**AMERICAN INTERNATIONAL UNIVERSITY–BANGLADESH (AIUB)**

**FACULTY OF SCIENCE & TECHNOLOGY**

**DEPARTMENT OF COMPUTER SCIENCE**

**DATA WAREHOUSING AND DATA MINING**

**Summer 2021-2022**

**Section: C**

**Supervised By**

Akinul Islam Jony

**Submitted By**

| Name | ID |
|------|-----|
| Either Rahman | 19-39750-1 |
| Hasan Sanjary Islam | 19-39589-1 |
| Farahn Sadik Ferdous | 20-42072-1 |
| Tapu Biswas | 20-42073-1 |

**Date of Submission: August 7, 2022**

## Project Overview

Data mining is a process where data are extracted from huge datasets according to one's usefulness. Data mining is basically done to gather knowledge from data. Extracted data can be used for more effective marketing strategies, increasing sales, decreasing costs, fraud detection, product design, weather forecasting, etc. Data mining generally breaks down patterns and connections in data based on what information users request. KNN, Naïve Bayes, and Decision Tree are some of the algorithms used in data mining. Naïve Bayes and KNN are classification method that only deals with categorical attributes or classified dataset. Naïve Bayes is a supervised learning algorithm that is used for solving classification problems. Naïve Bayes uses probability theory to find the most likely of the possible classifications. KNN is a classifier algorithm where the classification is based on the nearest Neighbor's data (distance between instances). A decision tree is a type of probability tree that enables you to make a decision on a process followed by some decision rule. Decision rules can be combined together to form a tree structure. Prediction accuracy defines the predicted values match the actual values of the target field. Prediction accuracy is used to estimate the performance of a classifier. Prediction accuracy determines the proportion of a set of unseen instances that it correctly classifies. The confusion matrix gives an overview of prediction accuracy thoroughly like instances of the correct class being correctly classified as the correct class or misclassified as some other class and its breakdown. In this project, we have determined the accuracy of a dataset. We have collected a dataset that is Bank marketing response prediction. We have applied three algorithms KNN, Naïve Bayes, and Decision Tree to WEKA software. To predict the accuracy and its breakdown we have used prediction accuracy and confusion matrix.

## Dataset Overview

i.   Data source with valid URL

Bank Marketing Response Predict

https://www.kaggle.com/datasets/kukuroo3/bank-marketing-response-predict

ii.   Description about dataset

The dataset is collected from Kaggle. We have taken a classified dataset from Kaggle. The data is from a marketing campaign (phone calls) of a Portuguese banking institution. The classification goal is to predict if the client will subscribe to a term deposit. This dataset has 12870 instances and 17 attributes. The class attribute here is the attribute y. There is no missing data here. We have applied three algorithms KNN, Naïve Bayes, and Decision Tree on WEKA software in this dataset. To predict the accuracy and its breakdown we have used a prediction accuracy and confusion matrix.

Figure: Bank Marketing Response Predict Dataset

## Model Development

Classification in data mining is a common technique that separates data points into different classes. It allows you to organize data sets of all sorts, including complex and large datasets as well as small and simple ones. It primarily involves using algorithms that you can easily modify to improve the data quality. Here, we have applied three algorithms Naïve Bayes, KNN and Decision Tree on WEKA software in this dataset. To predict the accuracy and its breakdown we have used a prediction accuracy and confusion matrix.

Before inserting the dataset into the WEKA software, the .csv file was converted into .arff file. Procedure of inserting dataset into the WEKA software is given below:

1. First, Weka 3.8.6 was opened and 'Explorer' option was chosen and new window named Weka Explorer was opened.



Figure: Weka Explorer

2. Then, train dataset was selected from the device.



Figure: Train dataset

## Naïve Bayes classifier:

1. Firstly, classify was selected, 'Choose' option was pressed then NaiveBayes was chosen, and at last start button was pressed.



Figure: Naïve Bayes Algorithm



Figure: Naïve Bayes Algorithm

Figure: Naïve Bayes Algorithm



Figure: Naïve Bayes Algorithm

## Figure 1 (Weka Explorer - Naïve Bayes Algorithm)

Weka Explorer

RConsole | Visualize 3D | Forecast | Projection Plot | RVines | Dl4j Inference | CPython Scripting

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize | Auto-WEKA | Interactive Parallel Coordinates Plot

Classifier
Choose | NaiveBayes

Test options
- Use training set
- Supplied test set [Set...]
- Cross-validation  Folds [10]
- Percentage split  % [66]
- More options...

(Nom) y

[Start] [Stop] [Run on server]

Result list (right-click for options)
22:33:43 - bayes.NaiveBayes

Classifier output

```
campaign
  mean                    2.9603      2.2687
  std. dev.               3.0977      1.8763
  weight sum              8903        3967
  precision               1.2         1.2

pdays
  mean                   34.5724     67.9996
  std. dev.              95.9654    118.9524
  weight sum              8903        3967
  precision               1.8587      1.8587

previous
  mean                    0.4291      1.1021
  std. dev.               1.7053      2.6808
  weight sum              8903        3967
  precision               2           2

poutcome
  unknown              7535.0      2537.0
  failure               903.0       479.0
  other                 342.0       235.0
  success               127.0       720.0
  [total]              8907.0      3971.0



Time taken to build model: 0.07 seconds

=== Stratified cross-validation ===
=== Summary ===
```
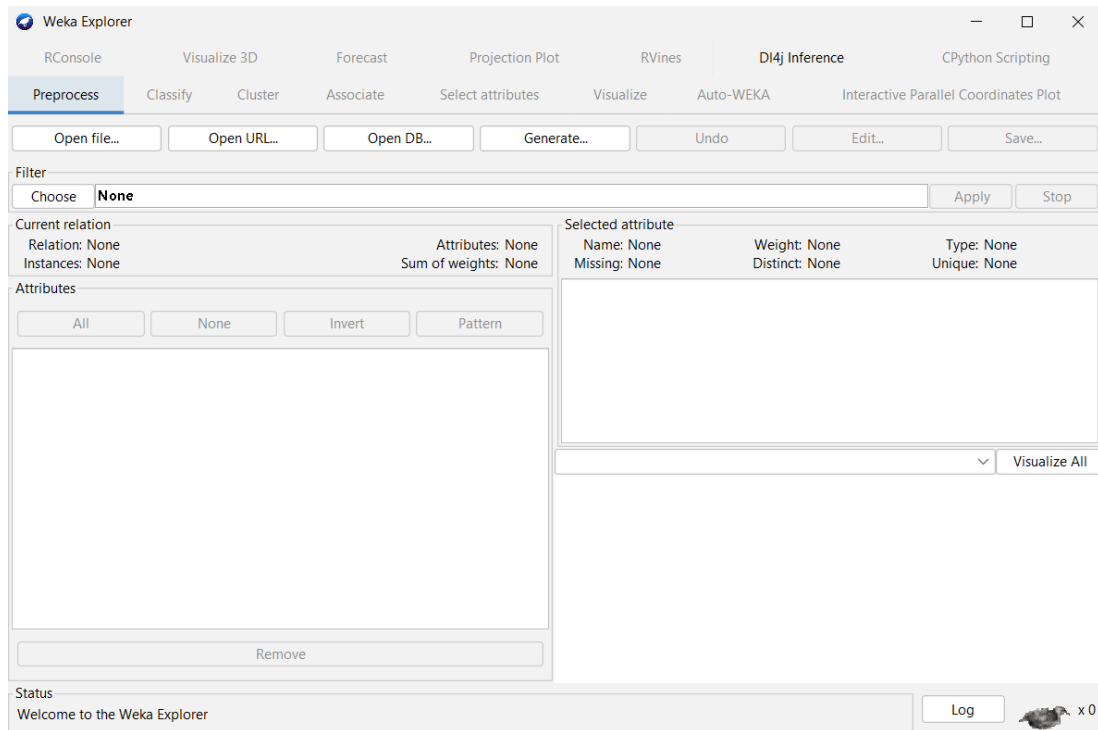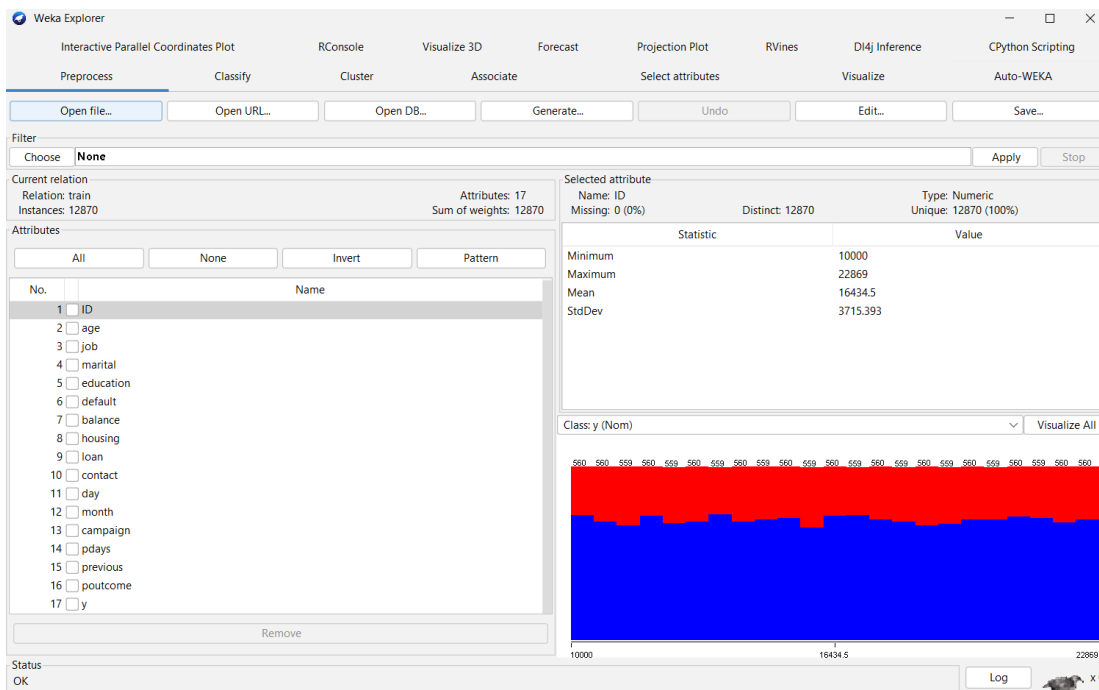
Status
OK

Log | x 0

Figure: Naïve Bayes Algorithm

## Figure 2 (Weka Explorer - Naïve Bayes Algorithm)

Weka Explorer

RConsole | Visualize 3D | Forecast | Projection Plot | RVines | Dl4j Inference | CPython Scripting

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize | Auto-WEKA | Interactive Parallel Coordinates Plot

Classifier
Choose | NaiveBayes

Test options
- Use training set
- Supplied test set [Set...]
- Cross-validation  Folds [10]
- Percentage split  % [66]
- More options...

(Nom) y

[Start] [Stop] [Run on server]

Result list (right-click for options)
22:33:43 - bayes.NaiveBayes

Classifier output

```
Time taken to build model: 0.07 seconds

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

Correctly Classified Instances        9448            73.411 %
Incorrectly Classified Instances      3422            26.589 %
Kappa statistic                          0.3539
Mean absolute error                      0.3058
Root mean squared error                  0.4529
Relative absolute error                 71.716 %
Root relative squared error             98.0894 %
Total Number of Instances            12870

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 0.836    0.494    0.791      0.836   0.813      0.356  0.736     0.839     no
                 0.506    0.164    0.579      0.506   0.540      0.356  0.736     0.549     yes
Weighted Avg.    0.734    0.393    0.726      0.734   0.729      0.356  0.736     0.750

=== Confusion Matrix ===

    a    b   <-- classified as
 7442 1461 |   a = no
 1961 2006 |   b = yes
```

Status
OK

Log | x 0

Figure: Naïve Bayes Algorithm

## K-Nearest Neighbour Classification:

1. Firstly, classify was selected, 'Choose' option was pressed then IBK was chosen, and at last start button was pressed.



Figure: KNN Algorithm



Figure: KNN Algorithm

## Decision Trees:

1. Firstly, classify was selected, 'Choose' option was pressed then J48 was chosen, and at last start button was pressed.



Figure: Decision Tree Algorithm
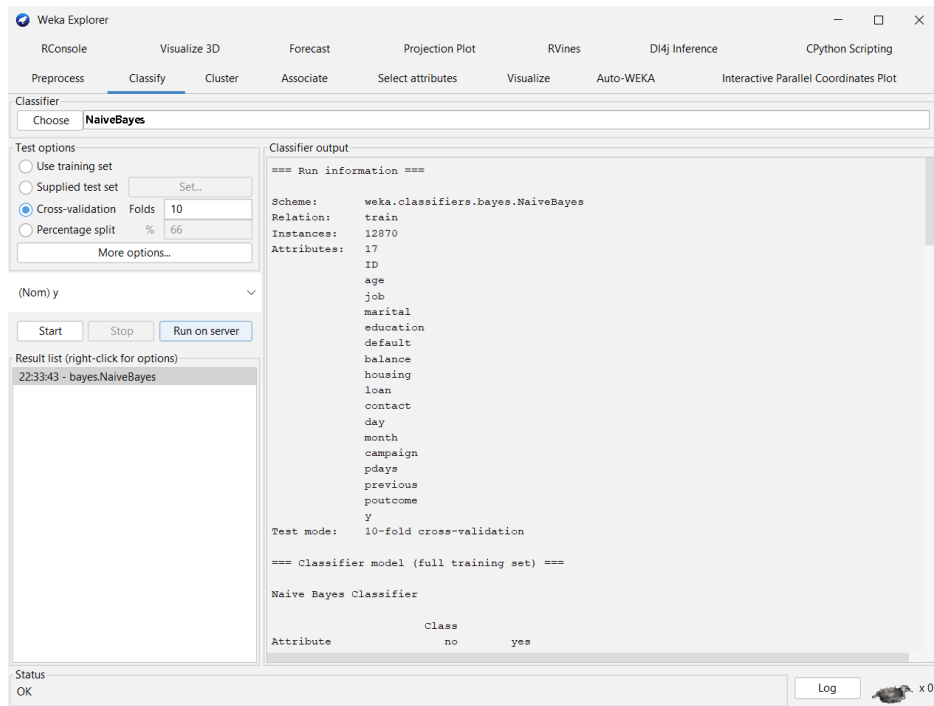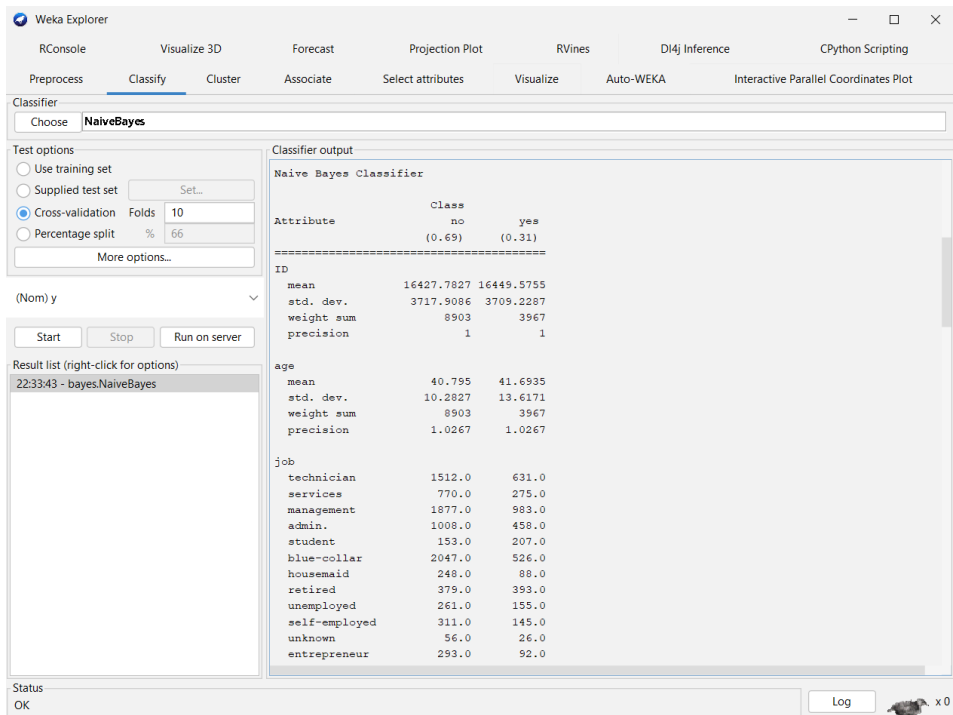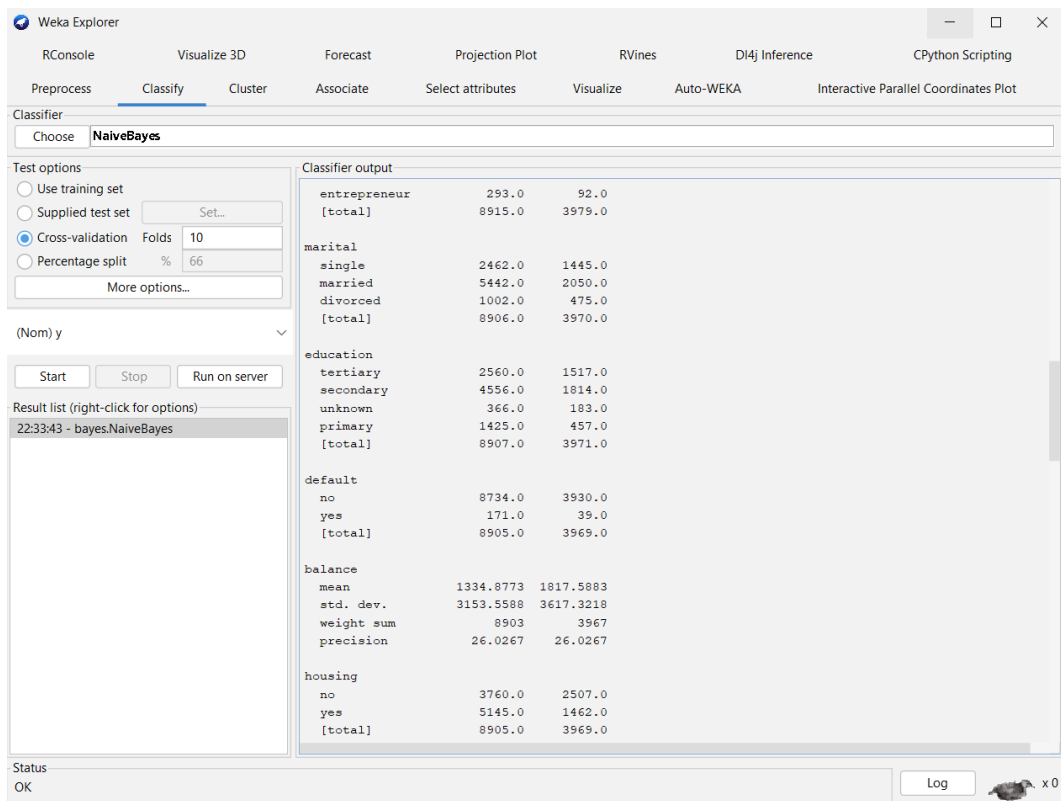


Figure: Decision Tree Algorithm

RConsole    Visualize 3D    Forecast    Projection Plot    RVines    Dl4j Inference    CPython Scripting

Preprocess    Classify    Cluster    Associate    Select attributes    Visualize    Auto-WEKA    Interactive Parallel Coordinates Plot

Classifier

Choose    J48 -C 0.25 -M 2

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

More options...

(Nom) y

Start    Stop    Run on server

Result list (right-click for options)
15:21:27 - trees.J48

Classifier output

```
|  |  |  |  |       education = unknown
|  |  |  |  |         marital = single: yes (6.0)
|  |  |  |  |         marital = married: no (3.0/1.0)
|  |  |  |  |         marital = divorced: yes (0.0)
|  |  |  |  |       education = primary
|  |  |  |  |         ID <= 18469: no (5.0)
|  |  |  |  |         ID > 18469: yes (3.0/1.0)
|  |  |  |     housing = yes: no (708.0/184.0)
|  |  |     month = jun
|  |  |  |     job = technician
|  |  |  |  |       education = tertiary
|  |  |  |  |         age <= 36: yes (8.0)
|  |  |  |  |         age > 36
|  |  |  |  |         |  day <= 4: no (7.0)
|  |  |  |  |         |  day > 4: yes (3.0/1.0)
|  |  |  |  |       education = secondary
|  |  |  |  |         day <= 11: yes (14.0)
|  |  |  |  |         day > 11
|  |  |  |  |         |  age <= 31: yes (2.0)
|  |  |  |  |         |  age > 31: no (2.0)
|  |  |  |  |       education = unknown: yes (2.0)
|  |  |  |  |       education = primary: yes (0.0)
|  |  |  |     job = services: yes (10.0)
|  |  |  |     job = management
|  |  |  |  |     balance <= 243: no (6.0/1.0)
|  |  |  |  |     balance > 243
|  |  |  |  |     |  balance <= 935: yes (16.0)
|  |  |  |  |     |  balance > 935
|  |  |  |  |     |  |  housing = no
|  |  |  |  |     |  |  |  day <= 4
|  |  |  |  |     |  |  |  |  balance <= 4210
```
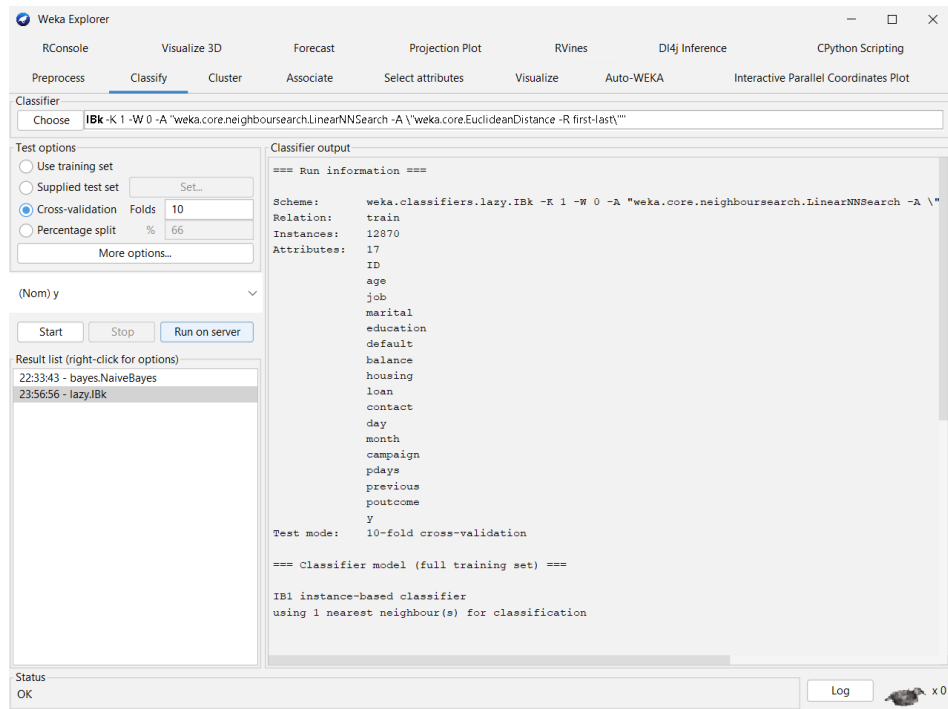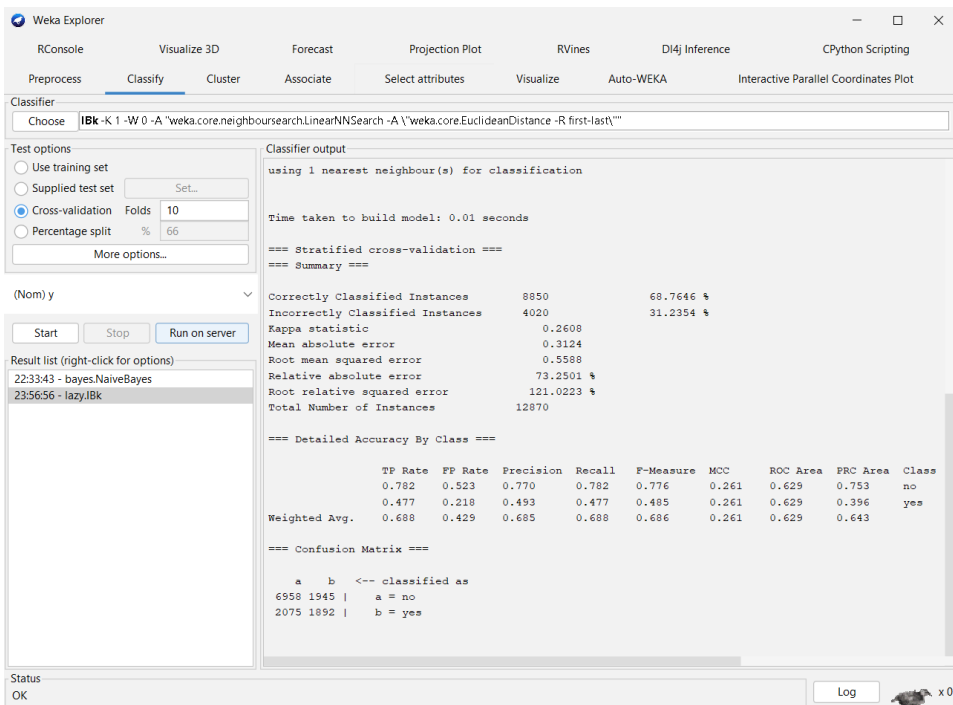
Status
OK

Log    x 0

Figure: Decision Tree Algorithm

---

RConsole    Visualize 3D    Forecast    Projection Plot    RVines    Dl4j Inference    CPython Scripting

Preprocess    Classify    Cluster    Associate    Select attributes    Visualize    Auto-WEKA    Interactive Parallel Coordinates Plot

Classifier

Choose    J48 -C 0.25 -M 2

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

More options...

(Nom) y

Start    Stop    Run on server

Result list (right-click for options)
15:21:27 - trees.J48

Classifier output

```
|  |  |  |  |  |  |  |     balance <= 4210
|  |  |  |  |  |  |  |       age <= 54: yes (14.0/2.0)
|  |  |  |  |  |  |  |       age > 54: no (2.0)
|  |  |  |  |  |  |     balance > 4210: no (3.0)
|  |  |  |  |  |     day > 4: yes (11.0)
|  |  |  |  |     housing = yes
|  |  |  |  |       campaign <= 1: yes (3.0/1.0)
|  |  |  |  |       campaign > 1: no (2.0)
|  |  |     job = admin.
|  |  |  |     loan = no
|  |  |  |  |     housing = no
|  |  |  |  |     |  ID <= 10711: no (2.0)
|  |  |  |  |     |  ID > 10711: yes (16.0/2.0)
|  |  |  |  |     housing = yes
|  |  |  |  |     |  ID <= 14810: yes (2.0)
|  |  |  |  |     |  ID > 14810: no (3.0)
|  |  |  |     loan = yes: no (4.0/1.0)
|  |  |     job = student
|  |  |  |     campaign <= 2: yes (8.0/1.0)
|  |  |  |     campaign > 2: no (3.0)
|  |  |     job = blue-collar
|  |  |  |     housing = no: yes (11.0/1.0)
|  |  |  |     housing = yes: no (3.0/1.0)
|  |  |     job = housemaid: yes (3.0/1.0)
|  |  |     job = retired: yes (6.0/1.0)
|  |  |     job = unemployed
|  |  |  |     campaign <= 2: no (5.0/1.0)
|  |  |  |     campaign > 2
|  |  |  |  |     age <= 29: no (2.0)
|  |  |  |  |     age > 29: yes (5.0/1.0)
|  |  |     job = self-employed: yes (7.0/1.0)
```
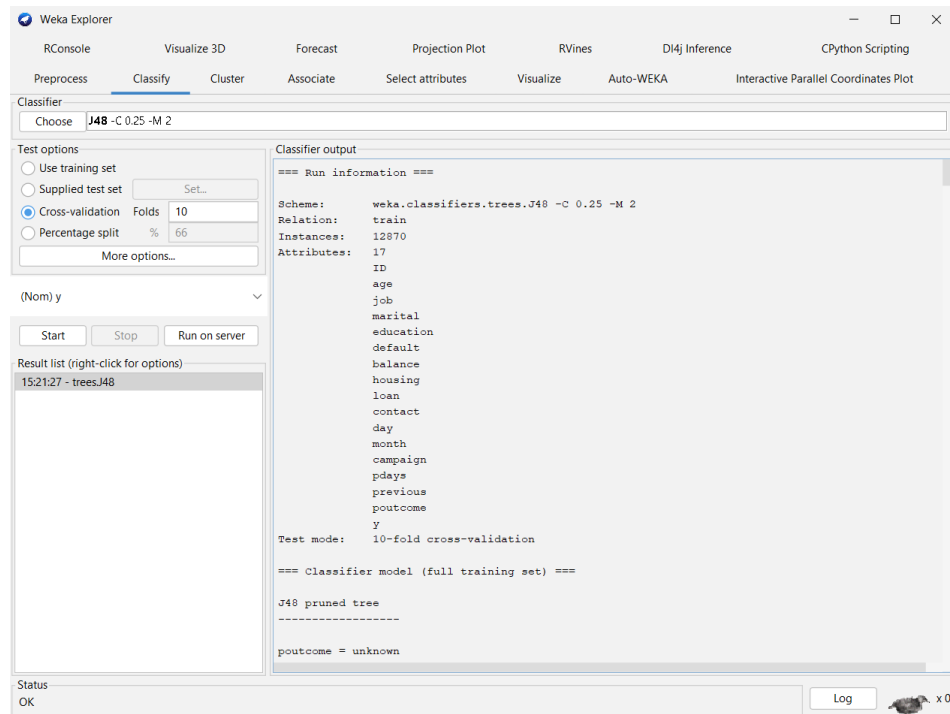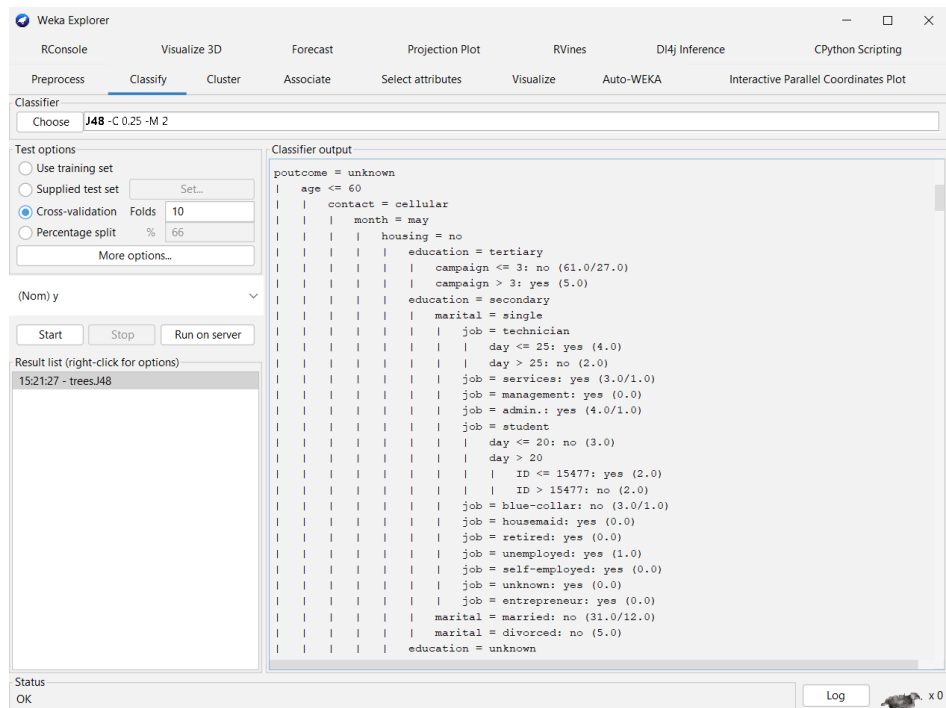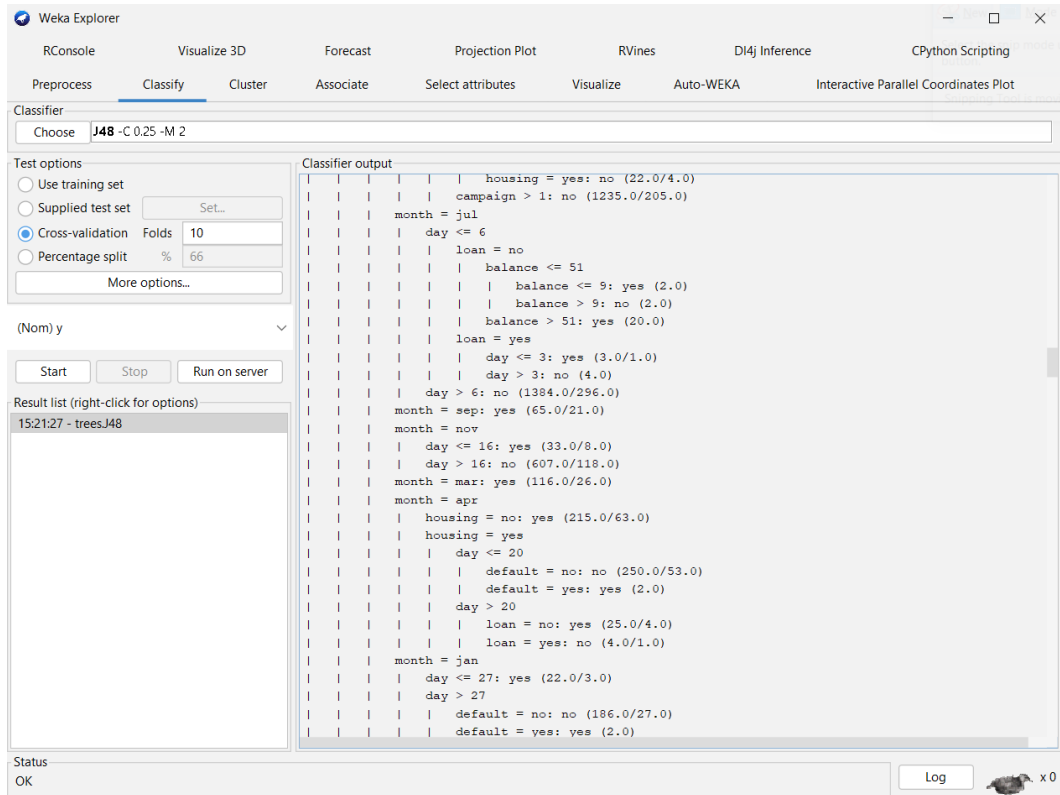
Status
OK

Log    x 0

Figure: Decision Tree Algorithm

Figure: Decision Tree Algorithm



Figure: Decision Tree Algorithm

Figure: Decision Tree Algorithm



Figure: Decision Tree Algorithm
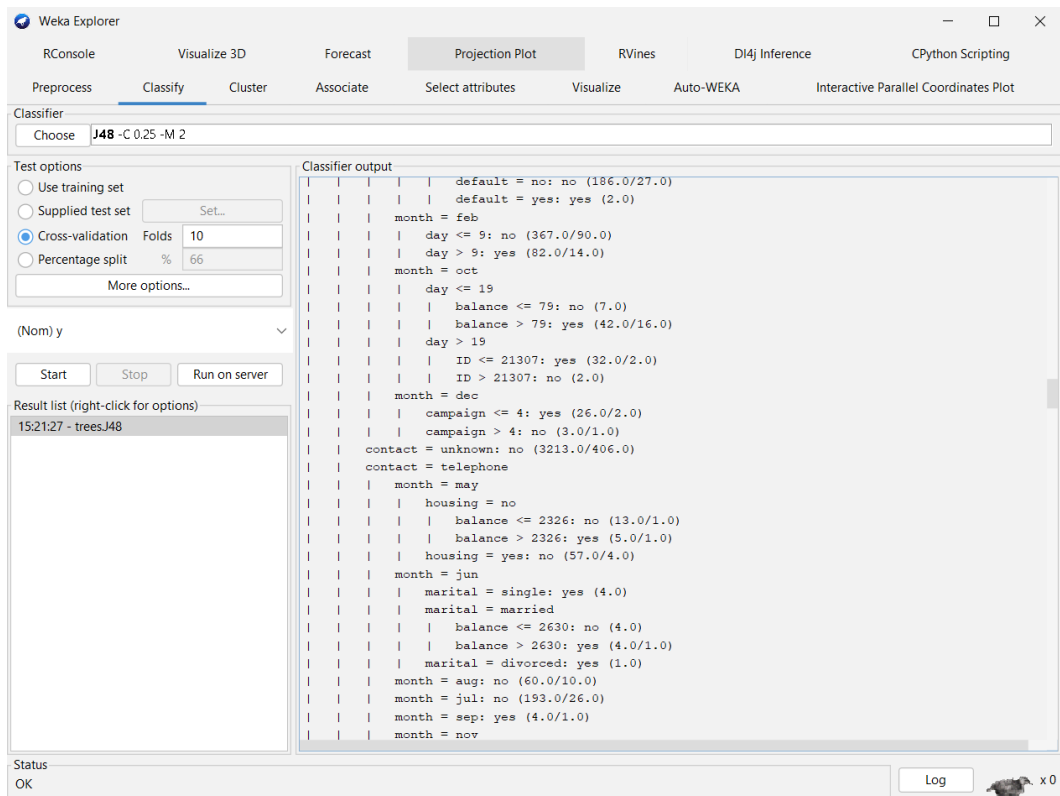
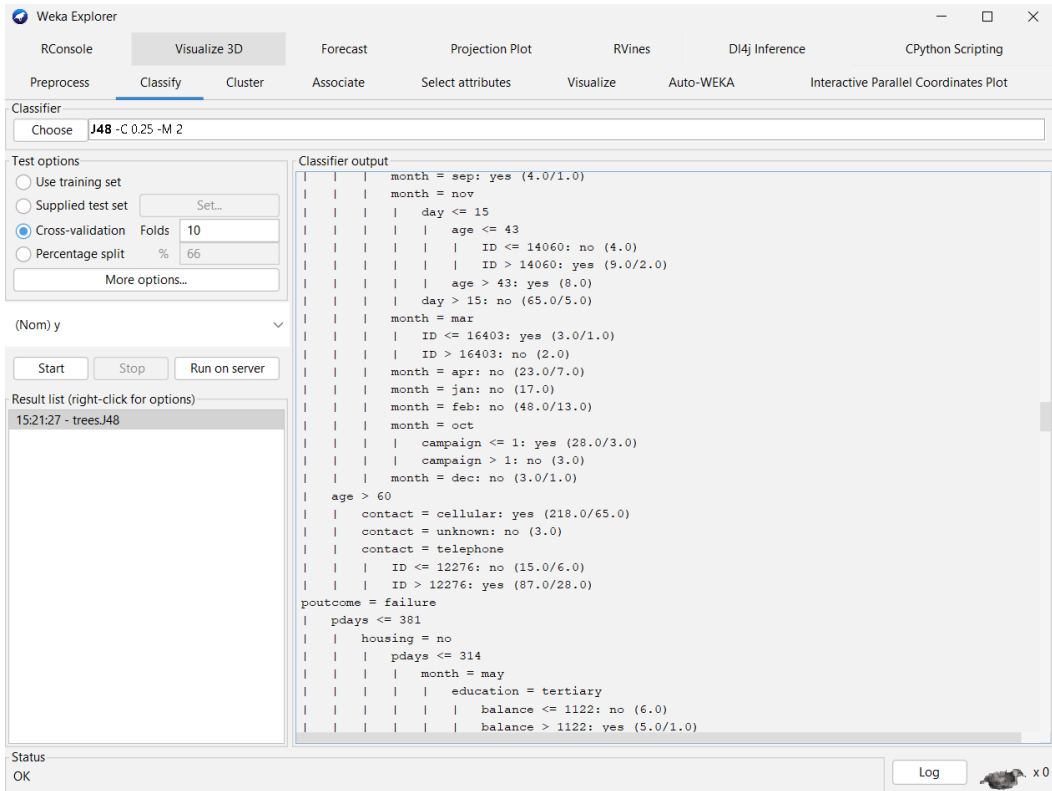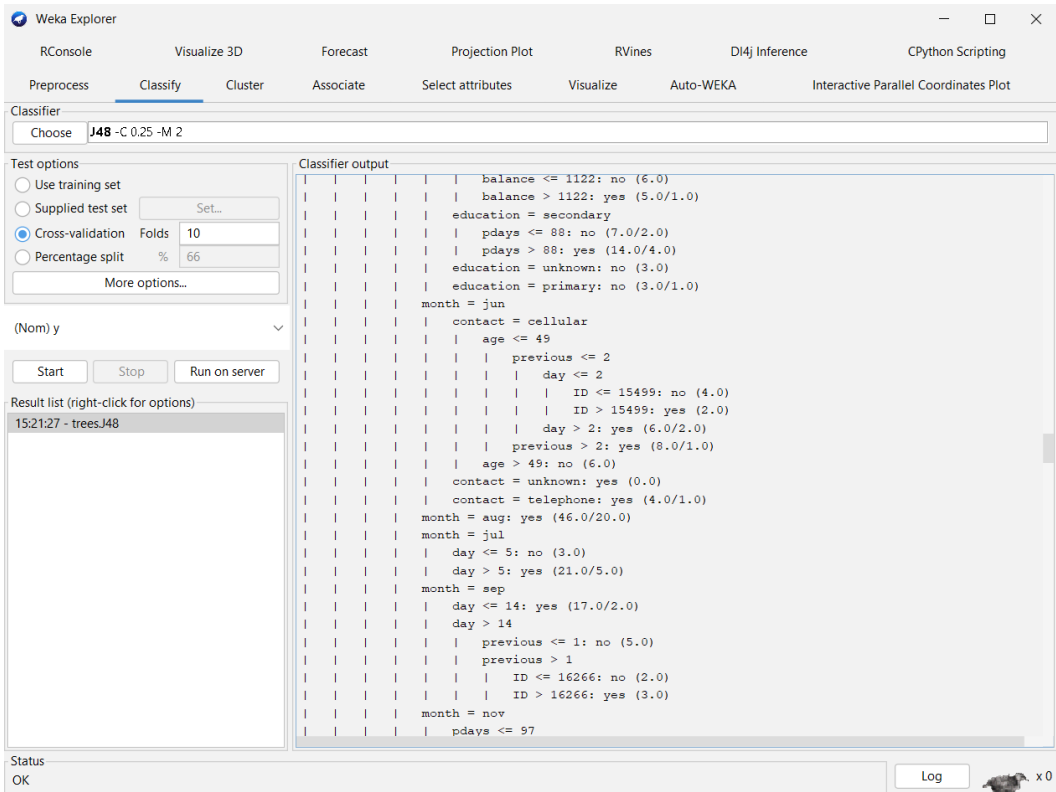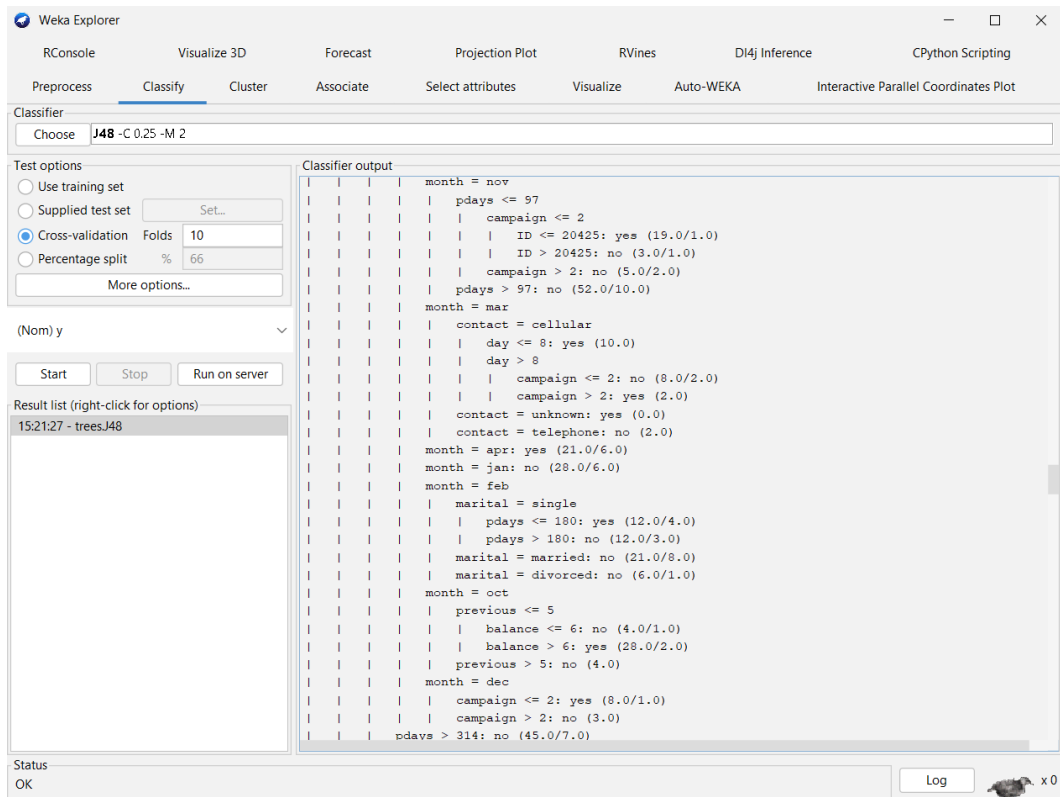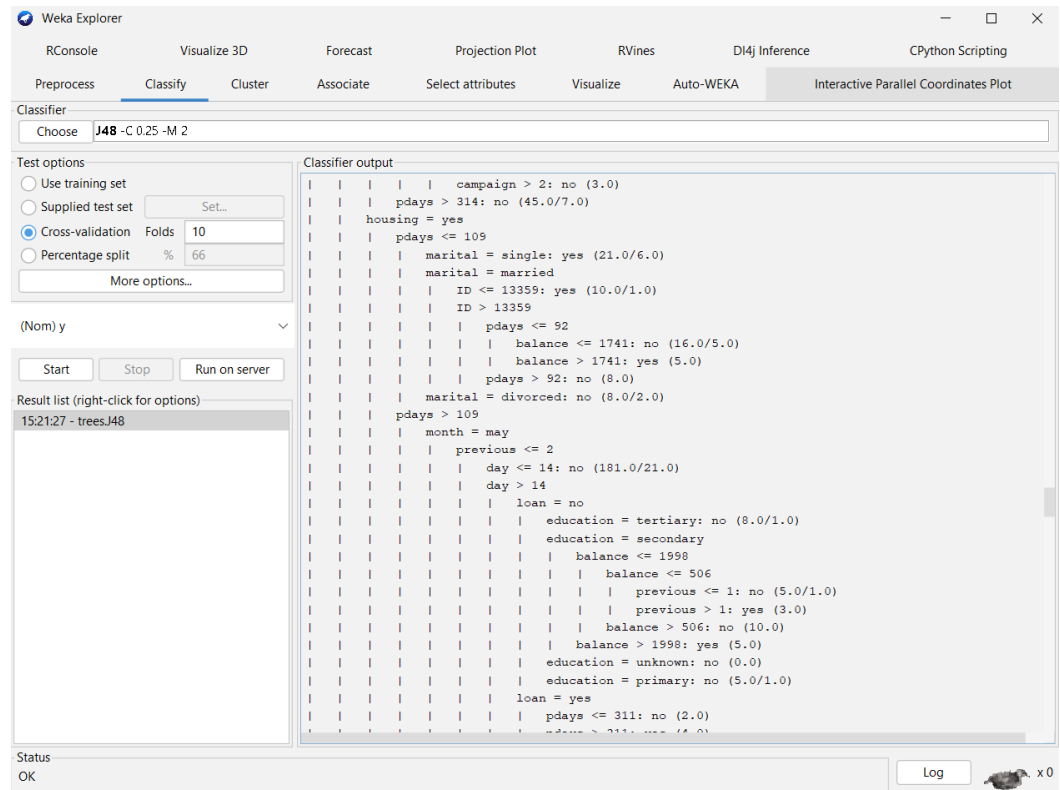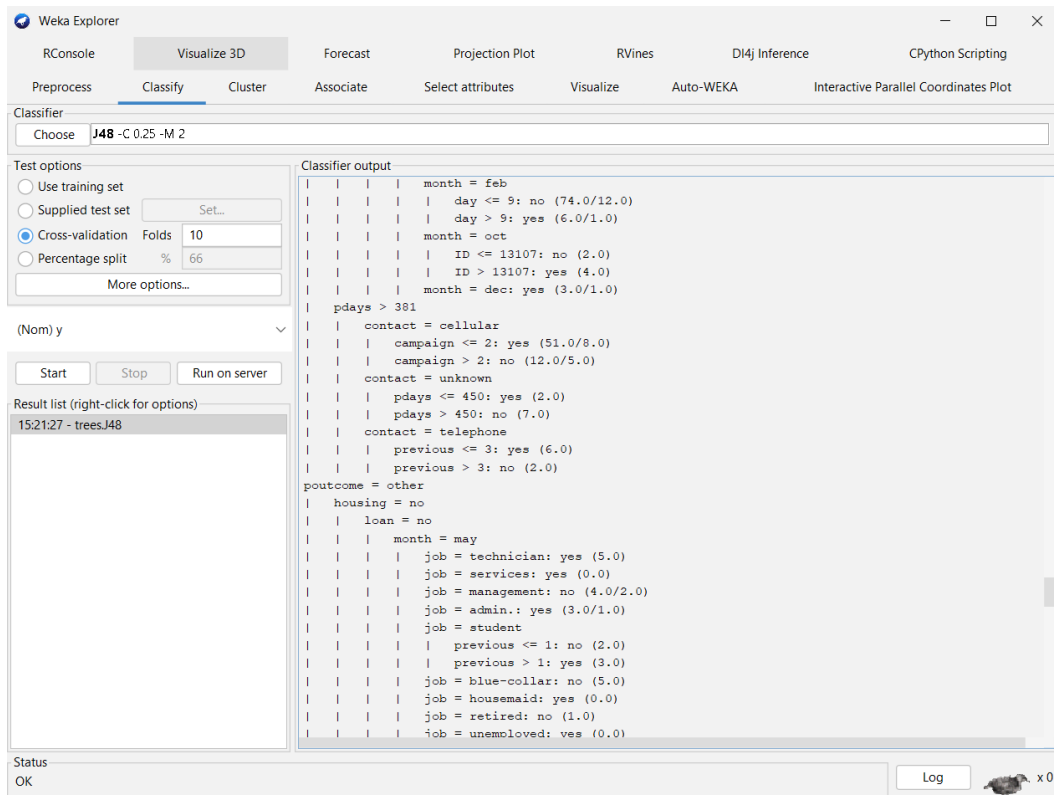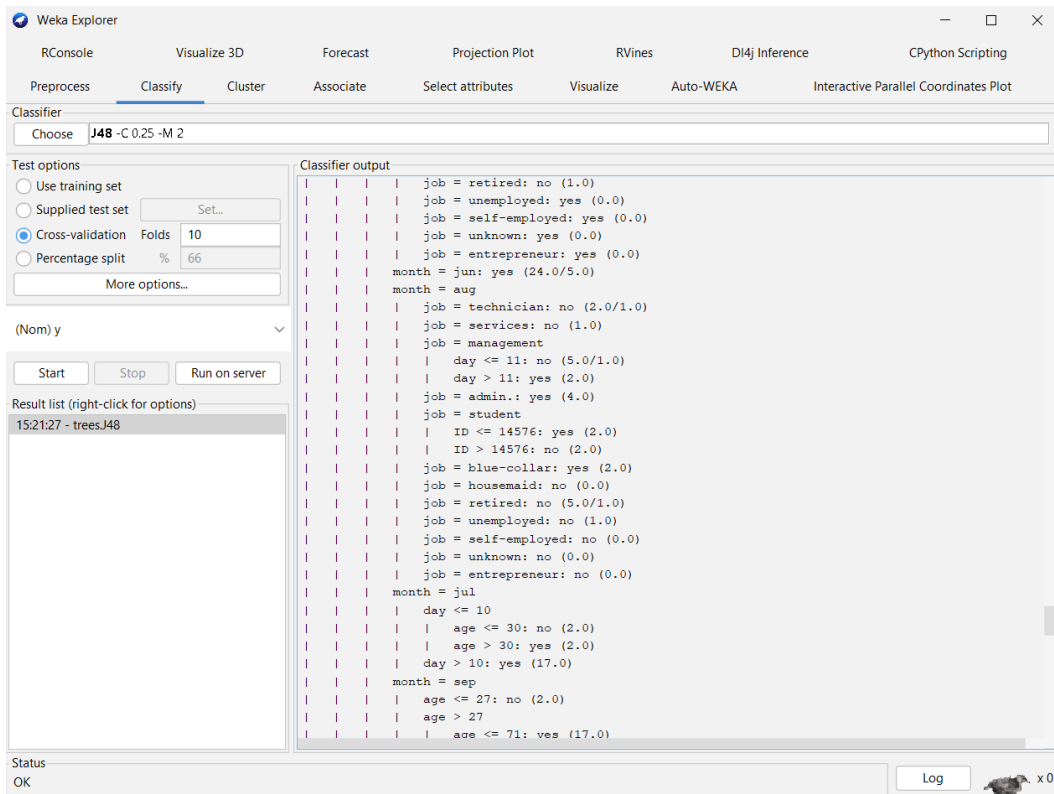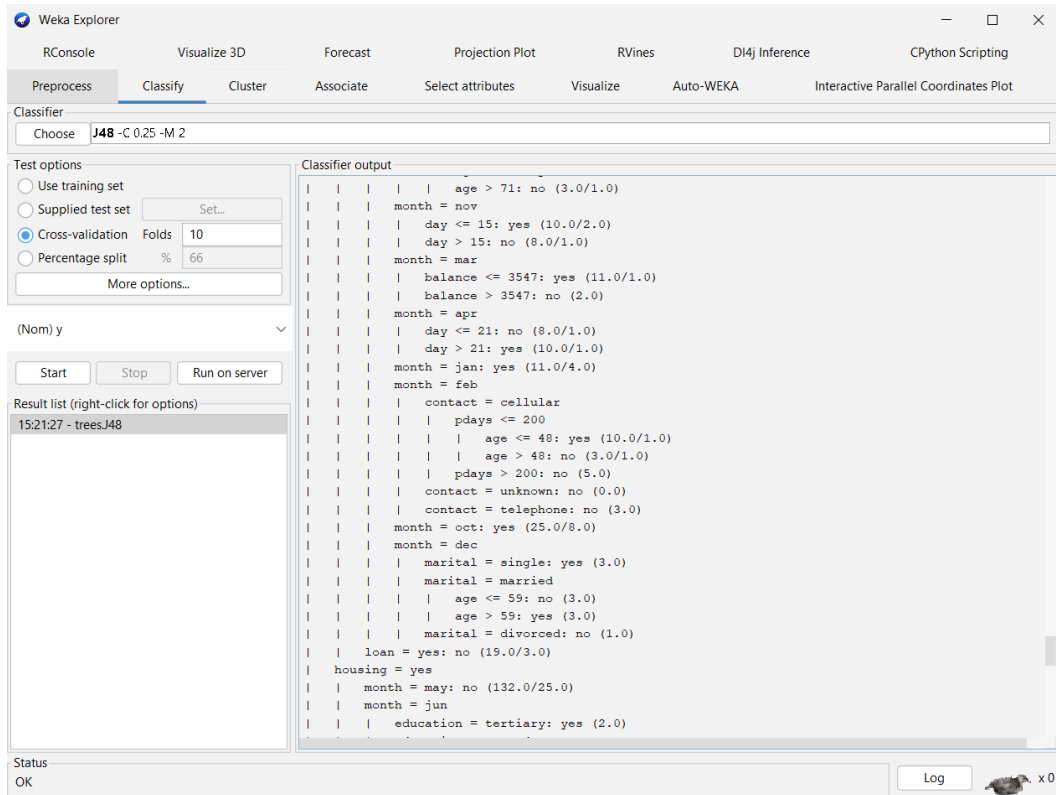Figure: Decision Tree Algorithm



Figure: Decision Tree Algorithm

Figure: Decision Tree Algorithm



Figure: Decision Tree Algorithm

Figure: Decision Tree Algorithm



Figure: Decision Tree Algorithm

## Weka Explorer

RConsole | Visualize 3D | Forecast | Projection Plot | RVines | DI4j Inference | CPython Scripting

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize | Auto-WEKA | Interactive Parallel Coordinates Plot

**Classifier**

Choose | J48 -C 0.25 -M 2

**Test options**
- ○ Use training set
- ○ Supplied test set — Set...
- ● Cross-validation  Folds  10
- ○ Percentage split  %  66

More options...

(Nom) y

Start | Stop | Run on server

Result list (right-click for options)

15:21:27 - trees.J48

**Classifier output**

```
|   |   |   |   | month = feb
|   |   |   |   |   day <= 9: no (74.0/12.0)
|   |   |   |   |   day > 9: yes (6.0/1.0)
|   |   |   |   month = oct
|   |   |   |   |   ID <= 13107: no (2.0)
|   |   |   |   |   ID > 13107: yes (4.0)
|   |   |   |   month = dec: yes (3.0/1.0)
|   |   |   pdays > 381
|   |   |   contact = cellular
|   |   |   |   campaign <= 2: yes (51.0/8.0)
|   |   |   |   campaign > 2: no (12.0/5.0)
|   |   |   contact = unknown
|   |   |   |   pdays <= 450: yes (2.0)
|   |   |   |   pdays > 450: no (7.0)
|   |   |   contact = telephone
|   |   |   |   previous <= 3: yes (6.0)
|   |   |   |   previous > 3: no (2.0)
poutcome = other
|   housing = no
|   |   loan = no
|   |   |   month = may
|   |   |   |   job = technician: yes (5.0)
|   |   |   |   job = services: yes (0.0)
|   |   |   |   job = management: no (4.0/2.0)
|   |   |   |   job = admin.: yes (3.0/1.0)
|   |   |   |   job = student
|   |   |   |   |   previous <= 1: no (2.0)
|   |   |   |   |   previous > 1: yes (3.0)
|   |   |   |   job = blue-collar: no (5.0)
|   |   |   |   job = housemaid: yes (0.0)
|   |   |   |   job = retired: no (1.0)
|   |   |   |   job = unemployed: yes (0.0)
```

**Status**
OK

Log | x 0

Figure: Decision Tree Algorithm

## Weka Explorer

RConsole | Visualize 3D | Forecast | Projection Plot | RVines | DI4j Inference | CPython Scripting

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize | Auto-WEKA | Interactive Parallel Coordinates Plot

**Classifier**

Choose | J48 -C 0.25 -M 2

**Test options**
- ○ Use training set
- ○ Supplied test set — Set...
- ● Cross-validation  Folds  10
- ○ Percentage split  %  66

More options...

(Nom) y

Start | Stop | Run on server

Result list (right-click for options)

15:21:27 - trees.J48

**Classifier output**

```
|   |   |   |   job = retired: no (1.0)
|   |   |   |   job = unemployed: yes (0.0)
|   |   |   |   job = self-employed: yes (0.0)
|   |   |   |   job = unknown: yes (0.0)
|   |   |   |   job = entrepreneur: yes (0.0)
|   |   |   month = jun: yes (24.0/5.0)
|   |   |   month = aug
|   |   |   |   job = technician: no (2.0/1.0)
|   |   |   |   job = services: no (1.0)
|   |   |   |   job = management
|   |   |   |   |   day <= 11: no (5.0/1.0)
|   |   |   |   |   day > 11: yes (2.0)
|   |   |   |   job = admin.: yes (4.0)
|   |   |   |   job = student
|   |   |   |   |   ID <= 14576: yes (2.0)
|   |   |   |   |   ID > 14576: no (2.0)
|   |   |   |   job = blue-collar: yes (2.0)
|   |   |   |   job = housemaid: no (0.0)
|   |   |   |   job = retired: no (5.0/1.0)
|   |   |   |   job = unemployed: no (1.0)
|   |   |   |   job = self-employed: no (0.0)
|   |   |   |   job = unknown: no (0.0)
|   |   |   |   job = entrepreneur: no (0.0)
|   |   |   month = jul
|   |   |   |   day <= 10
|   |   |   |   |   age <= 30: no (2.0)
|   |   |   |   |   age > 30: yes (2.0)
|   |   |   |   day > 10: yes (17.0)
|   |   |   month = sep
|   |   |   |   age <= 27: no (2.0)
|   |   |   |   age > 27
|   |   |   |   |   age <= 71: yes (17.0)
```

**Status**
OK

Log | x 0

Figure: Decision Tree Algorithm

Weka Explorer

RConsole  Visualize 3D  Forecast  Projection Plot  RVines  DI4j Inference  CPython Scripting

Preprocess  Classify  Cluster  Associate  Select attributes  Visualize  Auto-WEKA  Interactive Parallel Coordinates Plot

Classifier

Choose  J48 -C 0.25 -M 2

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

More options...

(Nom) y

Start  Stop  Run on server

Result list (right-click for options)
15:21:27 - trees.J48

Classifier output
```
|   |   |   |   |   |   age > 71: no (3.0/1.0)
|   |   |   month = nov
|   |   |   |   day <= 15: yes (10.0/2.0)
|   |   |   |   day > 15: no (8.0/1.0)
|   |   |   month = mar
|   |   |   |   balance <= 3547: yes (11.0/1.0)
|   |   |   |   balance > 3547: no (2.0)
|   |   |   month = apr
|   |   |   |   day <= 21: no (8.0/1.0)
|   |   |   |   day > 21: yes (10.0/1.0)
|   |   |   month = jan: yes (11.0/4.0)
|   |   |   month = feb
|   |   |   |   contact = cellular
|   |   |   |   |   pdays <= 200
|   |   |   |   |   |   age <= 48: yes (10.0/1.0)
|   |   |   |   |   |   age > 48: no (3.0/1.0)
|   |   |   |   |   pdays > 200: no (5.0)
|   |   |   |   contact = unknown: no (0.0)
|   |   |   |   contact = telephone: no (3.0)
|   |   |   month = oct: yes (25.0/8.0)
|   |   |   month = dec
|   |   |   |   marital = single: yes (3.0)
|   |   |   |   marital = married
|   |   |   |   |   age <= 59: no (3.0)
|   |   |   |   |   age > 59: yes (3.0)
|   |   |   |   marital = divorced: no (1.0)
|   |   loan = yes: no (19.0/3.0)
|   housing = yes
|   |   month = may: no (132.0/25.0)
|   |   month = jun
|   |   |   education = tertiary: yes (2.0)
```

Status
OK

Log  x 0

Figure: Decision Tree Algorithm

Weka Explorer

RConsole  Visualize 3D  Forecast  Projection Plot  RVines  DI4j Inference  CPython Scripting

Preprocess  Classify  Cluster  Associate  Select attributes  Visualize  Auto-WEKA  Interactive Parallel Coordinates Plot

Classifier

Choose  J48 -C 0.25 -M 2

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

More options...

(Nom) y

Start  Stop  Run on server

Result list (right-click for options)
15:21:27 - trees.J48

Classifier output
```
|   |   |   education = secondary
|   |   |   |   pdays <= 89: yes (2.0)
|   |   |   |   pdays > 89: no (3.0)
|   |   |   education = unknown: yes (0.0)
|   |   |   education = primary: yes (0.0)
|   |   month = aug
|   |   |   pdays <= 102: no (2.0)
|   |   |   pdays > 102: yes (8.0)
|   |   month = jul: yes (5.0/2.0)
|   |   month = sep
|   |   |   contact = cellular
|   |   |   |   marital = single: yes (3.0)
|   |   |   |   marital = married: no (2.0)
|   |   |   |   marital = divorced: yes (0.0)
|   |   |   contact = unknown: no (3.0)
|   |   |   contact = telephone: yes (1.0)
|   |   month = nov
|   |   |   pdays <= 192: no (35.0/3.0)
|   |   |   pdays > 192: yes (5.0/1.0)
|   |   month = mar: no (3.0)
|   |   month = apr
|   |   |   loan = no
|   |   |   |   education = tertiary
|   |   |   |   |   previous <= 4
|   |   |   |   |   |   previous <= 1: no (3.0/1.0)
|   |   |   |   |   |   previous > 1: yes (8.0/1.0)
|   |   |   |   |   previous > 4: no (2.0)
|   |   |   |   education = secondary
|   |   |   |   |   pdays <= 342: no (20.0/4.0)
|   |   |   |   |   pdays > 342: yes (2.0)
|   |   |   |   education = unknown: no (2.0/1.0)
```

Status
OK

Log  x 0

Figure: Decision Tree Algorithm

Figure: Decision Tree Algorithm



Figure: Decision Tree Algorithm

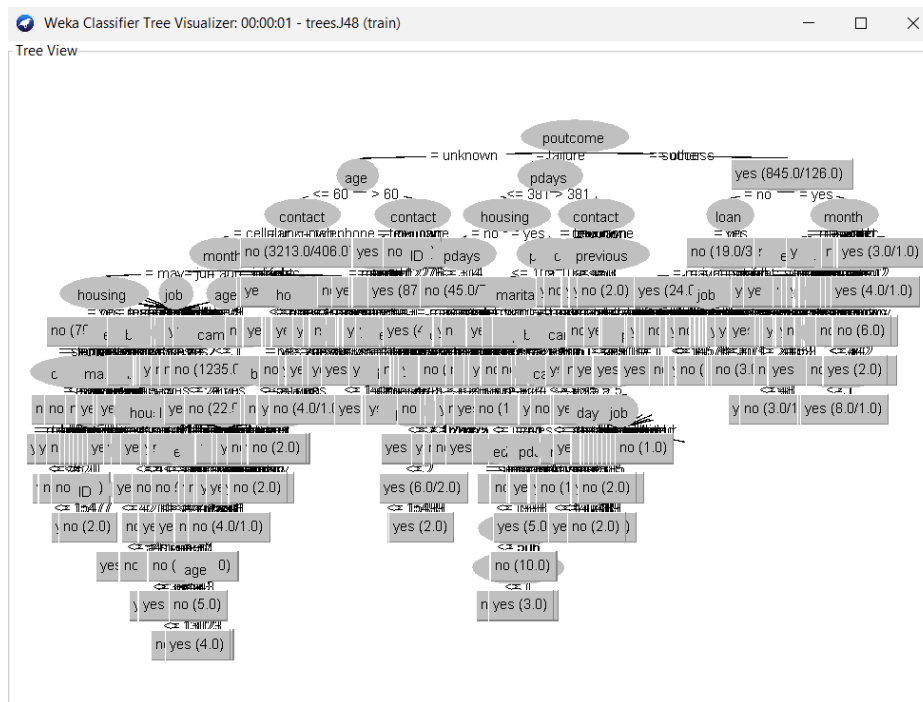2. Then trees.J48 was left click and Visualize tree was selected.



Figure: Decision Tree

## Visualization Plotting:

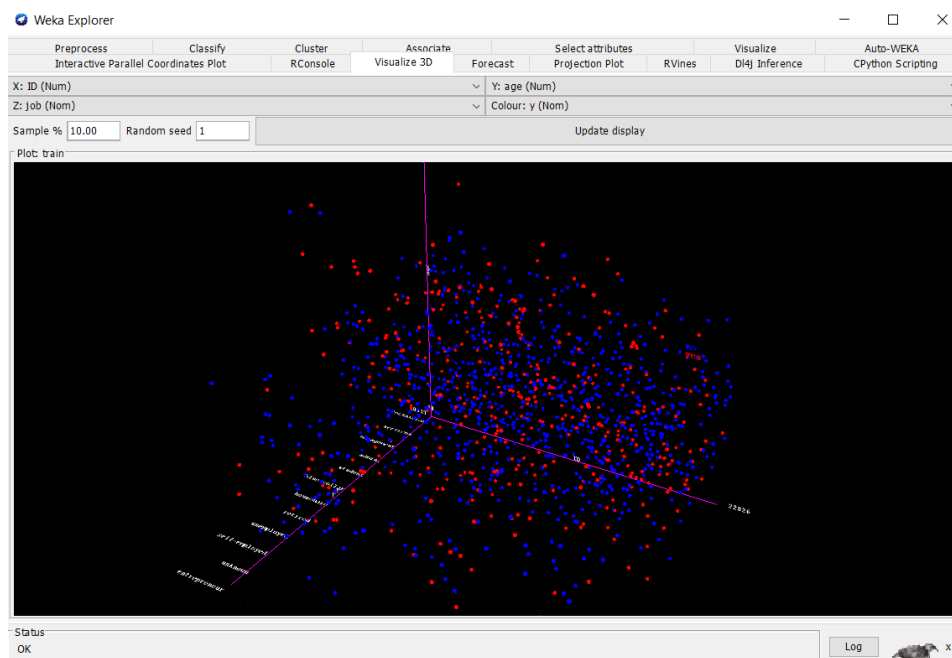1. Firstly, visualization 3D was selected, 'Sample %' was set to10.00% then update display was button was pressed.



Figure: Weka Explorer Naïve bayes classification

# Conclusion & Discussion

## Naïve Bayes classifier:

There are 12870 instances in the Naive Bayes model. In this case, there are 9448 correctly classified instances and 3422 incorrectly classified instances.

```
Correctly Classified Instances        9448              73.411  %
Incorrectly Classified Instances      3422              26.589  %
Kappa statistic                          0.3539
Mean absolute error                      0.3058
Root mean squared error                  0.4529
Relative absolute error                 71.716  %
Root relative squared error             98.0894 %
Total Number of Instances            12870

=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC     ROC Area  PRC Area  Class
                0.836    0.494    0.791      0.836   0.813      0.356   0.736     0.839     no
                0.506    0.164    0.579      0.506   0.540      0.356   0.736     0.549     yes
Weighted Avg.   0.734    0.393    0.726      0.734   0.729      0.356   0.736     0.750

=== Confusion Matrix ===

    a    b    <-- classified as
 7442 1461 |    a = no
 1961 2006 |    b = yes
```

Figure: Naive Bayes model accuracy with confusion matrix

## K-Nearest Neighbour Classification:

In the KNN model, there are 12970 occurrences. In this situation, there are 8850 cases that are correctly classified and 4020 instances that are incorrectly classified.

```
Correctly Classified Instances        8850              68.7646 %
Incorrectly Classified Instances      4020              31.2354 %
Kappa statistic                          0.2608
Mean absolute error                      0.3124
Root mean squared error                  0.5588
Relative absolute error                 73.2501 %
Root relative squared error            121.0223 %
Total Number of Instances            12870

=== Detailed Accuracy By Class ===

                TP Rate  FP Rate  Precision  Recall  F-Measure  MCC     ROC Area  PRC Area  Class
                0.782    0.523    0.770      0.782   0.776      0.261   0.629     0.753     no
                0.477    0.218    0.493      0.477   0.485      0.261   0.629     0.396     yes
Weighted Avg.   0.688    0.429    0.685      0.688   0.686      0.261   0.629     0.643

=== Confusion Matrix ===

    a    b    <-- classified as
 6958 1945 |    a = no
 2075 1892 |    b = yes
```

Figure: Naive KNN accuracy with confusion matrix

## Decision Trees:

There are 12870 instances in the Decision Tree. There are 10058 instances that are correctly classified and 2812 instances that are incorrectly classified in this circumstance.

```
Correctly Classified Instances        10058              78.1507 %
Incorrectly Classified Instances       2812              21.8493 %
Kappa statistic                           0.4424
Mean absolute error                       0.31
Root mean squared error                   0.4163
Relative absolute error                  72.7021 %
Root relative squared error              90.1497 %
Total Number of Instances             12870

=== Detailed Accuracy By Class ===

               TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Class
               0.907    0.501    0.803      0.907   0.852      0.455    0.737     0.810     no
               0.499    0.093    0.706      0.499   0.585      0.455    0.737     0.576     yes
Weighted Avg.  0.782    0.375    0.773      0.782   0.769      0.455    0.737     0.738

=== Confusion Matrix ===

   a    b    <-- classified as
 8078  825 |    a = no
 1987 1980 |    b = yes
```

Figure: Decision Tree model accuracy with confusion matrix

As we can see, the percentage of cases of Naive Bayes that are successfully categorized is 73.41%, the percentage of instances of KNN that are correctly classified is 68.76%, and the percentage of instances of a Decision tree that is correctly classified is 78.15%. We may infer that the Decision Tree performs better in this dataset since it correctly classifies more instances than Naive Bayes and KNN.

Here the purpose of this report was to identify an appropriate classifier for bank marketing response prediction that will be able to categorize the bank marketing response as correctly as possible and be able to forecast the class from the test data set. Following the application of three different classifiers—KNN, naive Bayes, and decision tree the dataset's best classifier with a 78.15% accuracy rate is the decision tree classifier. Then, a training set was taken from the original dataset to create a machine learning model. The model was tested using a prepared test dataset, and the results showed that the model's accuracy was 73.83% for the prepared test dataset. Creating training and the testing dataset is an important concept in data science as it is used to improve generalization and minimize overfitting. This also helps to give an unbiased evaluation of the accuracy of the model itself.