# DATA WAREHOSING & DATA MINING

## SEMESTER: Spring 2021-2022

## FINAL-TERM ASSIGNMENT

## DATA MINING WITH WEKA

## SUBMITTED BY:

| STUDENT NAME | STUDENT ID |
|---|---|
| KHANAM, FAHMIDA | 19-40351-1 |

## SUBMITTED TO:

## COURSE TEACHER: TOHEDUL ISLAM

**INTRODUCTION:**

Data mining is the process of uncovering patterns and other valuable information from large data sets. It is also known as knowledge discovery in data (KDD). Data mining is used in many areas of research and business, including healthcare, education, sales and marketing, product development, etc. It is a computer science and statistics multidisciplinary topic with the general purpose of extracting information from data collection and transforming the information into an accessible structure for subsequent use. KNN, Naive Bayes, and Decision Tree are some of the classification algorithms used in data mining. I have chosen the "Heart-Diseases dataset" to classify the type by using two different classifiers and find the best-suited classifier for the dataset. I have chosen Naive Bayes and Decision Tree to classify the dataset.

**Information about the dataset:** In this report, the used "Heart-failure dataset", a CSV dataset file [ Later converted into .arff ], collected from Kaggle.com was used to predict the outcome of the heart failure that might be accurate for the patient according to their health condition.
The targeted feature is:
- HeartDisease

The other feature sets are:
- Age
- Sex
- ChestPainType
- RestingBP
- Cholesterol
- FastingBS
- RestingECG
- MaxHR
- ExerciseAngina
- Oldpeak
- ST_Slope

There is a total of 918 instances of these 12 attributes and all these instances were used for classification. Here are the graphical details of the attributes:
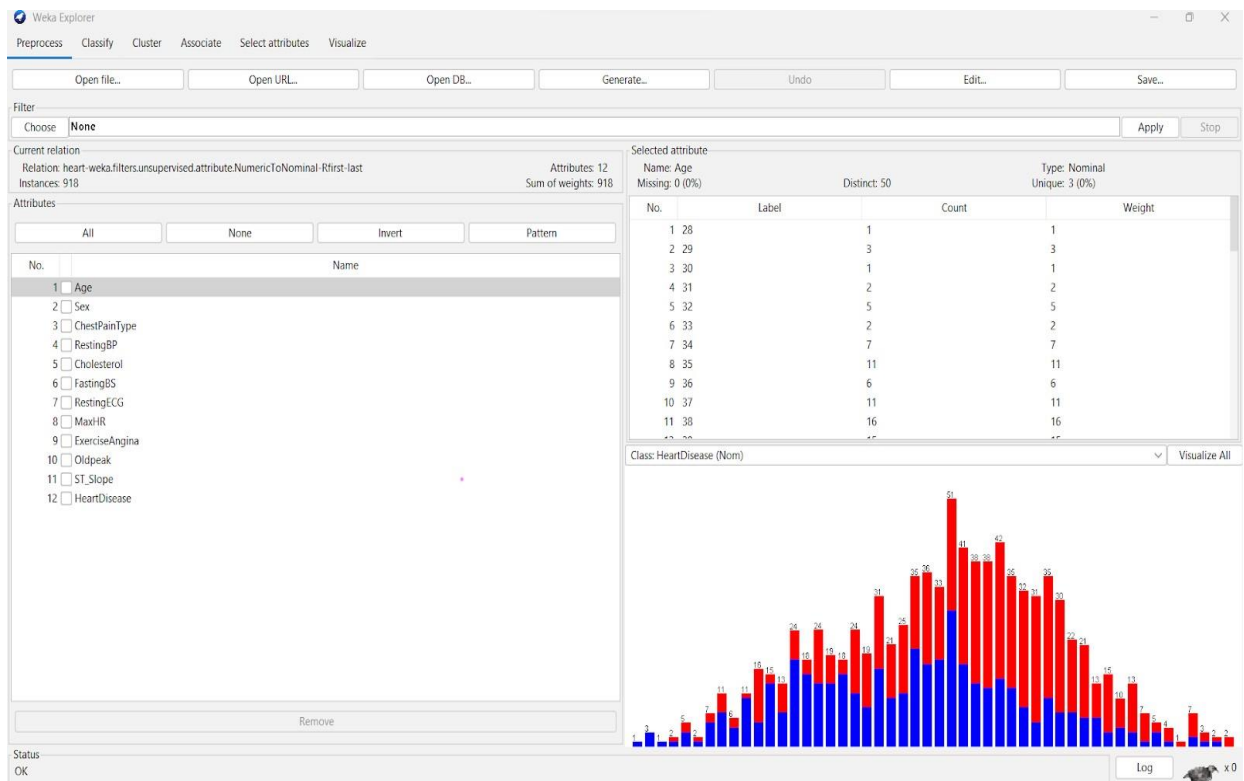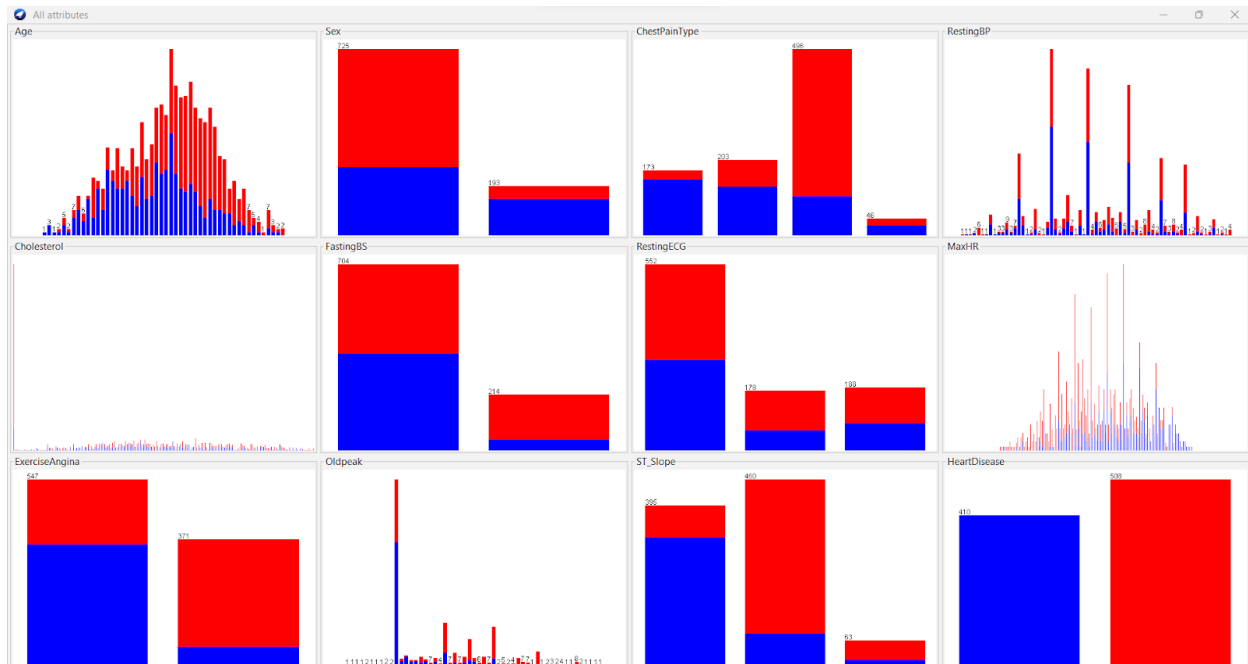


**Figure 1: Selected dataset**

**Figure 2: Details of all attribute**

**Classifier:** A classifier is a machine learning model that is used to discriminate different objects based on certain features. Two kinds of classification have been used with the same data to compare the result. In this process, Naïve Bayes, and Decision Tree Classifiers were used.

## RESULT OF THE CLASSIFIERS:

Weka 3.8.6 version software was used to construct the classifier. Weka is a collection of machine learning algorithms for data mining tasks. Weka contains tools for data pre-processing, classification, regression, clustering, association rules, and visualization.

**Applying naïve Bayes classifier**: Naïve Bayes algorithm is a supervised learning algorithm, which is based on the Bayes theorem and is used for solving classification problems. It is mainly used in text classification that includes a high-dimensional training dataset. It is a probabilistic classifier, which means it predicts on the basis of the probability of an object. While classifying the selected dataset, the NaiveBayes format was selected from the Bayes folder.
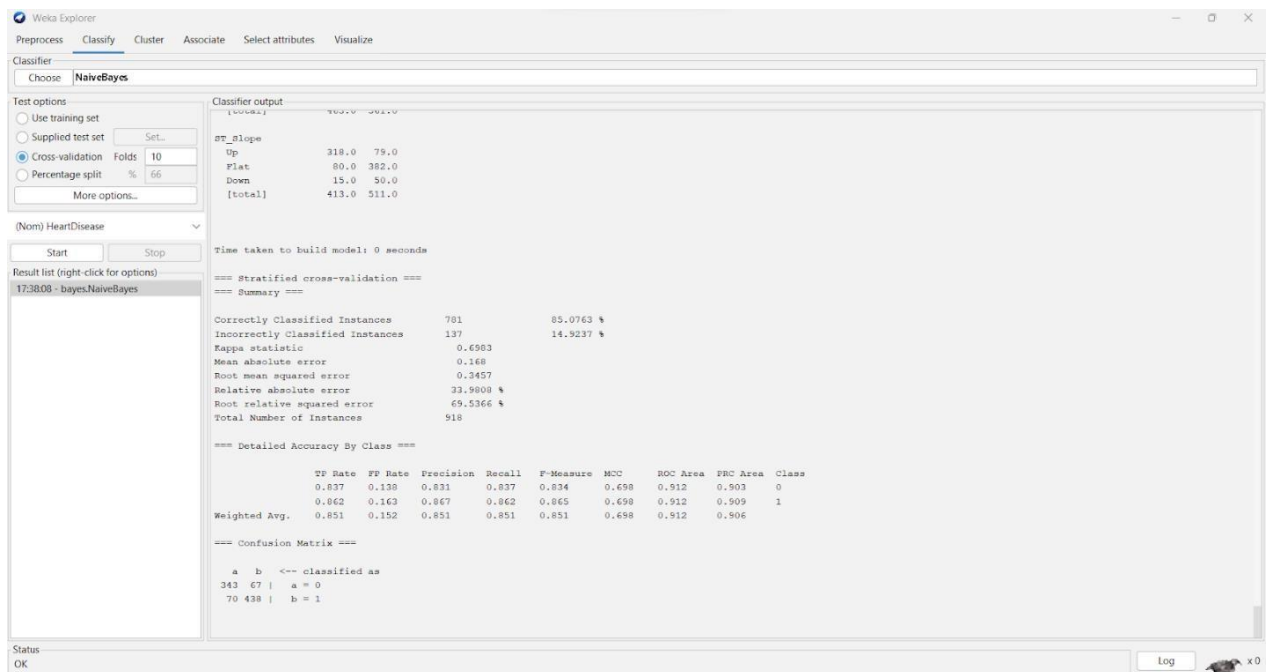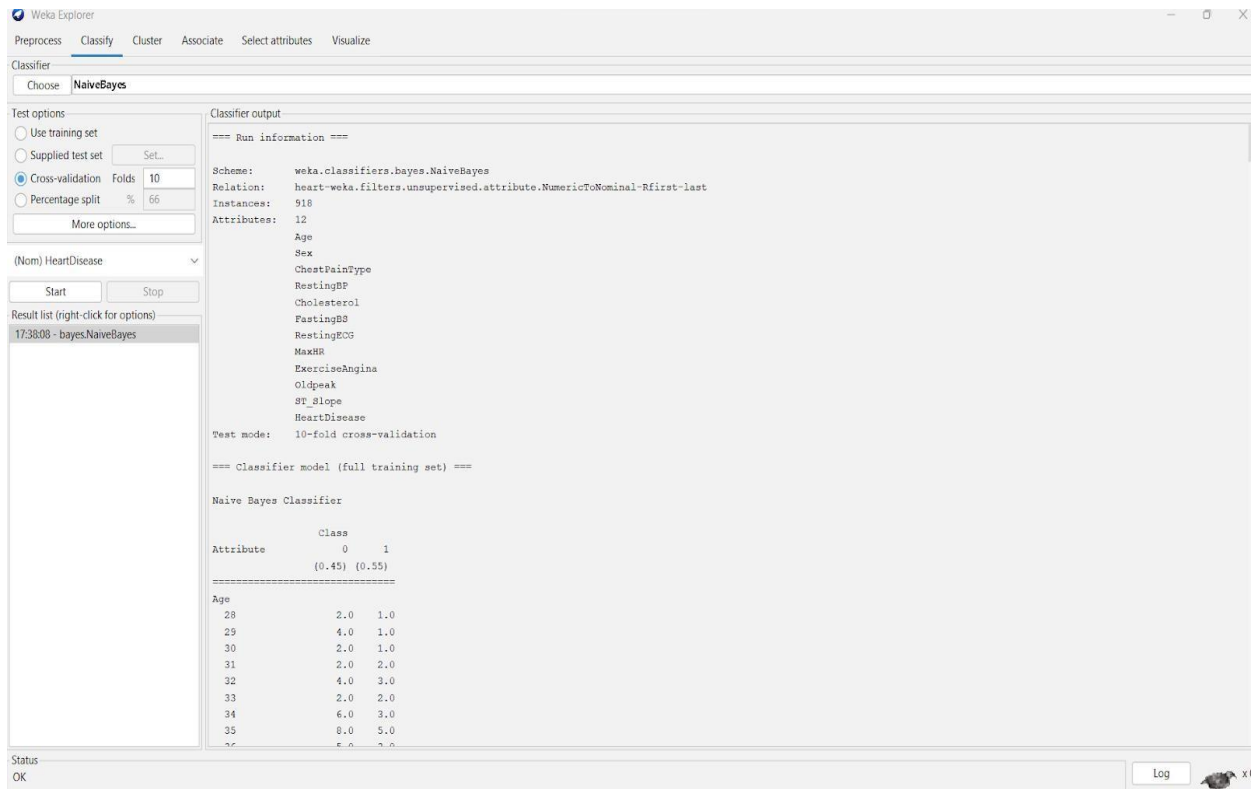
## Figure 3 (top screenshot)

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  %  66

More options...

(Nom) HeartDisease

Start | Stop

Result list (right-click for options)
17:38:08 - bayes.NaiveBayes

**Classifier output**

```
=== Run information ===

Scheme:       weka.classifiers.bayes.NaiveBayes
Relation:     heart-weka.filters.unsupervised.attribute.NumericToNominal-Rfirst-last
Instances:    918
Attributes:   12
              Age
              Sex
              ChestPainType
              RestingBP
              Cholesterol
              FastingBS
              RestingECG
              MaxHR
              ExerciseAngina
              Oldpeak
              ST_Slope
              HeartDisease
Test mode:    10-fold cross-validation

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

Naive Bayes Classifier

                 Class
Attribute          0       1
                 (0.45) (0.55)
===============================
Age
  28              2.0    1.0
  29              4.0    1.0
  30              2.0    1.0
  31              2.0    2.0
  32              4.0    3.0
  33              2.0    2.0
  34              6.0    3.0
  35              8.0    5.0
```

Status
OK

## Figure 4 (bottom screenshot)

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  %  66

More options...

(Nom) HeartDisease

Start | Stop

Result list (right-click for options)
17:38:08 - bayes.NaiveBayes

**Classifier output**

```
  [total]        463.0  561.0

ST_Slope
  Up             318.0   79.0
  Flat            80.0  382.0
  Down            15.0   50.0
  [total]        413.0  511.0


Time taken to build model: 0 seconds

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

Correctly Classified Instances         781               85.0763 %
Incorrectly Classified Instances       137               14.9237 %
Kappa statistic                          0.6983
Mean absolute error                      0.168
Root mean squared error                  0.3457
Relative absolute error                 33.9808 %
Root relative squared error             69.5366 %
Total Number of Instances              918

=== Detailed Accuracy By Class ===

              TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
              0.837    0.138    0.831      0.837   0.834      0.698  0.912     0.903     0
              0.862    0.163    0.867      0.862   0.865      0.698  0.912     0.909     1
Weighted Avg. 0.851    0.152    0.851      0.851   0.851      0.698  0.912     0.906

=== Confusion Matrix ===

   a    b   <-- classified as
 343   67 |  a = 0
  70  438 |  b = 1
```

Status
OK

**Figure 3& 4: Naïve bayes classification**

Applying decision tree classifier: A decision Tree is a Supervised learning technique that can be used for both classification and Regression problems, but mostly it is preferred for solving Classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome. While classifying the selected dataset, the J48 format was selected for Decision tree classification.
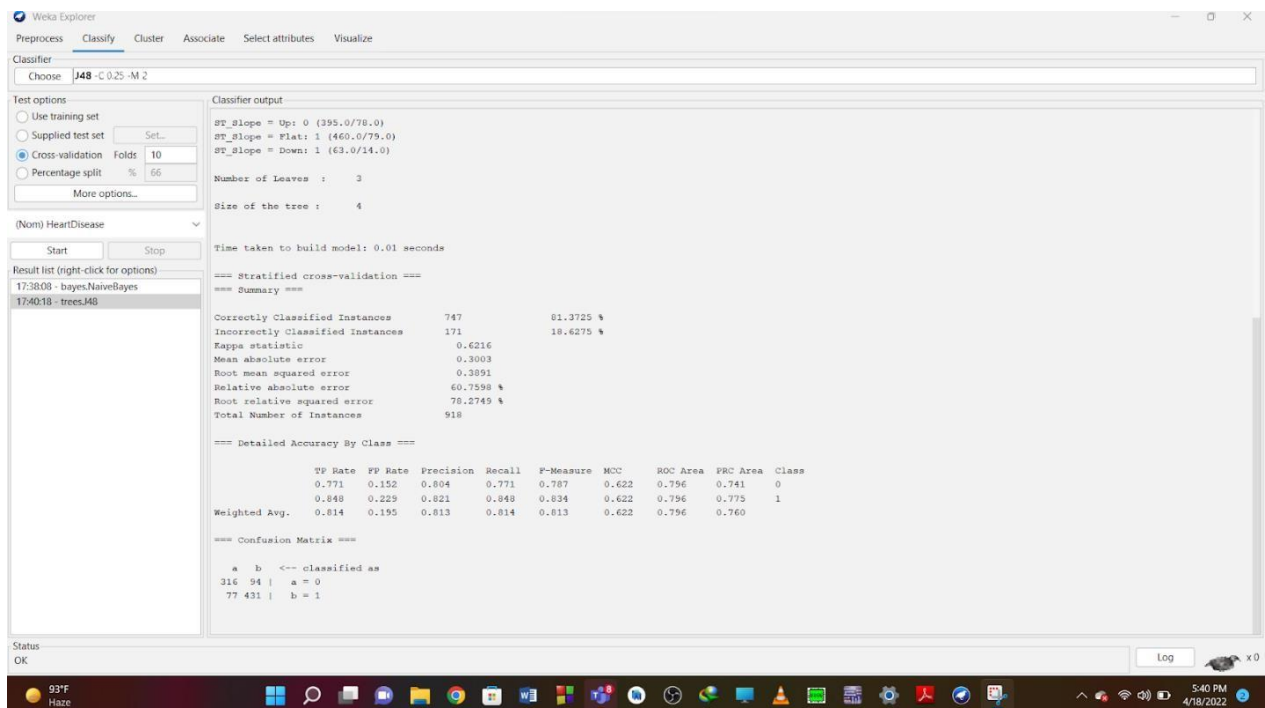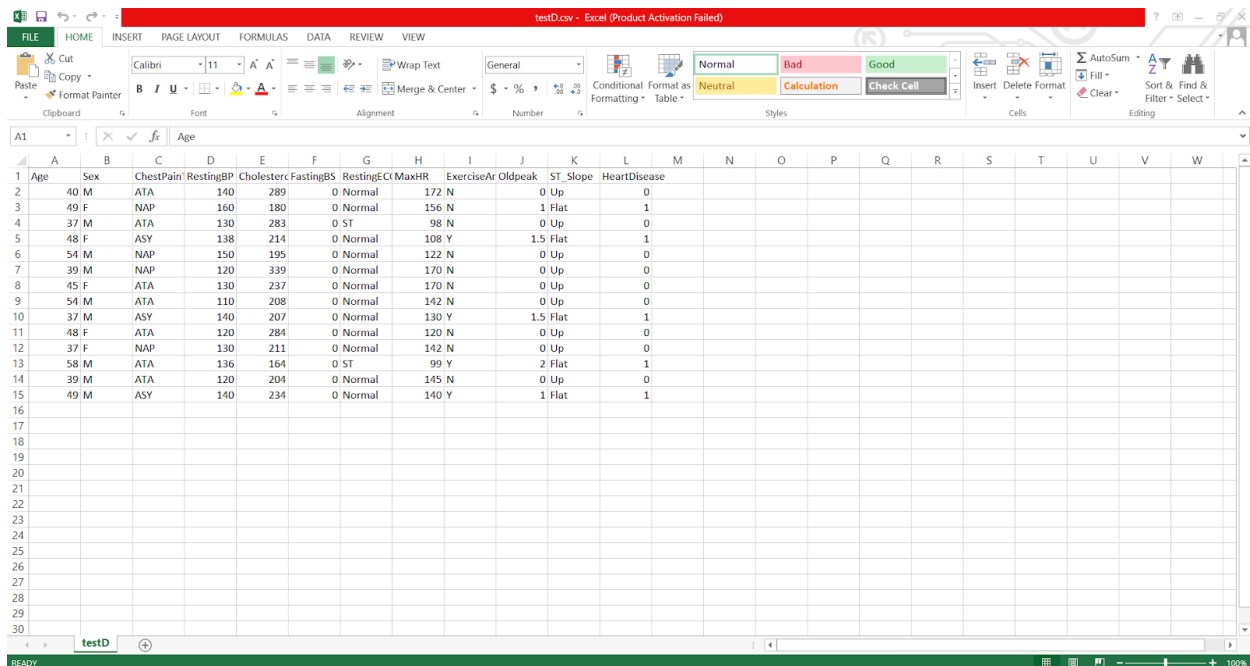


**Figure 5: Decision tree classification**

After applying two types of classifiers, the highest percentage of correctly classified instances is for the naïve Bayes classifier with 85.0763%. After that comes the Decision tree classifier with 81.3725%. The naïve Bayes classifier is considered the best classifier for the dataset.

Reason to choose naïve Bayes classifier: From the obtained result of the two classifiers, it is clearly seen that naïve Bayes has the highest percentage of correctly classified instances which is 85.0763%. As it has the most accurate value so it would be a more suitable classifier for the dataset. the advantages of naïve Bayes is

- It is simple and easy to implement.
- It doesn't require as much training data.

- It handles both continuous and discrete data.
- It is highly scalable with a number of predictors and data points.
- It is fast and can be used to make real-time predictions.

Here is the summary of the naïve Bayes classifiers result:

- Correctly Classified Instances  :   781          85.0763 %
- Incorrectly Classified Instances :  137          14.9237 %
- Kappa statistic                        : 0.6983
- Mean absolute error               :   0.168
- Root mean squared error          :   0.3457
- Relative absolute error             :   33.9808 %
- Root relative squared error        :   69.5366 %
- Total Number of Instances         :    918

**DISCUSSION:**

In this task 1 section at first, I have chosen a proper dataset that has proper context and content and is also applicable for two of my selected classifiers. After that I selected two classifiers from weka software that is applicable to my dataset with proper classification and accuracy. after applying these classifiers to my dataset, I got two different correctly classified instances, which are 85.0763%, 81.3725%. Also, with the value of some parameters like IP rate fp rate recall f measure, etc., and with the help of the confusion matrix I have selected my best classifier. confusion matrix helps us to identify his good is my algorithm doing. that's how I have completed task 1 and found the best classifier which is the naïve Bayes classifier.

## INTRODUCTION:

One of the most important mechanisms in machine learning is to train your algorithm on a training set that is separate and distinct from the test set for which will be gauged its accuracy. To detect a machine learning behavior, a training dataset has been set which was extracted subset from the referred dataset. Then this model had been tested with a test dataset which is a subset to test the trained model. While preparing the test dataset, things that were made sure are that the dataset was large enough to yield statistically meaningful results. Also, it was representative of the data set as a whole. In other words, the test sets with unusual characteristics than the training set were not chosen. The suitable classifier is then used to predict the classification for the instances in the test set. If the test set contains N instances of which C is correctly classified, C is correctly classified Predictive accuracy, $P = C/N$. There are 14 instances in this prepared test dataset.



**Figure 6: test dataset**

**PROCEDURE OF TESTING THE TEST DATASET:**

1. First, Weka 3.8.6 was opened and the 'Explorer' option was chosen and such window named weka explorer was opened.



**Figure 7: Weka Explorer**

2. Then, the open file option was selected and the extracted .arff file, training dataset was selected from the device.
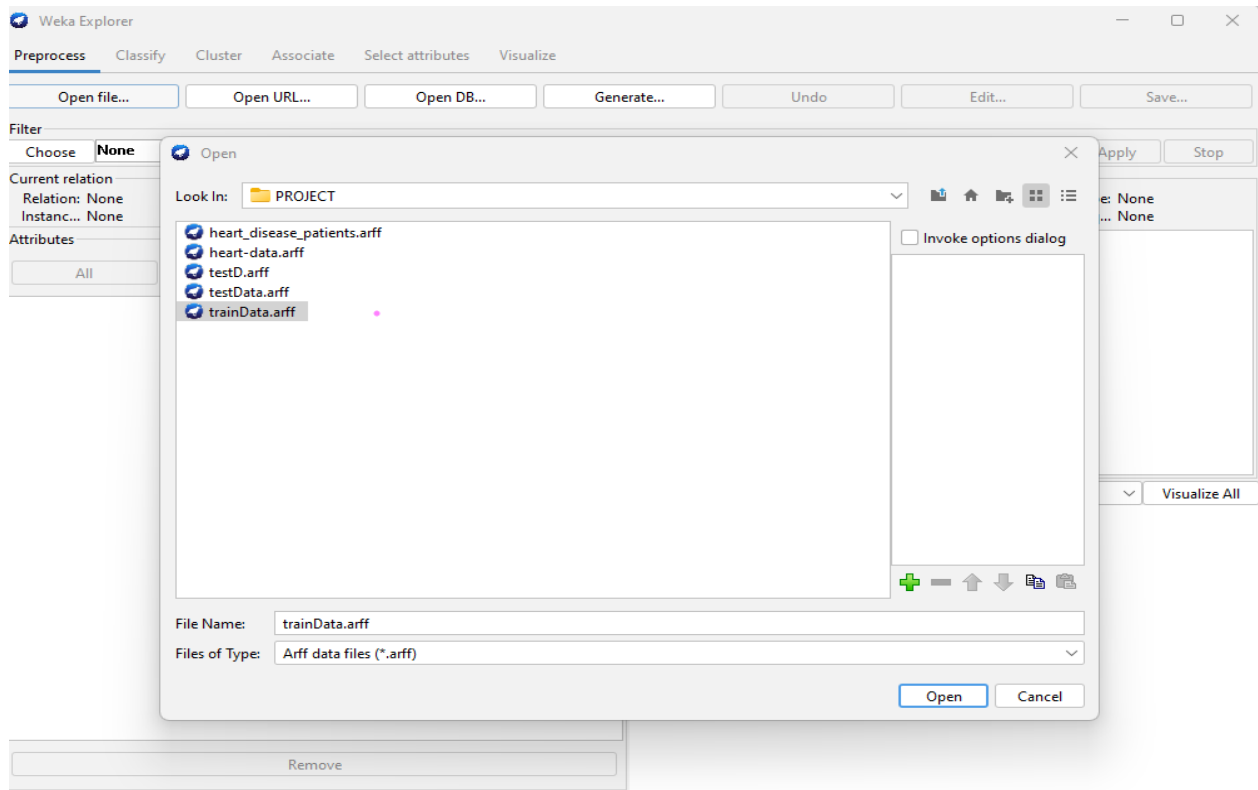


**Figure 8: Training Dataset**

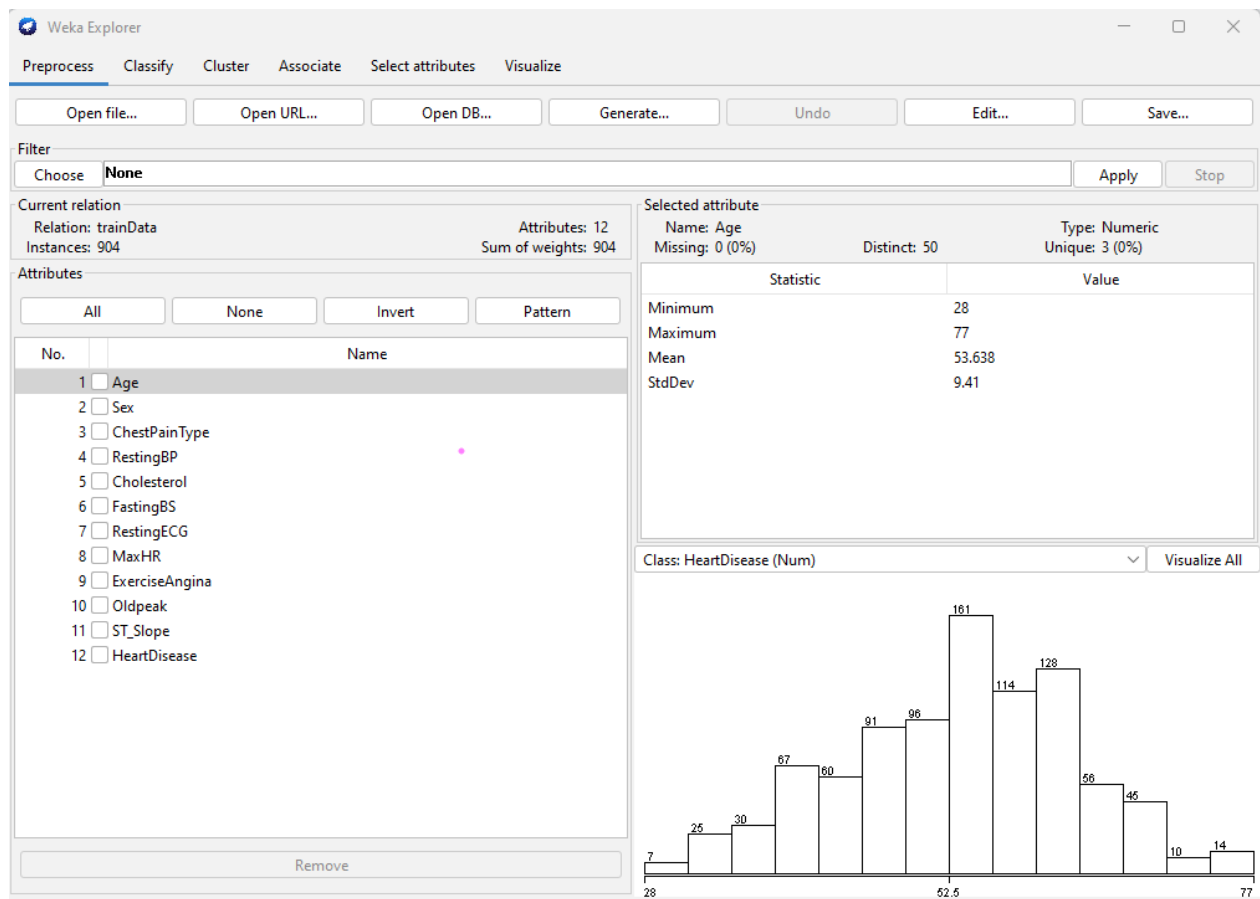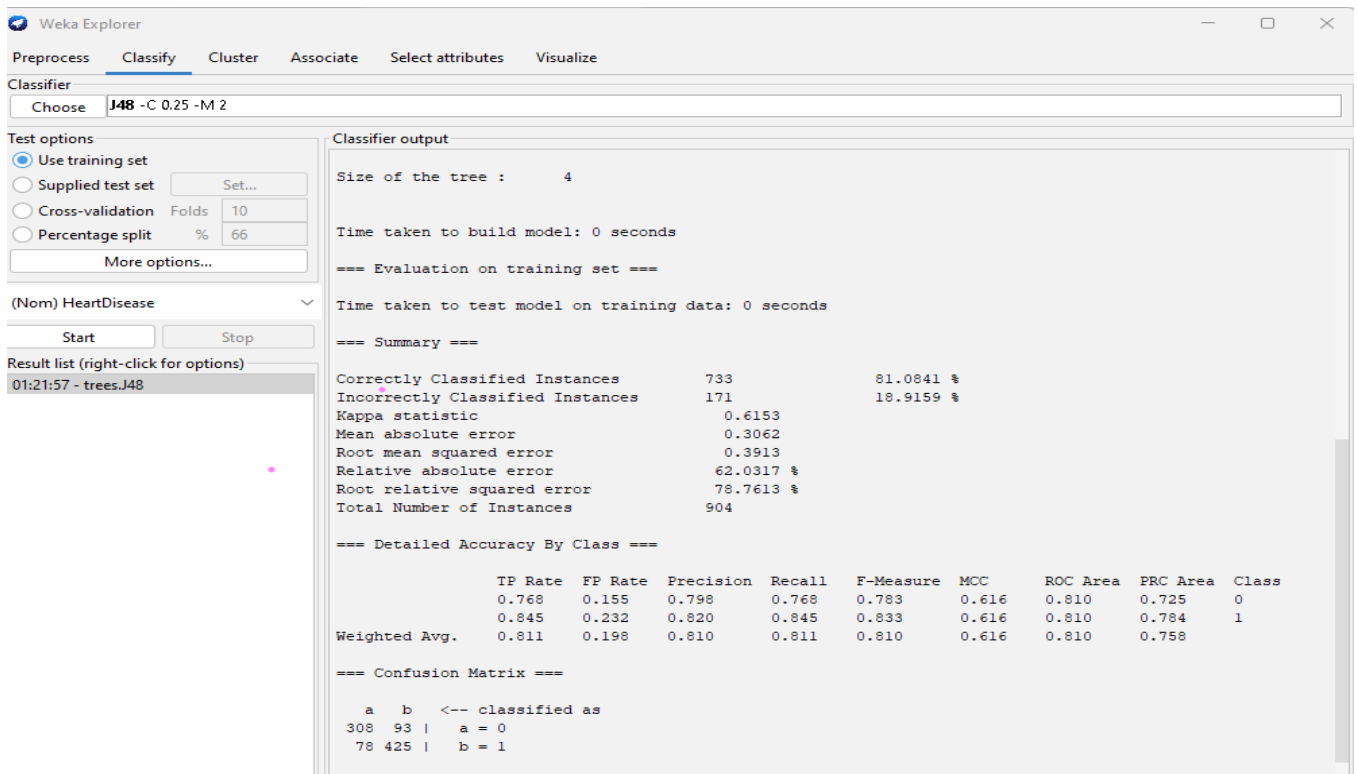3. After the open option was clicked, the details of the dataset popped.



**Figure 9: training dataset**

4. Then the preferred classifier (Decision tree classifier) for the training dataset was selected. Then from the test options, the use training set was chosen and the start option was selected.





**Figures 10 & 11: the result of the training set**

5. To input the test set, from the test options, supplied test data was selected and then a window named test instances came.
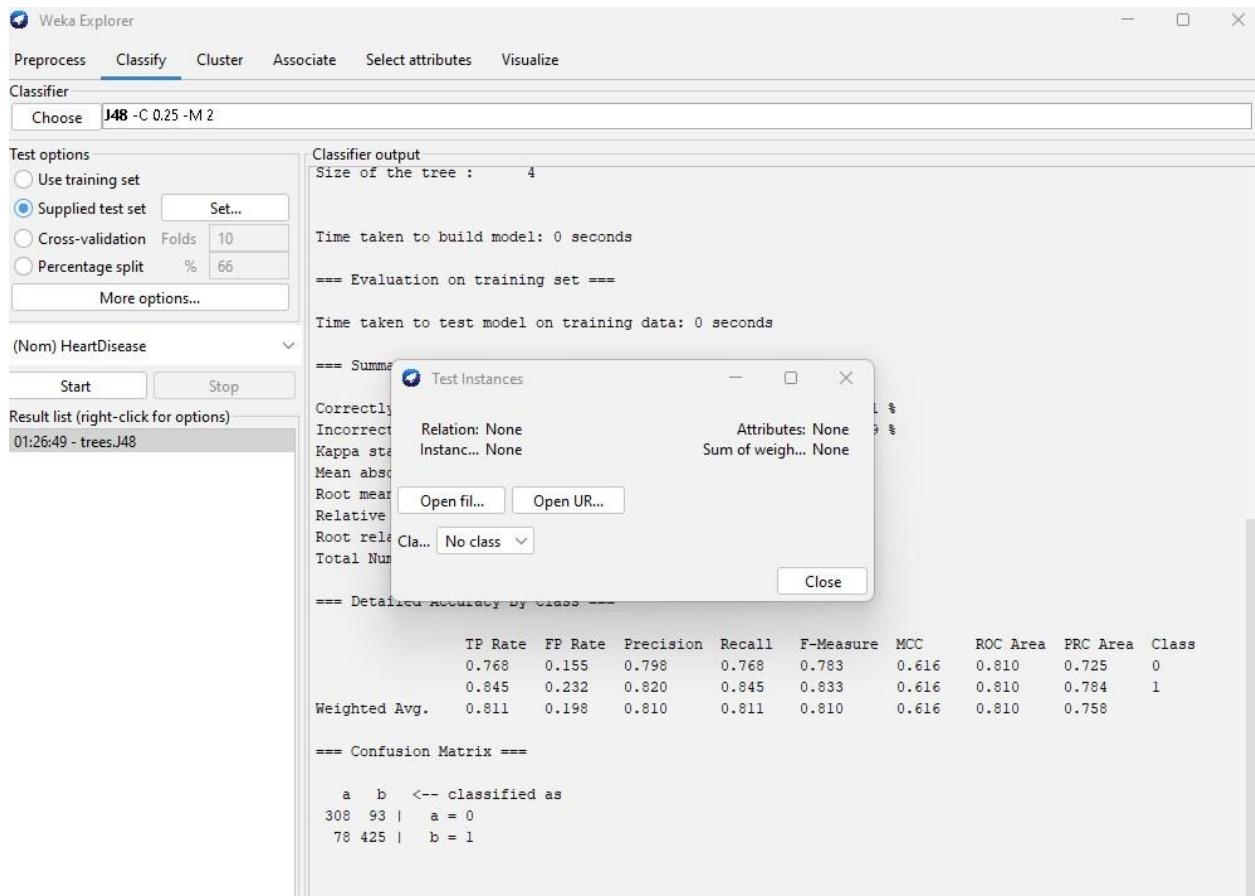


**Figure 12: supplied test set selected**

6. Then the open file option was chosen to insert the test data and then the dataset was opened.
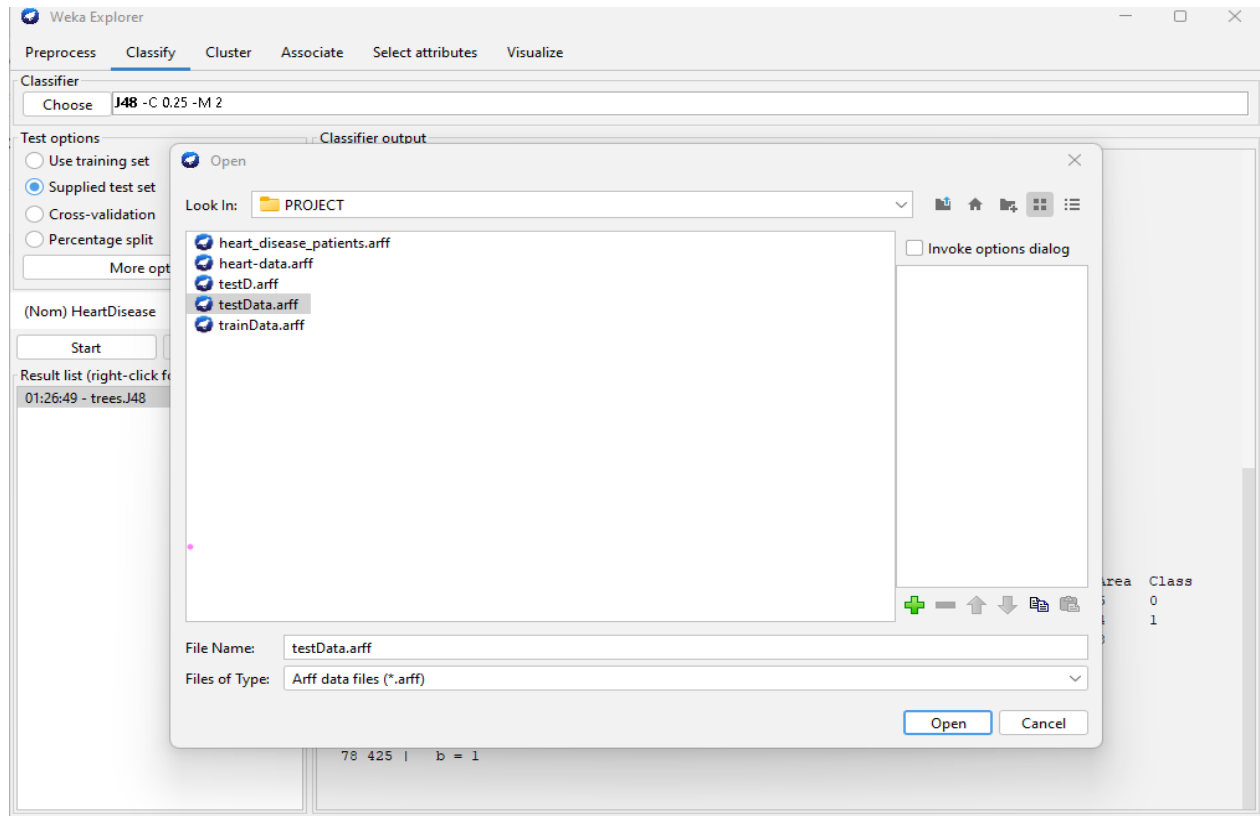


**Figure 13: test set opened**

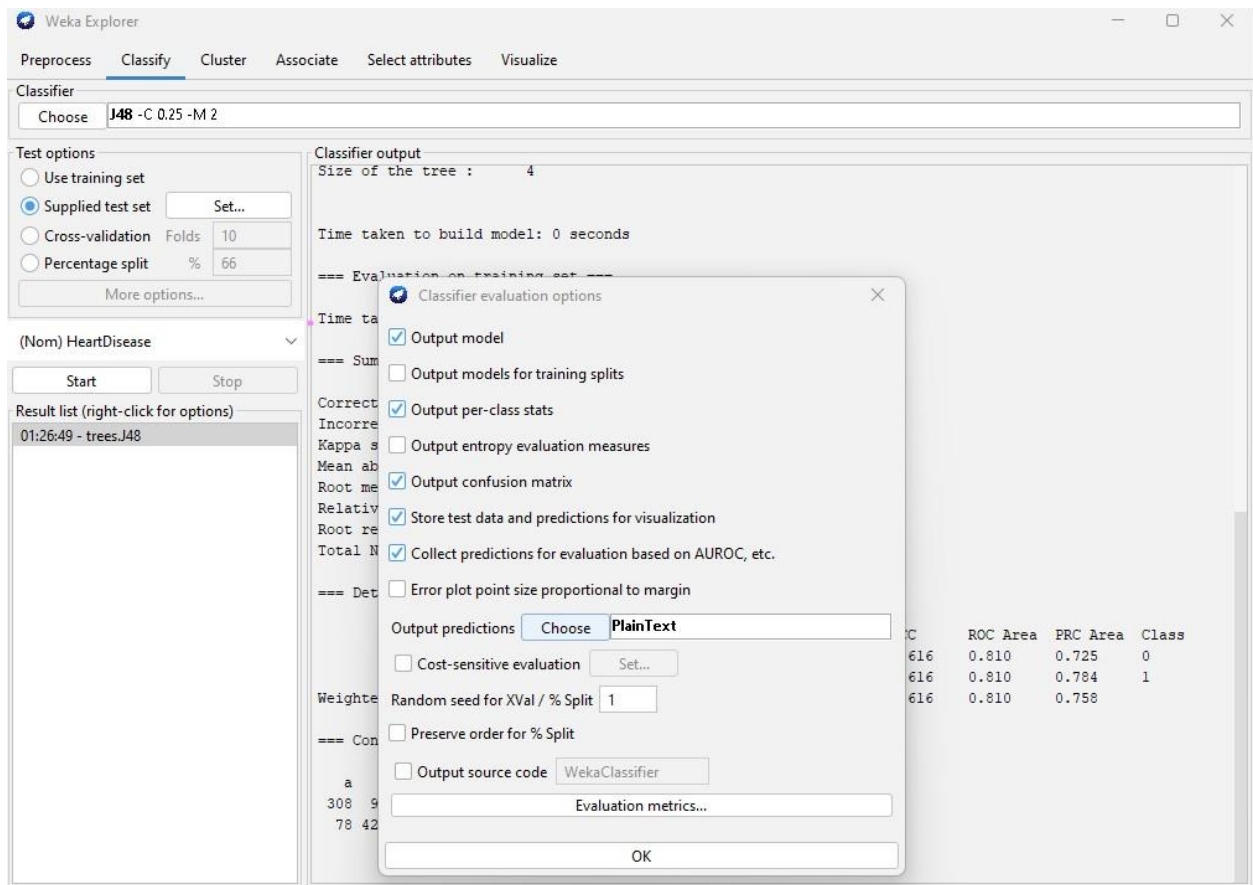7.  To make sure the test set works properly, it was made sure that the output predictions was in PlainTest form



**Figure 14: PlainText format chosen for output prediction**
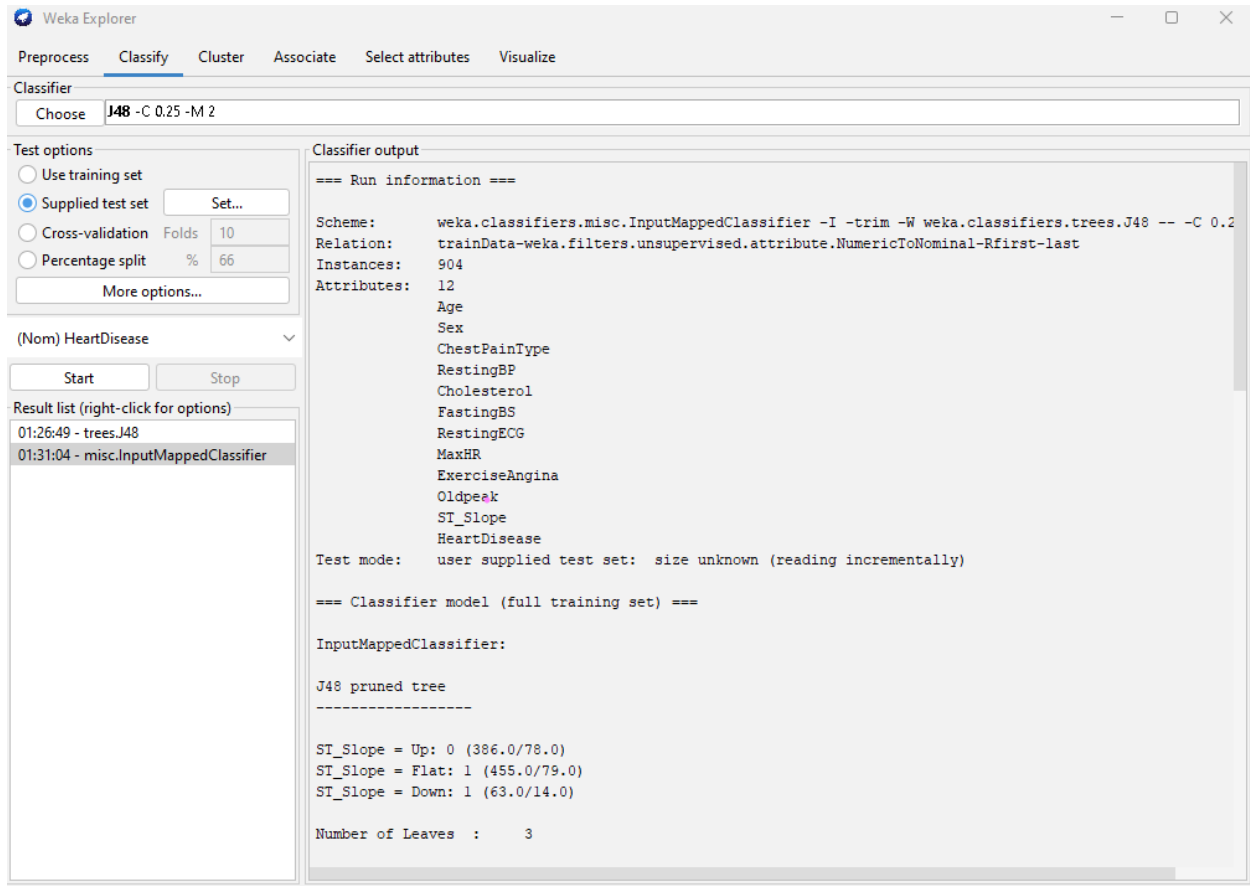
8. Then the start button was clicked.



**Figure 15: Test set run**

**RESULT OF TEST-DATASET MODEL:**

Once the start button was clicked, the output for the test dataset came. In the result, the test mode was 'user-supplied test set' which means the classifier is evaluated on how well it predicts the class of a set of instances loaded from a file that was inputted by the user. The total time taken to build the model was 0.0 seconds and among the 14 instances, there were 0 instances where the error was found. In this model, the Correctly classified instances % Value describes the amount of accuracy of correctly classified instances provided by the algorithm. In this case, the percentage is 100%. Incorrectly classified instances % Value describes how much incorrect instances are given by the algorithm. In this case, the percentage is 0%.

Mean Absolute Error (MAE): It can define as a statistical measure of how far an estimate is from actual values i.e. the average of the absolute magnitude of the individual errors. It is usually similar in magnitude but slightly smaller than the root mean squared error. In this model the MAE is 0.1919

Root Mean-Squared Error (RMSE): The Root Mean Square Error (RMSE) calculates the differences between values predicted by a model / an estimator and the values observed from the thing being modeled/ estimated. RMSE is used to measure the accuracy. It is ideal if it is small. In this case, the RMSE is 0.1924 which is ideal.

**DISCUSSION:**

The idea behind training and test sets is to test the generalization error. Separating data into training and testing sets is an important part of evaluating data mining models. Since the test set data already contains known values for the property I want to predict, it is easy to determine if the model's assumptions are correct. So here I carefully separate the data into test/training and apply them in weka after that it gives the error table and accuracy rate that's how I get my proper result and accuracy model.

**INTRODUCTION:**

Unsupervised Learning is a machine learning technique in which the users do not need to supervise the model. Instead, it allows the model to work on its own to discover patterns and information that was previously undetected. It mainly deals with the unlabelled data. Unsupervised learning is helpful for finding useful insights from the data. Unsupervised learning is much similar to human learning to think by their own experiences, which makes it closer to the real AI. Unsupervised learning works on unlabeled and uncategorized data which makes unsupervised learning more important.

I have chosen " Heart_Disease_Patients " dataset. I will also use K means clustering Algorithm.

About the dataset:n this report, the used " Heart_Disease_Patients ", a CSV dataset file [ Later converted into .arff ], collected from Kaggle.com.

**RESULT:**

Applying K means clustering Algorithm.

K-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. It allows us to cluster the data into different groups and is a convenient way to discover the categories of groups in the unlabeled dataset on its own without the need for any training.

## Weka Explorer — Window 1

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

**Clusterer**

Choose | SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

**Cluster mode**
- (●) Use training set
- ( ) Supplied test set — Set...
- ( ) Percentage split — % 66
- ( ) Classes to clusters evaluation
- (Num) slope
- [✓] Store clusters for visualization

Ignore attributes

Start | Stop

**Result list (right-click for options)**
02:42:14 - SimpleKMeans

**Clusterer output**

```
=== Run information ===

Scheme:       weka.clusterers.SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "w
Relation:     heart_disease_patients
Instances:    303
Attributes:   12
              id
              age
              sex
              cp
              trestbps
              chol
              fbs
              restecg
              thalach
              exang
              oldpeak
              slope
Test mode:    evaluate on training data


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


kMeans
======

Number of iterations: 4
Within cluster sum of squared errors: 332.45018999494914

Initial starting points (random):
```

**Status**
OK

Log | x 0

---

## Weka Explorer — Window 2

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

**Clusterer**

Choose | SimpleKMeans -init 0 -max-candidates 100 -periodic-pruning 10000 -min-density 2.0 -t1 -1.25 -t2 -1.0 -N 2 -A "weka.core.EuclideanDistance -R first-last" -I 500 -num-slots 1 -S 10

**Cluster mode**
- (●) Use training set
- ( ) Supplied test set — Set...
- ( ) Percentage split — % 66
- ( ) Classes to clusters evaluation
- (Num) slope
- [✓] Store clusters for visualization

Ignore attributes

Start | Stop

**Result list (right-click for options)**
02:42:14 - SimpleKMeans

**Clusterer output**

```
kMeans
======

Number of iterations: 4
Within cluster sum of squared errors: 332.45018999494914

Initial starting points (random):

Cluster 0: 49,65,0,3,140,417,1,2,157,0,0.8,1
Cluster 1: 285,61,1,4,148,203,0,0,161,0,0,1

Missing values globally replaced with mean/mode

Final cluster centroids:
                         Cluster#
Attribute    Full Data        0          1
              (303.0)      (70.0)     (233.0)
==============================================
id                152     144.8143    154.1588
age            54.4389     57.0857     53.6438
sex             0.6799      0.2571      0.8069
cp              3.1584      3.2286      3.1373
trestbps      131.6898    137.5571    129.927
chol          246.6931    268.9571    240.0043
fbs             0.1485      0.4286      0.0644
restecg         0.9901      1.8429      0.7339
thalach       149.6073    146.8286    150.4421
exang           0.3267      0.3143      0.3305
oldpeak         1.0396      1.1         1.0215
slope           1.6007      1.7286      1.5622
```
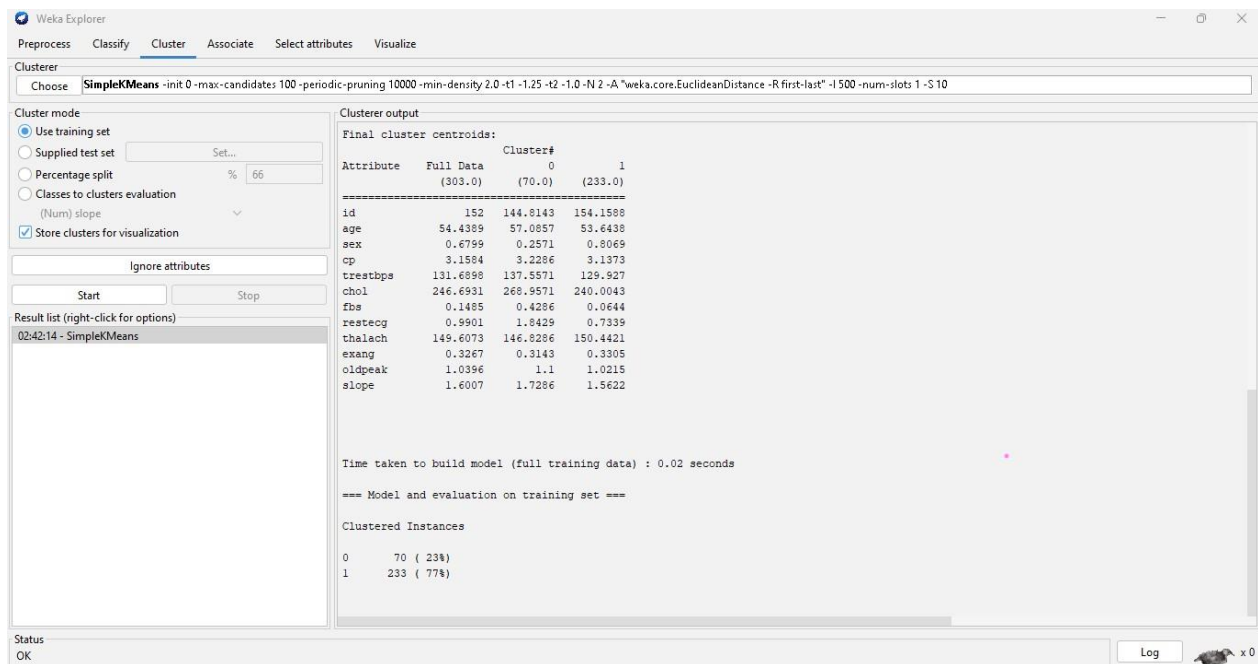
**Status**
OK

Log | x 0

**Figures 16, 17,18: K means clustering Algorithm**

Here is the summary of the K means clustering Algorithm result:

Number of iterations: 4

Initial starting points (random):

Cluster 0: 49,65,0,3,140,417,1,2,157,0,0.8,1

Cluster 1: 285,61,1,4,148,203,0,0,161,0,0,1

Time taken to build model (full training data) : 0.02 seconds

Clustered Instances

0      70 (23%)

1      233 (77%)

**DISCUSSION:**

I have chosen a proper unsupervised dataset. I convert csv file to arff file and then apply K means Clustering Algorithm. By default, the value of K is 2 so there are two results.

For 0 instances it is 23% whereas for 1 is 77%. I also find the final cluster of centroids also. By that I have completed K means clustering for unsupervised data.

**Reference:**

1. https://www.kaggle.com/datasets/kingabzpro/heart-disease-patients
2. https://www.kaggle.com/datasets/fedesoriano/heart-failure-prediction
3. https://en.wikipedia.org/wiki/Data_mining
4. https://www.javatpoint.com/k-means-clustering-algorithm-in-machine-learning