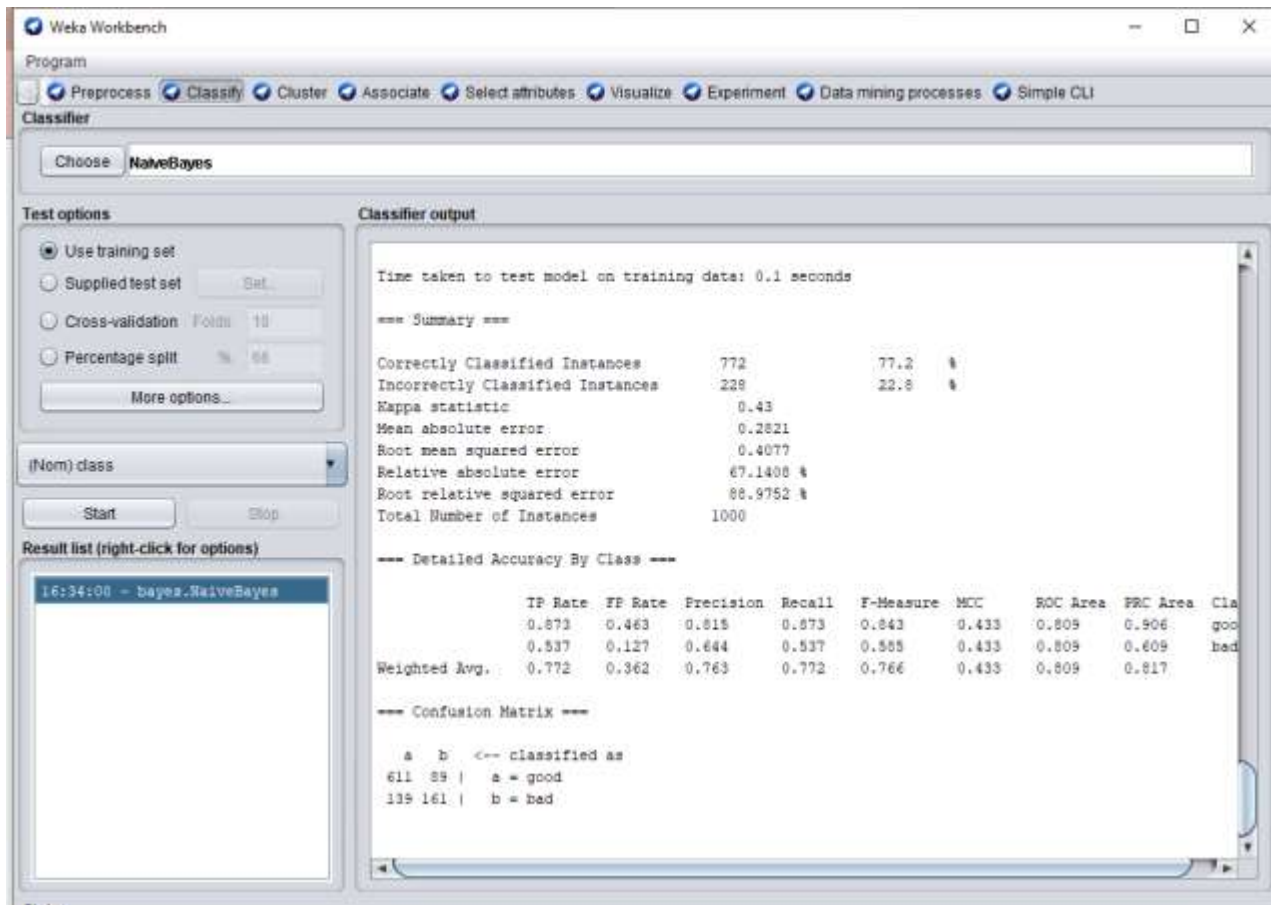# DATA MINING ASSIGNMENT 2

## Naive Bayes Classification

## TASK 1: Train the training dataset using naive bayes and observe classifier output.

PROCEDURE:

1) Open Weka GUI Chooser.

2) Select EXPLORER present in Applications.

3) Select Preprocess Tab.

4) Go to OPEN file and browse the file that is already stored in the system "credit-g.arff".

5) Go to Classify tab.

6) Select by clicking the button choose and select  bayes NaiveBayes.

7) Select Test options "Use training set".

8) Select class attribute.

9) Click Start.

10) Now we can see the output details in the Classifier output.

Program

Preprocess ● Classify ● Cluster ● Associate ● Select attributes ● Visualize ● Experiment ● Data mining processes ● Simple CLI

**Classifier**

Choose NaiveBayes

**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)**

16:34:00 - bayes.NaiveBayes

**Classifier output**

```
Time taken to test model on training data: 0.1 seconds

=== Summary ===

Correctly Classified Instances          772              77.2   %
Incorrectly Classified Instances        228              22.8   %
Kappa statistic                           0.43
Mean absolute error                       0.2821
Root mean squared error                   0.4077
Relative absolute error                  67.1408 %
Root relative squared error              86.9752 %
Total Number of Instances              1000

=== Detailed Accuracy By Class ===

              TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Cla
                0.873    0.463    0.815      0.873   0.843      0.433  0.809     0.906     goo
                0.537    0.127    0.644      0.537   0.585      0.433  0.809     0.609     bad
Weighted Avg.   0.772    0.362    0.763      0.772   0.766      0.433  0.809     0.817

=== Confusion Matrix ===

   a    b   <-- classified as
 611   89 |   a = good
 139  161 |   b = bad
```

---

**TASK 2:** Train a Decision Tree using percentage split and report your results. Increase percentage split by 5% upto 80% starting from 65% and check at which percentage split we are getting the best accuracy.

# 1. When percentage split is 65%, the accuracy is 77.4286%

Weka Workbench — □ ✕

Program

⊘ Preprocess ⊘ Classify ⊘ Cluster ⊘ Associate ⊘ Select attributes ⊘ Visualize ⊘ Experiment ⊘ Data mining processes ⊘ Simple CLI

**Classifier**

Choose | NaiveBayes

**Test options**

- ◯ Use training set
- ◯ Supplied test set    Set..
- ◯ Cross-validation  Folds  10
- ⦿ Percentage split    %  65

More options...

(Nom) class ▾

Start    Stop

**Result list (right-click for options)**

16:34:00 – bayes.NaiveBayes
16:35:53 – bayes.NaiveBayes

**Classifier output**

```
Time taken to test model on test split: 0.01 seconds

=== Summary ===

Correctly Classified Instances         271               77.4286 %
Incorrectly Classified Instances        79               22.5714 %
Kappa statistic                          0.4114
Mean absolute error                      0.2835
Root mean squared error                  0.4033
Relative absolute error                 68.3371 %
Root relative squared error             90.8424 %
Total Number of Instances              350

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Cla
                 0.853    0.446    0.843      0.853   0.848      0.411  0.796     0.918     goo
                 0.554    0.147    0.573      0.554   0.564      0.411  0.796     0.582     bad
Weighted Avg.    0.774    0.367    0.772      0.774   0.773      0.411  0.796     0.830

=== Confusion Matrix ===

   a   b   <-- classified as
 220  38 |   a = good
  41  51 |   b = bad
```

## 2. When percentage split is 70%, the accuracy is 75.3333%



Weka Workbench screenshot showing NaiveBayes classifier output:

```
Time taken to test model on test split: 0.01 seconds

=== Summary ===

Correctly Classified Instances         226              75.3333 %
Incorrectly Classified Instances        74              24.6667 %
Kappa statistic                          0.3537
Mean absolute error                      0.2851
Root mean squared error                  0.4116
Relative absolute error                 69.0347 %
Root relative squared error             92.7794 %
Total Number of Instances              300

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Cla
                 0.842    0.494    0.827      0.842   0.834      0.354  0.788     0.916     goo
                 0.506    0.158    0.533      0.506   0.519      0.354  0.788     0.547     bad
Weighted Avg.    0.753    0.405    0.749      0.753   0.751      0.354  0.788     0.819

=== Confusion Matrix ===

   a   b   <-- classified as
 186  35 |   a = good
  39  40 |   b = bad
```

## 3. When percentage split is 75%, the accuracy is 76.8%



```
Weka Workbench                                                           –   □   X

Program

  ○ Preprocess  ○ Classify  ○ Cluster  ○ Associate  ○ Select attributes  ○ Visualize  ○ Experiment  ○ Data mining processes  ○ Simple CLI
Classifier

  Choose   NaiveBayes

Test options                        Classifier output

  ○ Use training set                  Time taken to test model on test split: 0 seconds
  ○ Supplied test set      Set...
                                      === Summary ===
  ○ Cross-validation  Folds  10
                                      Correctly Classified Instances         192              76.8   %
  ● Percentage split   %  75          Incorrectly Classified Instances        58              23.2   %
                                      Kappa statistic                          0.403
         More options...              Mean absolute error                      0.2778
                                      Root mean squared error                  0.4029
                                      Relative absolute error                 67.5042 %
  (Nom) class                         Root relative squared error             90.8443 %
                                      Total Number of Instances              250
       Start            Stop
                                      === Detailed Accuracy By Class ===
Result list (right-click for options)
                                                  TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Cla
 16:34:00 - bayes.NaiveBayes                       0.842    0.439    0.842     0.842    0.842     0.403   0.806     0.924     goo
 16:35:53 - bayes.NaiveBayes                       0.561    0.158    0.561     0.561    0.561     0.403   0.806     0.567     bad
 16:36:04 - bayes.NaiveBayes          Weighted Avg.  0.768    0.365    0.768     0.768    0.768     0.403   0.806     0.830
 16:36:13 - bayes.NaiveBayes
                                      === Confusion Matrix ===

                                        a   b   <-- classified as
                                      155  29 |   a = good
                                       29  37 |   b = bad
```

4. When percentage split is 80%, the accuracy is 74.5%



CONCLUSION: When the percentage split is 65%, the accuracy is high(77.4286%).

# TASK 3: Train a Decision Tree using cross validation and report your results.

1. When cross validation folds : 10, accuracy is 75.4%

## 2. When cross validation folds : 8, accuracy is 75.9%

## 3. When cross validation folds : 6, accuracy is 75.4%



CONCLUSION: The accuracy is high(75.9%) when cross validation folds: 8