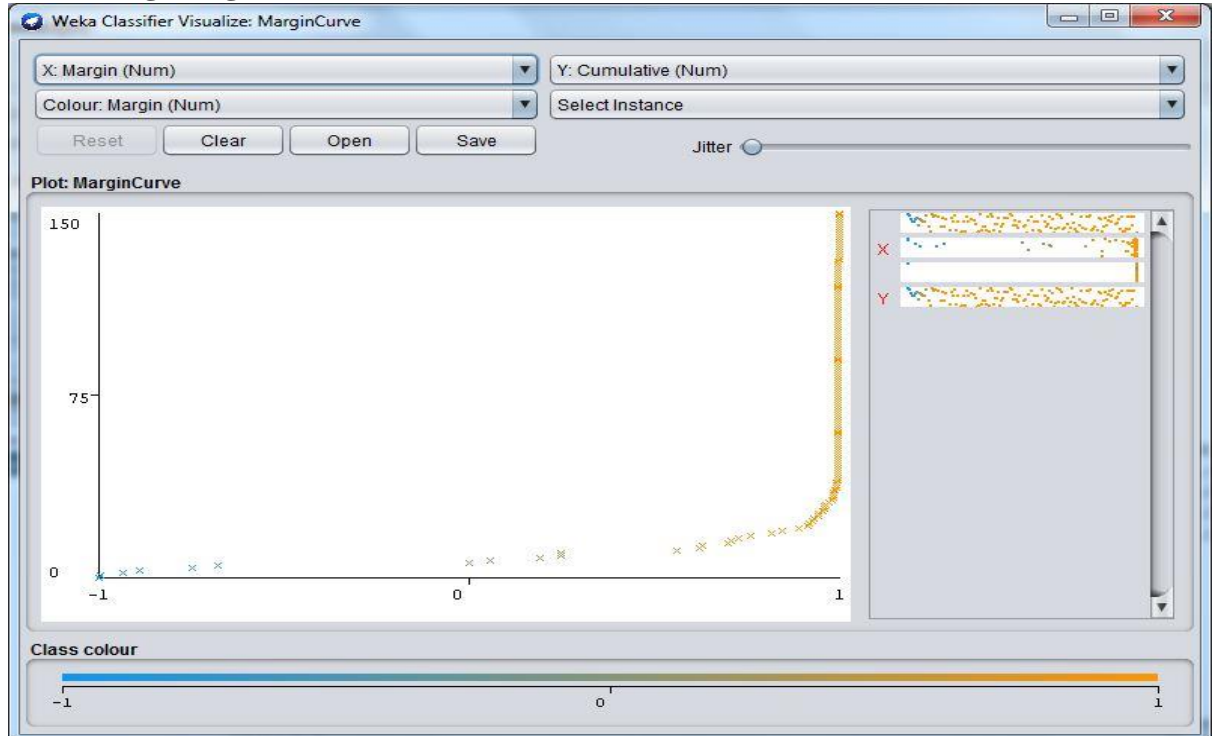**3. overview classification**



**4. Visualise tab screen**



**5. selecting classify tab, choosing naivebayes, with cross validation of 2 fold**
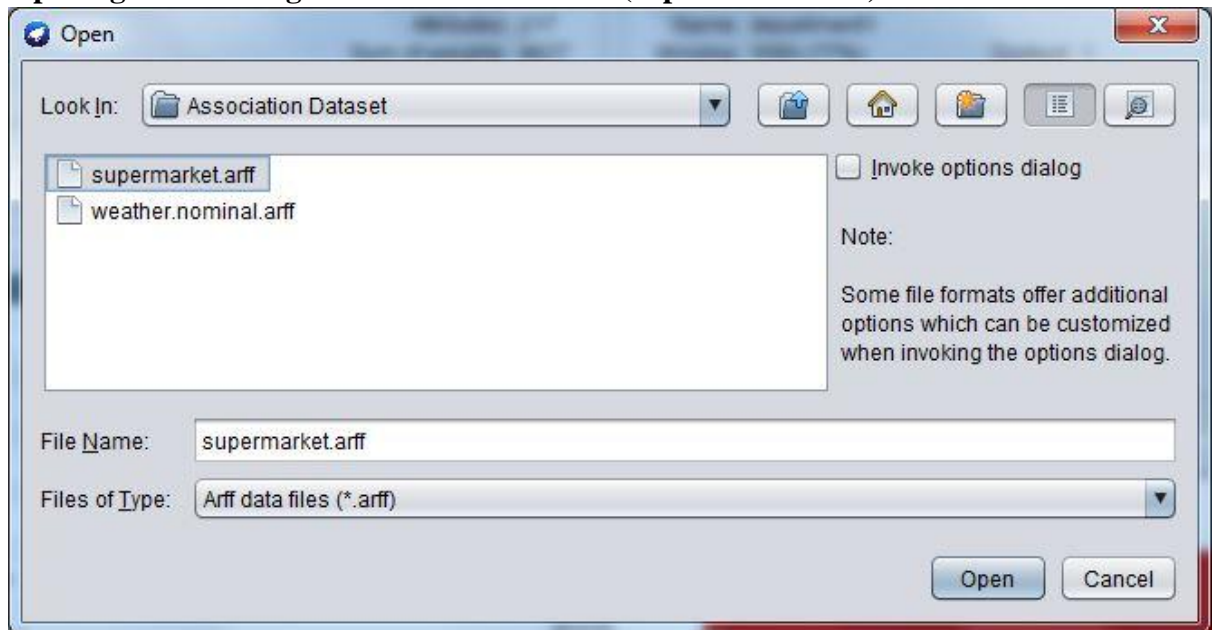
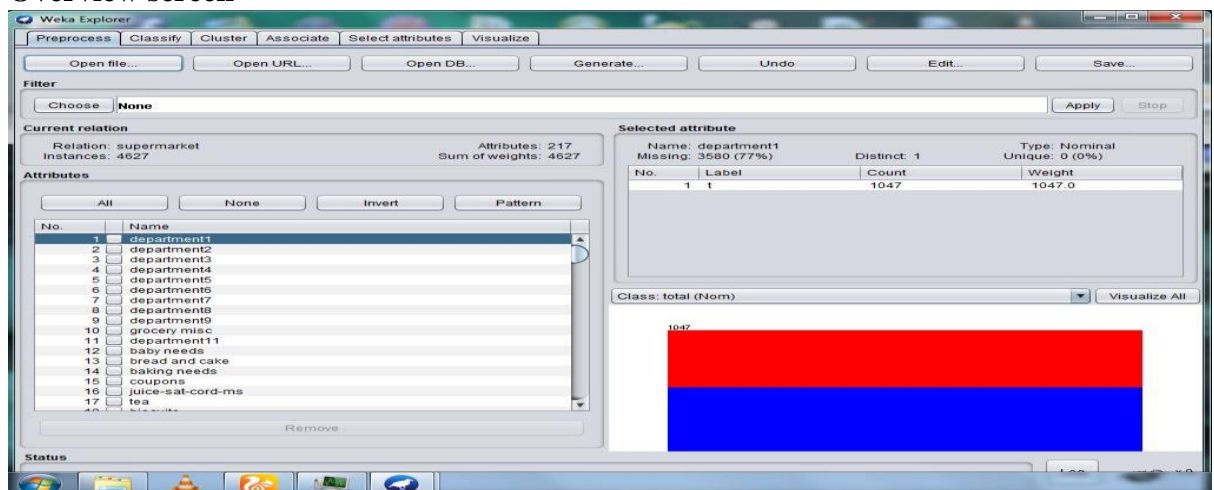## 6. Summary of result and confusion matrix



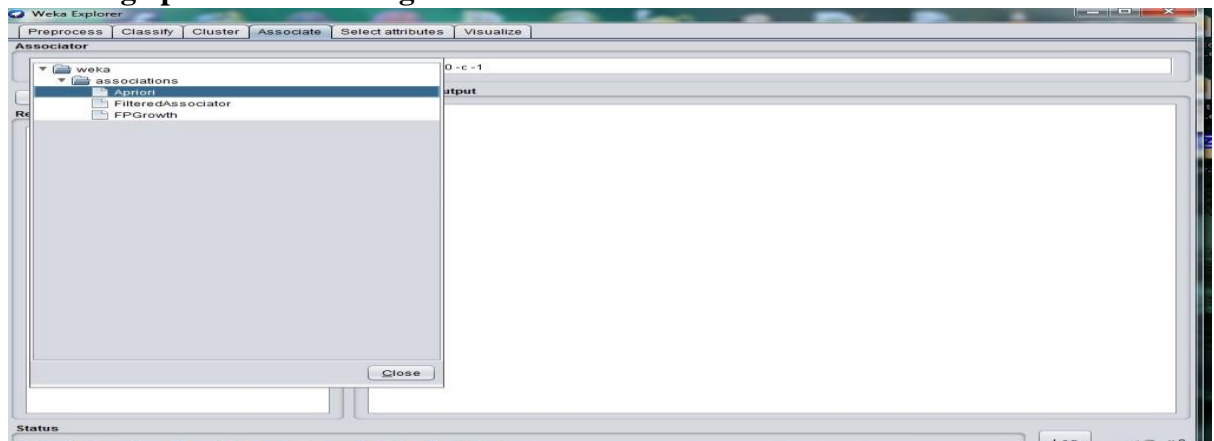## 7. Visualising margin curve

# #Performing Association

> **Opening file selecting dataset for association (supermarket.arff)**
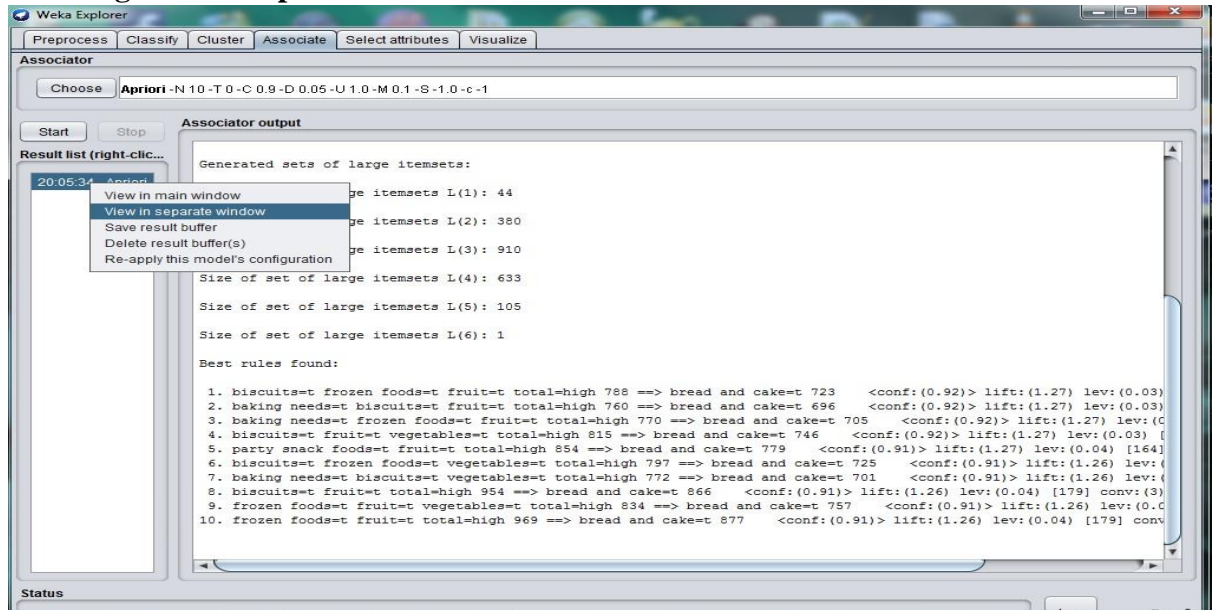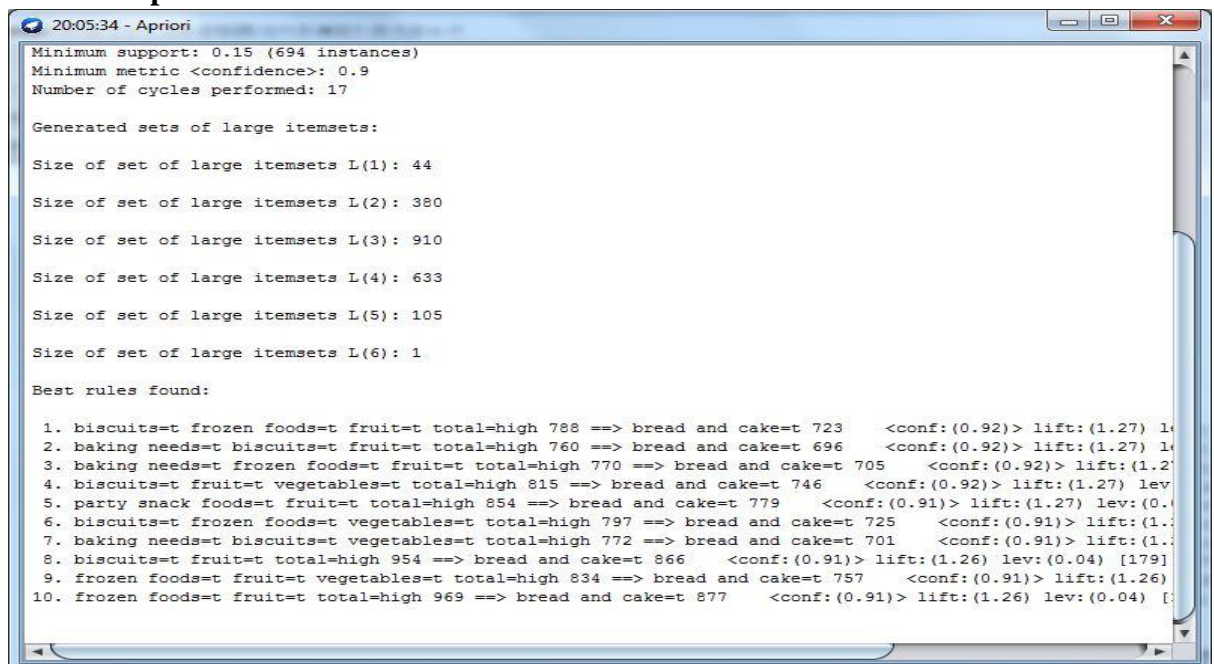


> **Overview screen**



> **Choosing apriori and starting**

➤ **Selecting View in separate window from context menu of result buffer**



➤ **Final rule pair of result**



# #Conclusion

From this experiment we learned some very essential concept related to any data mining process. For data pre-processing we used an open source weka tool which provides every result and ample of GUI based so the work task for data analysis becomes easier. Lastly, we implemented classification, clustering, and association using this tool and visualised it appropriately