# Analysis of supervised learning Classification Algorithms By using Weka 3.8

Name:Debobrata Chakraborty

ID:18-36103-1

Course: Data warehousing and Data Mining

Section:  B

Dept: Computer Engineering(CoE)

Email:

debobratachakraborty80@gmail.com

***Introduction—*** Heart disease is one of the most deadly and chronic illnesses, increasing the risk of heart failure and heart attack. The aim of this paper is to use supervised classification algorithms to classify heart disease. To detect heart disease, this experiment employs three machine learning classification algorithms: Naive Bayes, KNN, and Decision Tree.

Experiments are performed on the UCI machine learning repository's Heart Disease UCI dataset. Precision, Accuracy, F-Measure, and Recall are all used to test the performance of the three algorithms. The accuracy of a classification system is determined by the number of instances that are correctly classified and those that are incorrectly classified. According to the results, KNN outperforms with the highest precision of 70 percent. Those algorithms are also performed in a test dataset.

**Naive Bayes:** The Bayes Theorem is the basis for the Nave Bayes algorithm, which is used in a wide range of classification tasks.

**KNN:** Nearest Neighbour classification calculates the classification of an unseen instance using the classification of the instances "closest" to it, and Nave Bayes uses probability theory to find the most probable of the possible classifications. Since it does not learn from the training set right away, the K-Nearest Neighbor lazy learner algorithm stores the dataset and performs an operation on it at the time of classification.

**Decision Tree:** The Decision Tree is a supervised machine learning algorithm that solves classification problems using an intermediate tree-like structure. A decision tree is a diagram that depicts the various outcomes of a set of similar choices. It enables a person or organization to compare and contrast various options based on their costs, probabilities, and benefits. They can be used to ignite informal discussion or to create an algorithm that mathematically predicts the best option.

# I.    Result

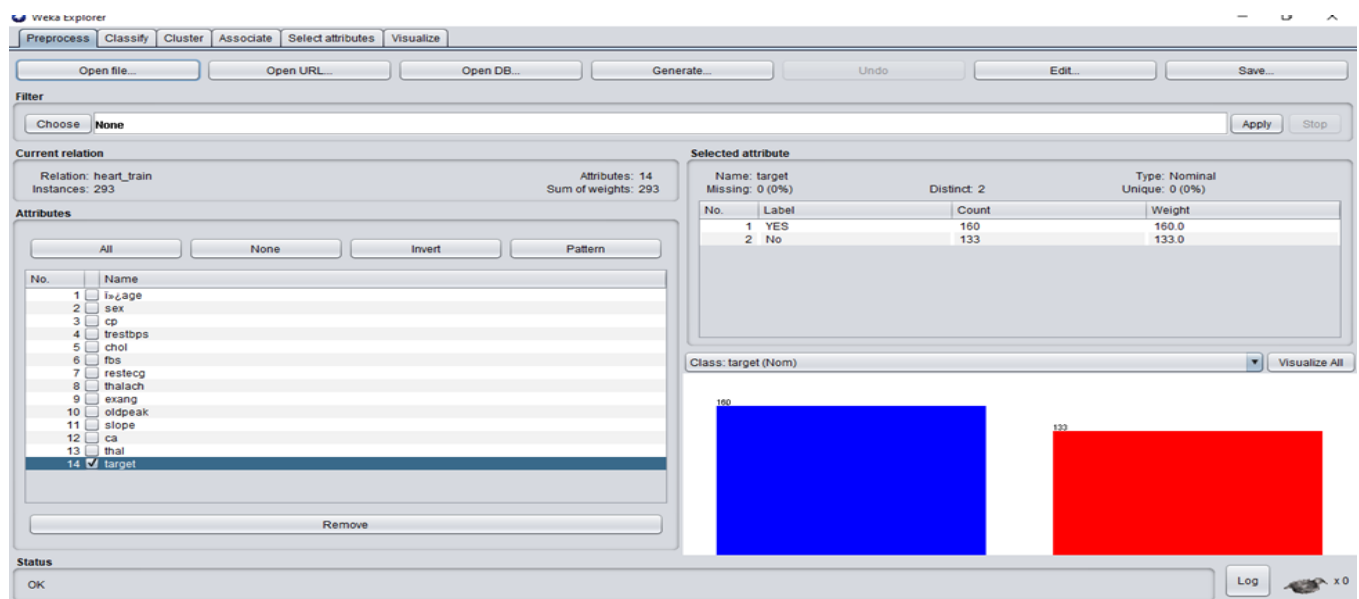At first, the heart disease dataset was loaded in Weka.



Figure 1: Preprocess the data

There are 293 instances in the dataset. Target indicates 2 classes. Those who have heart disease are labeled as 'Yes,' and those who do not have heart disease are labeled as 'No.' There are 14 attributes in a dataset. The classify algorithm was introduced after the data was pre-processed. Test dataset contain 10 instance.

For each classification algorithm, a 10-fold cross validation set was used.

Naive  Bayes: After applying naïve bayes algorithm on the training dataset  it classified 246 correctly and 47 incorrectly.
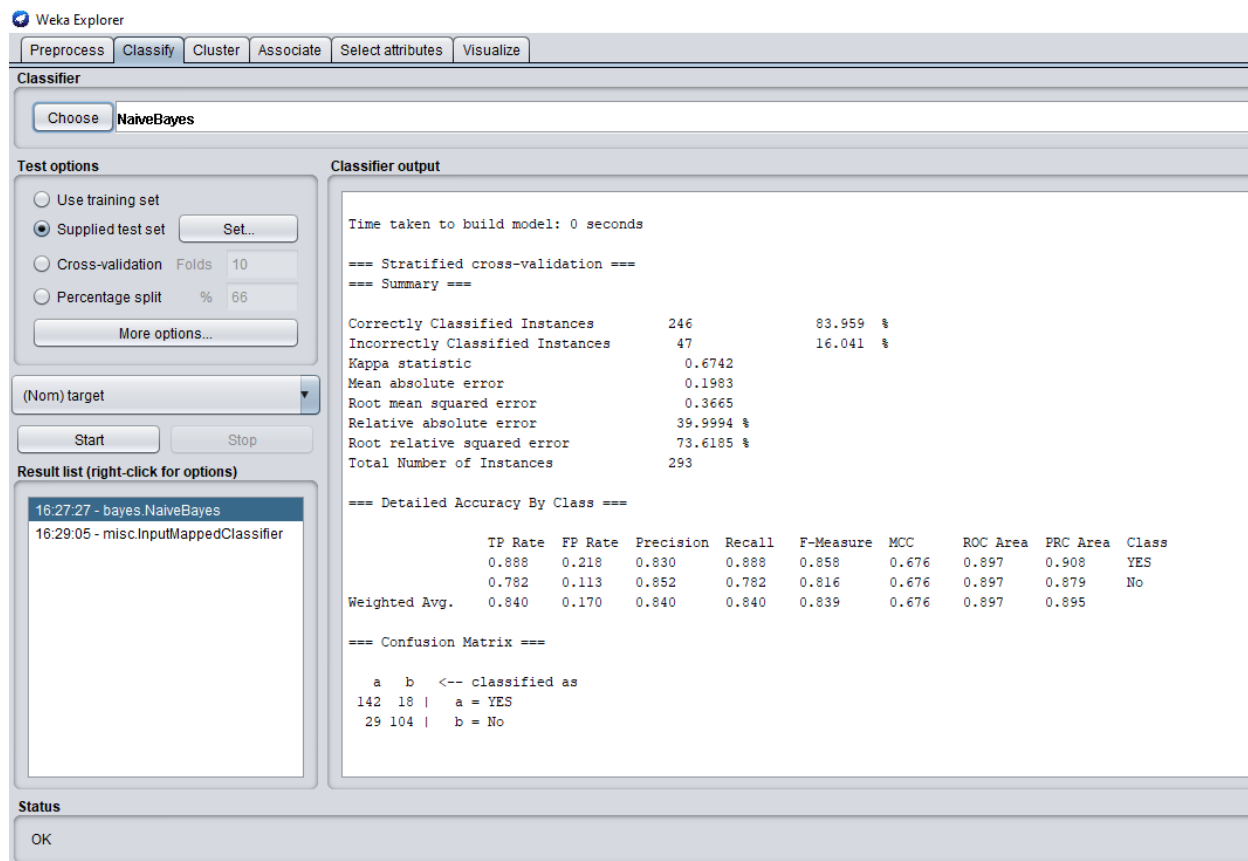


Figure 2: Naïve Bayes classification (training set)

**Confusion matrix:** A confusion matrix is a method of summarizing a classification algorithm's results.

Table 1:Confusion matrix of Naïve Bayes(training set):

| a | b | < classified as |
|---|---|---|
| 142 | 18 | a=Yes |
| 29 | 104 | b=No |

Table 2:Counfusion matrix of Naïve Bayes (test set)

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

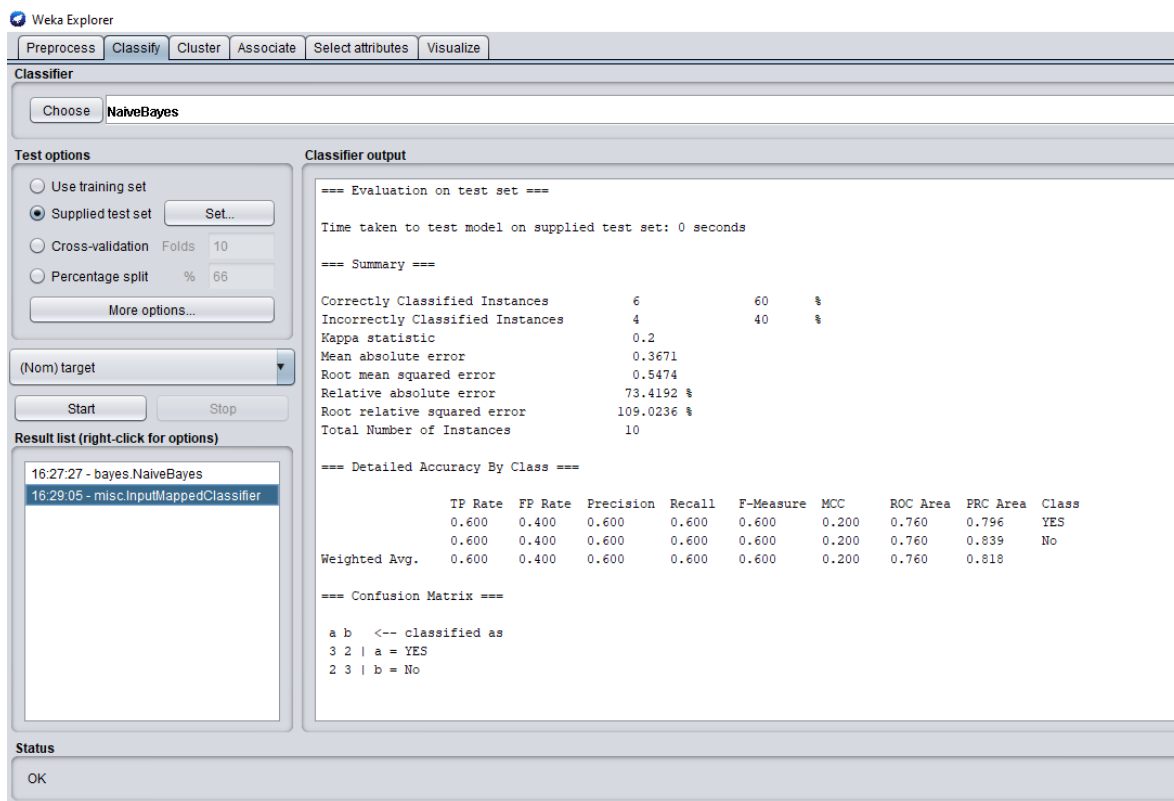Naïve bayes  correctly classified 6 instances from the test set.

Figure 3: Naïve Bayes classification (test set)

**KNN Algorithm:** To apply the knn algorithm, Lazy IBK was chosen from the classify option in weka.
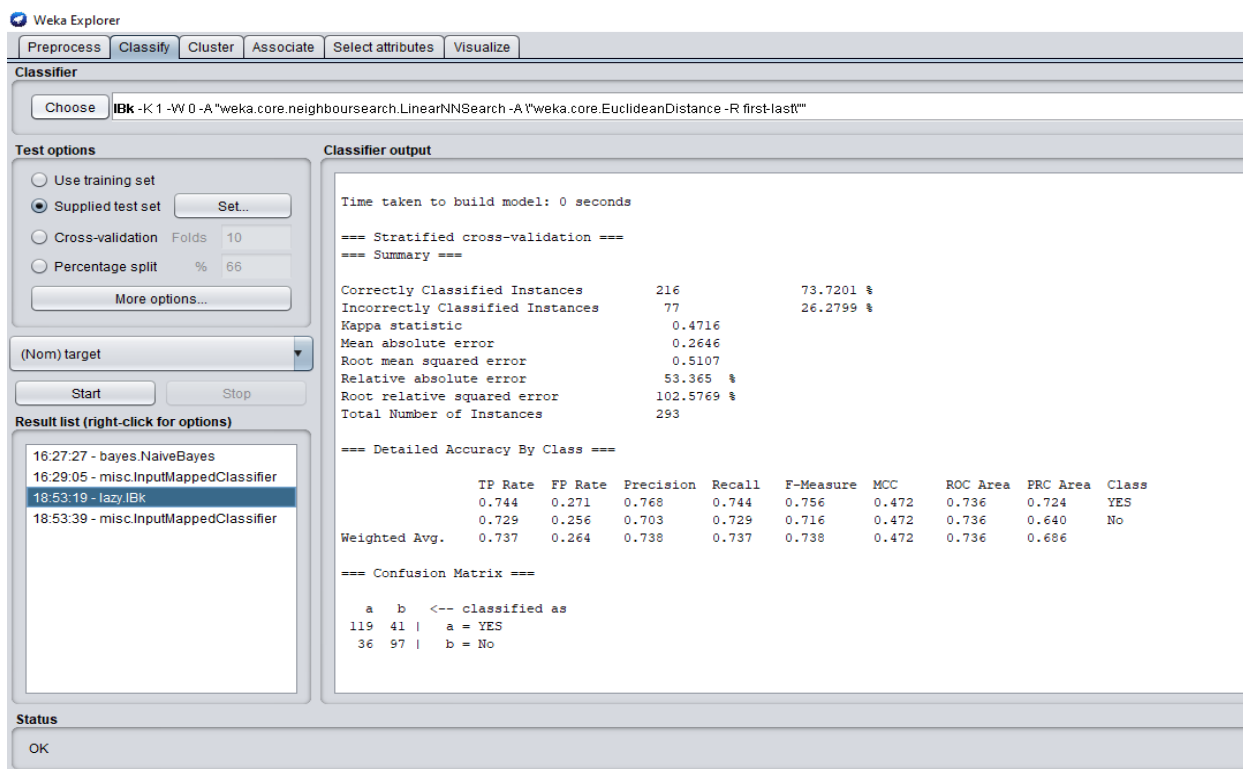


Figure 4: KNN classification (training set)

Table 3:Confusion matrix of KNN(training set):

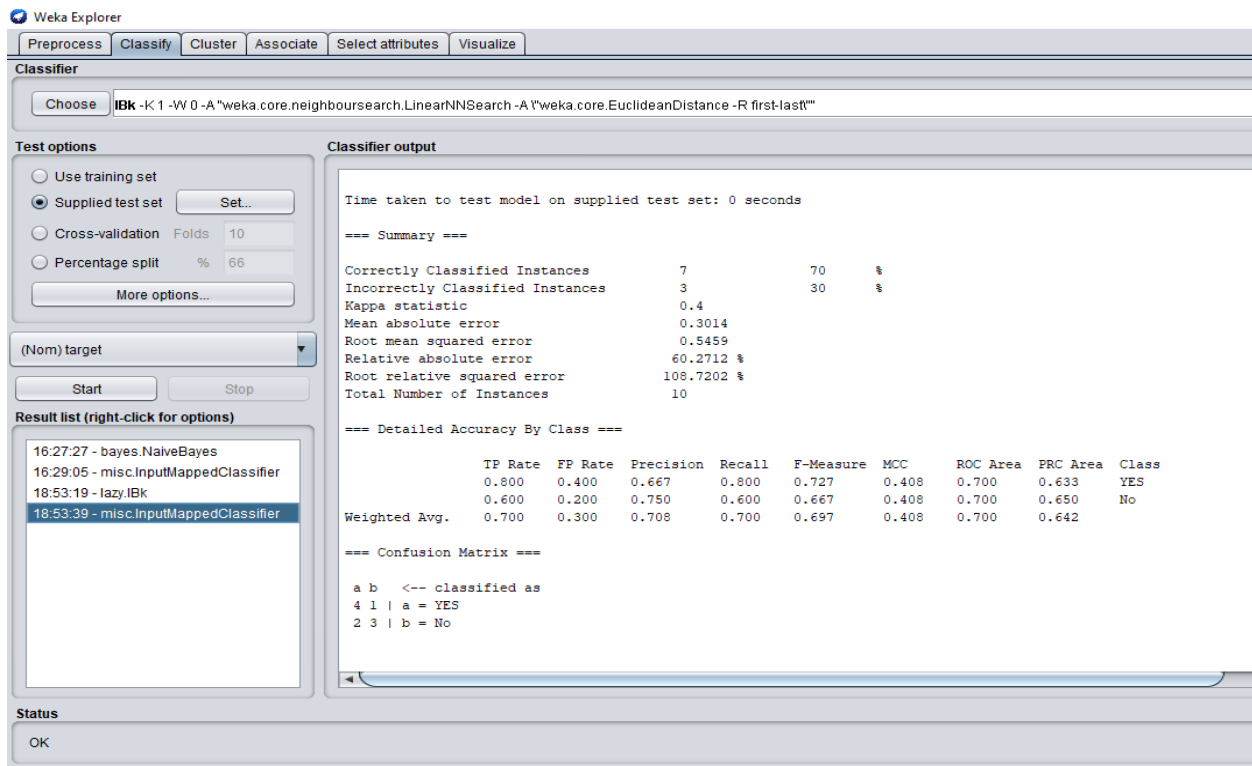| a | b | < classified as |
|---|---|---|
| 119 | 41 | a=Yes |
| 36 | 97 | b=No |

Figure 4: KNN classification (test set)

Table 4:Confusion matrix of KNN(Test set):

| a | b | < classified as |
|---|---|---|
| 4 | 1 | a=Yes |
| 2 | 3 | b=No |

KNN correctly classified 7 instances from the test set.

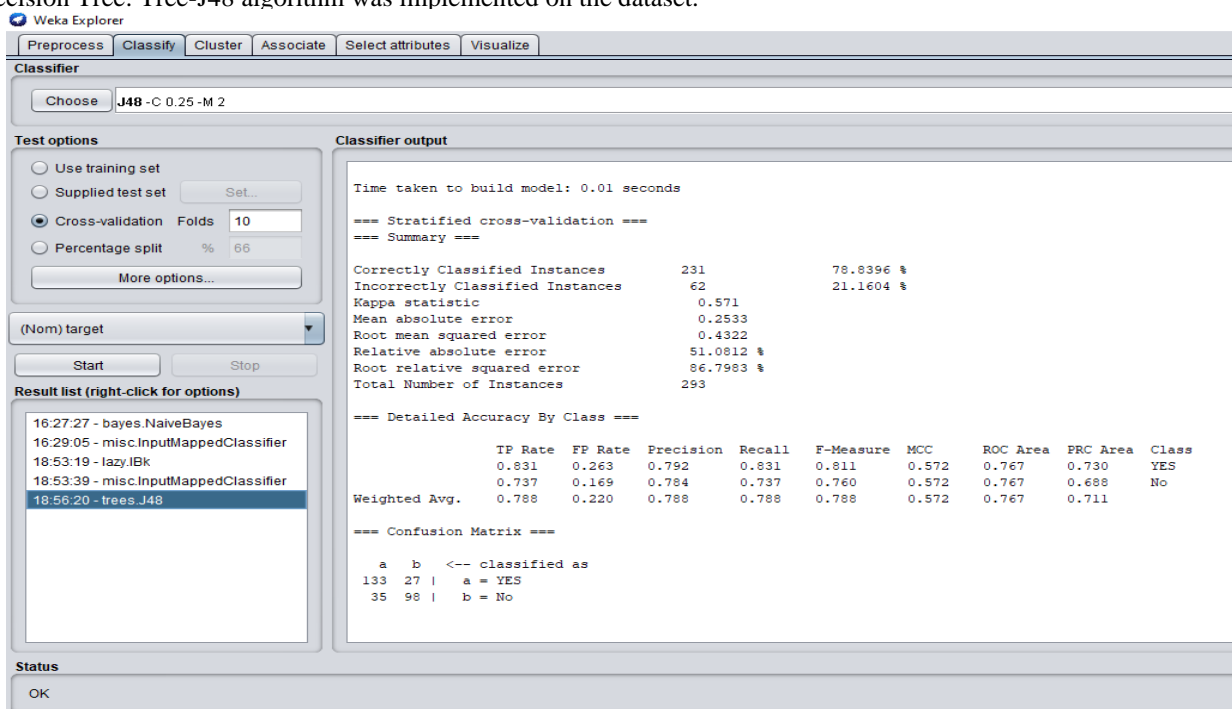Decision Tree: Tree-J48 algorithm was implemented on the dataset.



Figure 5: Decision Tree classification (training set)

Table 5:Confusion matrix of Decision Tree(training set):

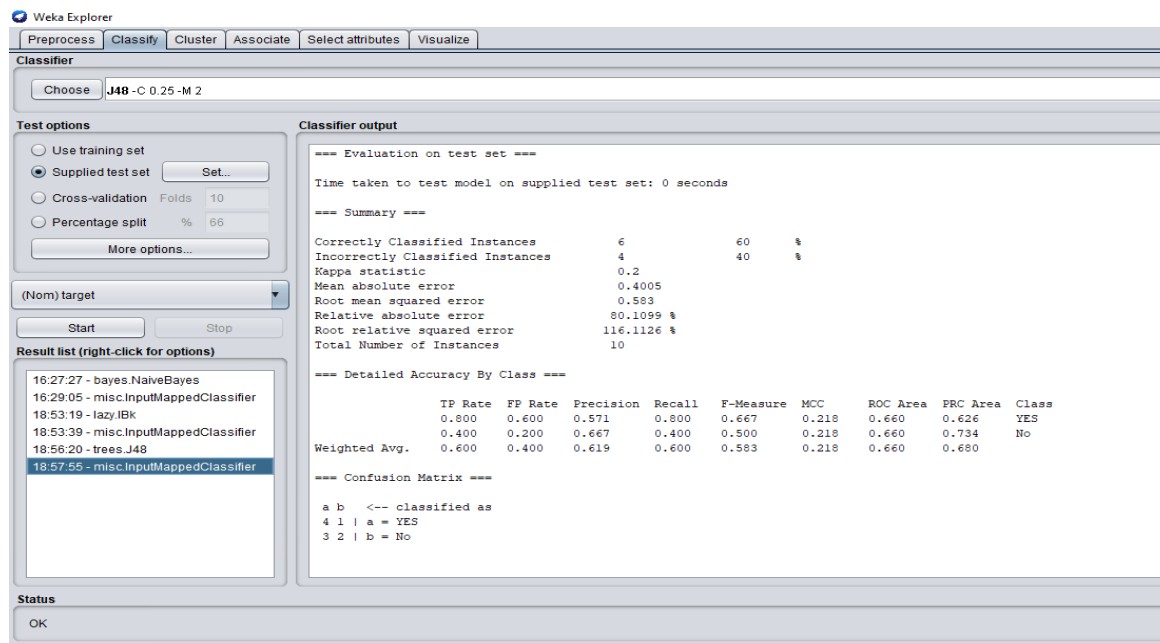| a | b | < classified as |
|---|---|---|
| 133 | 27 | a=Yes |
| 35 | 98 | b=No |



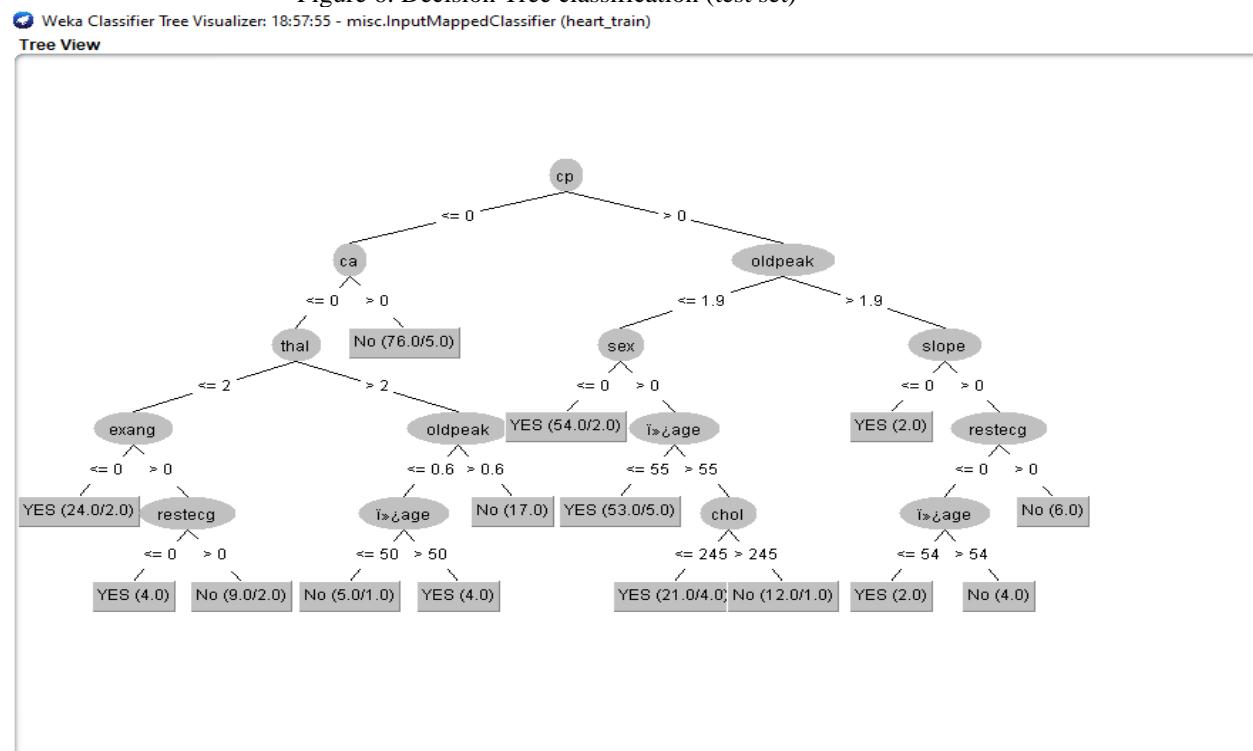Figure 6: Decision Tree classification (test set)



Figure 7: Tree

A decision tree usually begins with a single node and branches out into different outcomes. Each of those outcomes leads to new nodes, each of which leads to new possibilities. It takes on a tree-like form as a result of this.

Table 6:Confusion matrix of Decision Tree(training set):

| a | b | < classified as |
|---|---|---|
| 4 | 1 | a=Yes |
| 3 | 2 | b=No |

Table 7:Accuracy Measure

| Measures | Definitions | Formula |
|---|---|---|
| 1. Accuracy (A) | Accuracy determines the accuracy of the algorithm in predicting instances. | $A=(TP+TN) / \text{(Total no of samples)}$ |
| 2. Precision (P) | Classifiers correctness/accuracy is measured by Precision. | $P = TP / (TP+ FP)$ |
| 3. Recall (R) | To measure the classifiers completeness or sensitivity, Recall is used. | $R =TP / (TP+FN)$ |
| 4. F-Measure | F-Measure is the weighted average of precision and recall. | $F=2*(P*R) / (P+R)$ |
| 5. ROC | ROC(Receiver Operating Curve) curves are used to compare the usefulness of tests. | |

Test dataset was considered to find out the high accuracy classification.

Table 8: Comparative Performance of Classification Algorithms on Various Measures.

| Algorithm | Precision | Recall | F-Measure | ROC | Accuracy(%) |
|---|---|---|---|---|---|
| Naïve Bayes | 0.600 | 0.600 | 0.600 | 0.760 | 60 |
| KNN | 0.708 | 0.700 | 0.697 | 0.700 | 70 |
| Decision Tree | 0.619 | 0.600 | 0.583 | 0.660 | 60 |

Table-7 shows the performance of corresponding classifiers in terms of Accuracy, Precision, F-measure, Recall, and ROC values, while Table-8 shows the performance of classifiers based on classified instances. TP stands for True Positive, TN for True Negative, FP for False Positive, and FN for False Negative. The output of the corresponding classifiers is measured in terms of accuracy, precision, and recall.

After applying the model on test dataset the below data was found.

Table 9. Classifier's Performance on The Basis of Classified Instances

| Total Instance | Algorithm | Correctly Classified Instance | Incorrectly Classified Instances |
|---|---|---|---|
| 10 | Naïve Bayes | 6 | 4 |
| | KNN | 7 | 3 |
| | Decision Tree | 6 | 4 |

As a result, the KNN algorithm is the best supervised machine learning approach in this experiment because it has a higher accuracy than other classification algorithms, with a 70.00 percent accuracy.

## II.     Discussion

Three supervised machine learning algorithms were used in this experiment. After calculating all of the values, the KNN value has a maximum accuracy of 70%. On the evaluation dataset, it performs admirably.  This machine learning approach was applied to a dataset of heart diseases. This algorithm can also be used on other datasets to build a classifier model that can predict various diseases.

## III.     Reference

[1] https://www.kaggle.com/ronitf/heart-disease-uci

# Analysis of unsupervised learning Clustering Algorithm By using Weka 3.8

Name:Debobrata Chakraborty

ID:18-36103-1

Course: Data warehousing and Data Mining

Section:  B

Dept: Computer Engineering(CoE)

Email:

debobratachakraborty80@gmail.com

***Introduction-***  Diabetes is one of the most deadly and chronic diseases that causes blood sugar levels to rise. The aim of this paper is to use unsupervised learning algorithms to clustering diabetes disease. To detect diabetes, this experiment employs hierarchical cluster.

Experiments are performed on the Prima Indian Diabetes dataset. Clustering is a descriptive model that divides a set of objects into groups based on their relationships.

 Clustering is a method that is used in a variety of areas.Image analysis, pattern recognition, statistical data analysis, and so on are only a few examples. Clustering is the division of data into groups of items that are identical.  Clustering is the division of data into groups of items that are identical. Each cluster is made up of a variety of objects that are identical to one another but not to objects from other classes. To shape clusters, various clustering algorithms are available. To compare various clustering algorithms, the WEKA tool is used.

In this paper, hierarchical cluster is used,  an algorithm that groups related objects into clusters, also known as hierarchical cluster analysis. The endpoint is a set of clusters, each of which is distinct from the others while the artifacts within each cluster are broadly identical.

## I.      Result

At first, the diabetes test dataset was loaded in Weka. There are 223 instances and 9 attribute on the diabetes dataset.
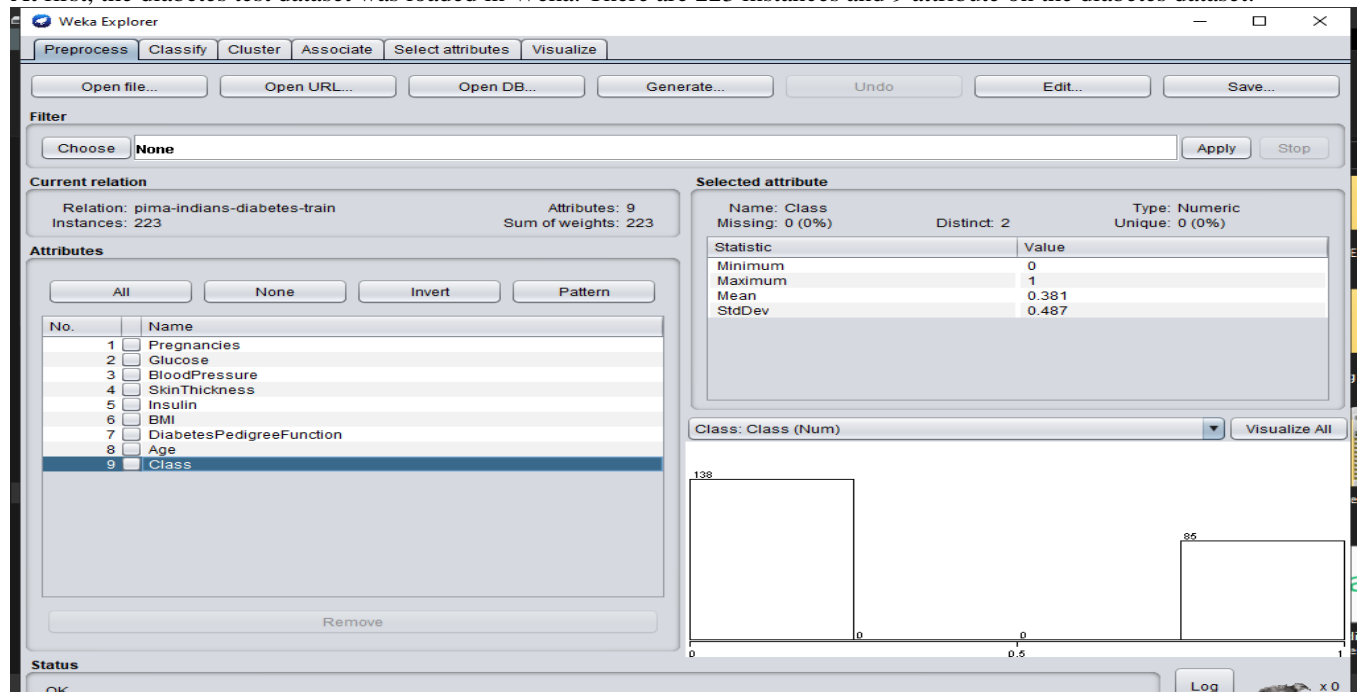


Figure 1: Preprocess the data

After that, the cluster option was selected and heirarchical clustering algorithm was implemented.

Figure 2: Clustering(training dataset)

After the clustering, the training set clustering into two part.

| 0 | 138 instance | 62% |
|---|---|---|
| 1 | 85 instance | 38% |

138 instance clustered into 0 category and 85 instances clustered into 1 category.



Figure 3: Clustering (test set)

By supplying the test set the model clustered the data. 179 instances(65%) are clustered into category 0 and 96 instances(35%) are clustered into category 1.
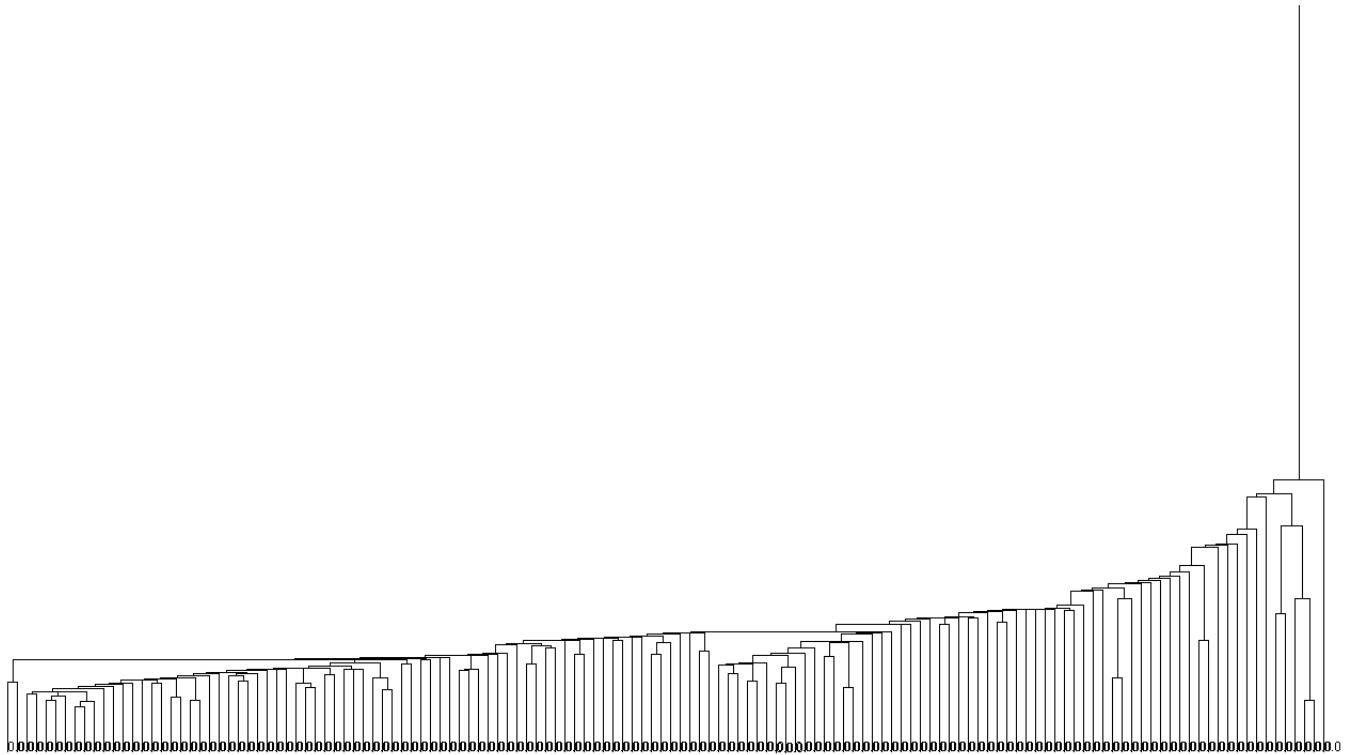


Figure 4: Dendrogram

This is the dendrogram of the hierarchical cluster. A dendrogram is a type of tree diagram that is commonly used to show how hierarchical clustering produces clusters. In computational biology, dendrograms are often used to depict the clustering of genes or samples.
The following table represents the analysis process of the clustering on test set.

| Algorithm | No. of cluster | Cluster instance | Time taken to build |
|---|---|---|---|
| hierarchical cluster | 2 | . 179 (65%) 96(35%) | 0.2s |

## II.    Discussion

Data mining is a computer science and information technology discipline. There is a vast amount of data scattered across the globe, much of which is in the form of raw data, which often includes relevant information. A mining method is used to obtain relevant data. A cluster is formed from a large amount of data by objects of a similar nature. Hierarchical cluster was performed in this paper. The algorithm generates clusters based on related objects and the amount of time it takes to construct them

## III.    Reference

[1]https://www.kaggle.com/kumargh/pimaindiansdiabetescsv