**Title:** Finding best Classifier and evaluating performance using test and training set for Glass Identification dataset

**Objective:**

In this Project, I will Choose five classifier algorithms and apply them into a dataset referred with http://archive.ics.uci.edu/ml/datasets/Glass+Identification . And also, I will find the ROC curve to make a precise decision that which is the best among the selected classifiers. Additionally, we will split our whole dataset into two segments. One is test set and another is training set. I will run this test data set with the training data set and trying to evaluate the performance of the test data. All the task will be done by Weka (a machine learning tool).

**Dataset description:**

In the dataset there are 214 instances and 11 attributes

Attribute Information:

| Attributes | values |
|---|---|
| Id: instances number | 1-214 |
| RI: refractive index | Numeric |
| Na: Sodium | Numeric |
| Mg: Magnesium | Numeric |
| Al: Aluminum | Numeric |
| Si: Silicon | Numeric |
| K: Potassium | Numeric |
| Ca: Calcium | Numeric |
| Ba: Barium | Numeric |
| Fe: Iron | Numeric |
| Type of glass: (class attribute) | 1: building_windows_float_processed (70 out of 214) |
| | 2: building_windows_non_float_processed (76 out of 214) |
| | 3: vehicle_windows_float_processed (17 out of 214) |
| | 4: vehicle_windows_non_float_processed (0 out of 214) |
| | 5: containers (13 out of 214) |
| | 6: tableware (9 out of 214) |
| | 7: headlamps (29 out of 214) |

Note: There is no missing value of attributes

**Output:** We select five Classification Algorithms and apply them into the mentioned training set. All the results are as follows:-

## 1. NaiveBayes

NaiveBayes Classifier works on Bayes theorem of probability to predict the class of unknown data sets. Here is output of it.



## 2.Logistic

Logistic regression is a statistical model that in its basic form uses a logistic function to model a binary dependent variable, although many more complex extensions exist. Here is output of it.

## 3.Kstar

K* is an instance-based classifier, that is the class of a test instance is based upon the class of those training instances similar to it, as determined by some similarity function. It differs from other instance-based learners in that it uses an entropy-based distance function. Here is output of it.

Weka Explorer

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

**Classifier**

Choose | KStar -B 20 -M a

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

More options...

(Nom) type

Start | Stop

**Result list (right-click for options)**
- 23:53:36 - bayes.NaiveBayes
- 00:27:39 - functions.Logistic
- 00:27:59 - lazy.KStar

**Classifier output**

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

Correctly Classified Instances         192               89.7196 %
Incorrectly Classified Instances        22               10.2804 %
Kappa statistic                          0.8612
Mean absolute error                      0.0367
Root mean squared error                  0.1724
Relative absolute error                 14.8641 %
Root relative squared error             49.174  %
Total Number of Instances              214


=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 0.971    0.042    0.919      0.971   0.944      0.917  0.997     0.994     1
                 0.855    0.036    0.929      0.855   0.890      0.835  0.984     0.972     2
                 0.941    0.041    0.667      0.941   0.780      0.772  0.981     0.848     3
                 0.769    0.010    0.833      0.769   0.800      0.788  0.986     0.885     5
                 1.000    0.000    1.000      1.000   1.000      1.000  1.000     1.000     6
                 0.828    0.005    0.960      0.828   0.889      0.876  0.964     0.947     7
Weighted Avg.    0.897    0.031    0.906      0.897   0.898      0.867  0.986     0.962


=== Confusion Matrix ===

  a  b  c  d  e  f   <-- classified as
 68  1  1  0  0  0 |  a = 1
  6 65  4  1  0  0 |  b = 2
  0  1 16  0  0  0 |  c = 3
  0  2  0 10  0  1 |  d = 5
  0  0  0  0  9  0 |  e = 6
  0  1  3  1  0 24 |  f = 7
```
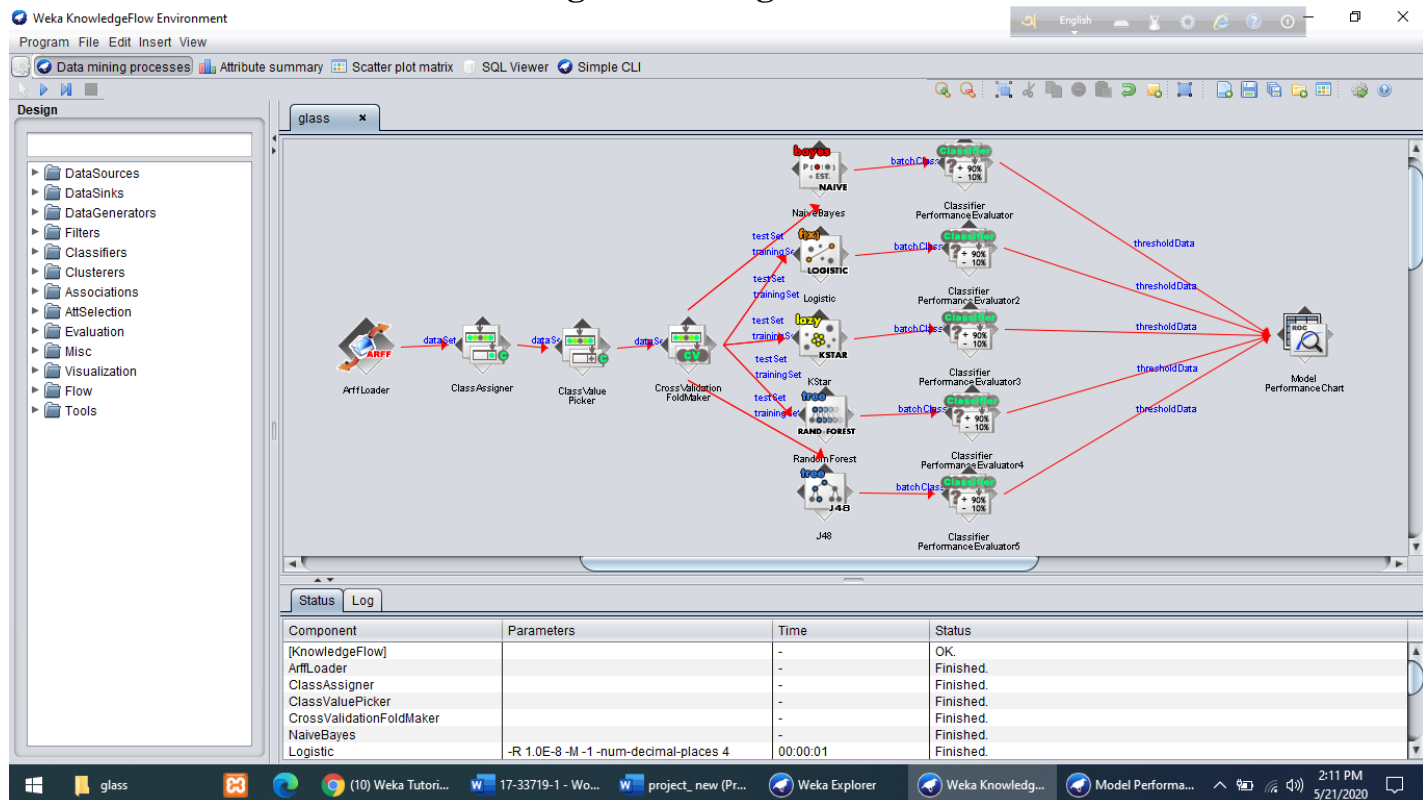
**Status**
OK

Log

glass | report - ashiqulhoq... | project_new (Prote... | 17-33719-1 - Word | Weka Explorer | 12:28 AM 5/21/2020

## 4.RandomForest

RandomForest Classifier Works like 'Decision tree' but the difference is it split tree randomly rather than traditional split of tree. Here is output of it.

## 5.J48

J48 algorithm is one of the best machine learning algorithms to examine the data categorically and continuously. Here is output of it.

## Performance Analysis:

| Classifiers | Correctly Classified instance | Total number of instances |
|---|---|---|
| NaiveBayes | 180 | |
| Logistic | 199 | |
| Kstar | 192 | 214 |
| **RandomForest** | **212** | |
| J48 | 207 | |

Here we see, among 214 instances Random Forest Correctly Classified 212 instance which is 99.0654 % and it's the highest percentage among the other Classifiers. So 'RandomForest' is the best classifier for the dataset.

## Knowledge Flow Diagram

# ROC Curve



This ROC curve is merged of all algorithm that we used, considering class value 2. We know that the curve goes nearer to the y-axis claimed as the better classifier. From the ROC curve we see that curve of RandomForest is nearer than the other curves. So, we can say that 'RandomForest' is best for the dataset.

## Additional task

In additionally, we split data into two segments. I take 70% of original data as training set(148 out of 214 instances) and 30% of data as test set(61 out of 214 instances). This task is done by weka, and run this test data set with the training data set. Here we evaluate the performance of test data by visualizing the classifiers error.

## Training set



## Run with test data



From the Classifiers output, I determined that the performance is 96.7213 %. It Classified 59 instances correctly out of 61 test data set.

# Classifiers error visualization