



Vidyavardhini's College of Engineering and Technology
Department of Artificial Intelligence & Data Science

Name:	Vinith Shetty
Roll No:	52
Class/Sem:	TE/V
Experiment No.:	4
Title:	Using open source tools Implement Classifiers
Date of Performance:	
Date of Submission:	
Marks:	
Sign of Faculty:	



Aim: To implement Naïve Bayes Classifier using open-source tool WEKA.

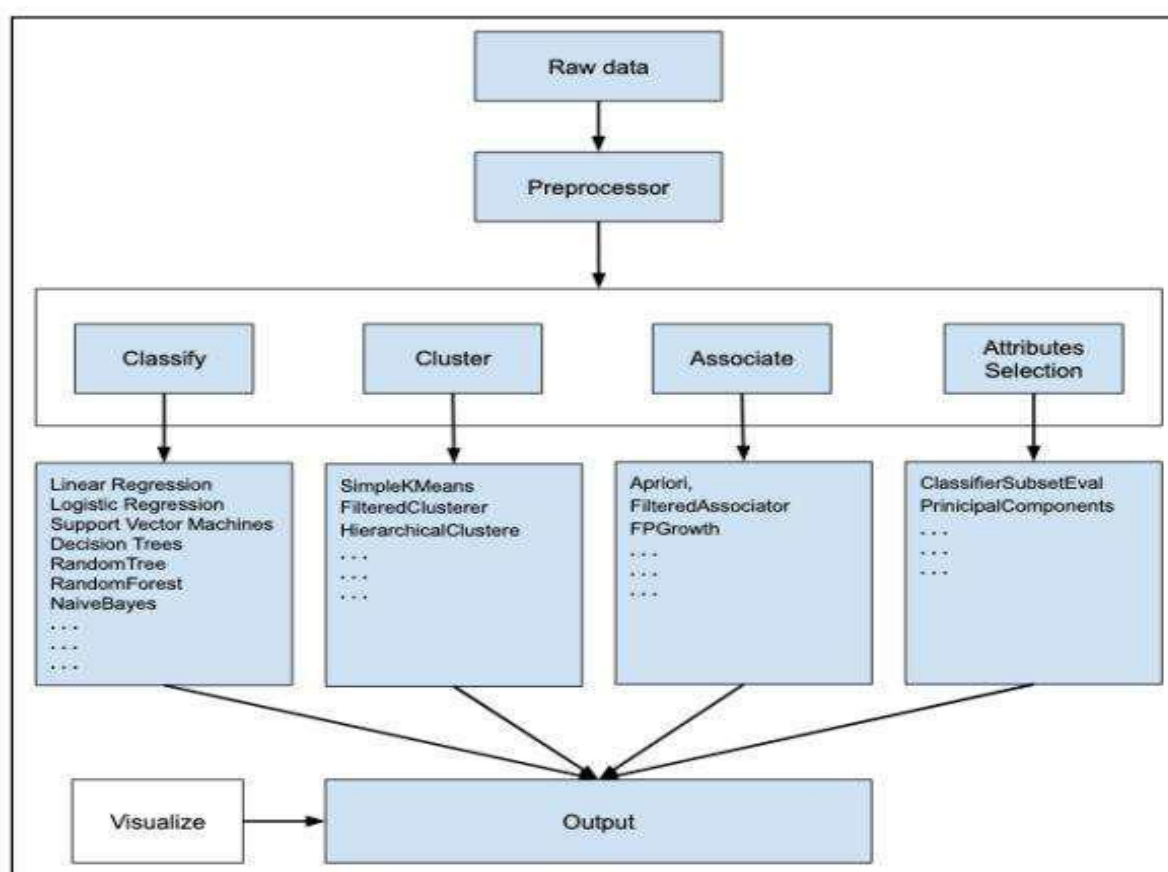
Objective: To make students well versed with open source tool like WEKA to implement Naïve Bayes Classifier.

Theory:

Classification is a data mining function that assigns items in a collection to target categories or classes. The goal of classification is to accurately predict the target class for each case in the data. For example, a classification model could be used to identify loan applicants as low, medium, or high credit risks.

WEKA:

WEKA – an open-source software provides tools for data preprocessing, implementation of several data Mining algorithms, and visualization tools so that you can develop data mining techniques and apply them to real-world data mining problems. Weka is summarized in the following diagram:



First, you will start with the raw data collected from the field. This data may contain several null values and irrelevant fields. You use the data preprocessing tools provided in WEKA to cleanse the data. Then, you would save the preprocessed data in your local storage for applying Data Mining algorithms.



Vidyavardhini's College of Engineering and Technology

Department of Artificial Intelligence & Data Science

Next, depending on the kind of Data Mining model that you are trying to develop you would select one of the options such as Classify, Cluster, or Associate. The Attributes Selection allows the automatic selection of features to create a reduced dataset. Note that under each category, WEKA provides the implementation of several algorithms. You would select an algorithm of your choice, set the desired parameters and run it on the dataset. Then, WEKA would give you the statistical output of the model processing. It provides you a visualization tool to inspect the data. The various models can be applied on the same dataset. You can then compare the outputs of different models and select the best that meets your purpose.

Output:

The screenshot shows the WEKA Classifier output window. On the left, the 'Test options' panel is visible with 'Cross-validation' selected (10 folds). The 'Result list' on the left shows '11:33:38 - bayes.NaiveBayes' selected. The main 'Classifier output' pane displays the following results:

```
[total] 11.0 7.0

Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===
Correctly Classified Instances      9      64.2857 %
Incorrectly Classified Instances    5      35.7143 %
Kappa statistic                    0.1026
Mean absolute error                 0.4649
Root mean squared error             0.543
Relative absolute error             97.6254 %
Root relative squared error         110.051 %
Total Number of Instances          14

=== Detailed Accuracy By Class ===
               TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Class
0.889    0.800    0.667    0.889    0.762    0.122    0.444    0.633    yes
0.200    0.111    0.500    0.200    0.286    0.122    0.444    0.397    no
Weighted Avg.   0.643    0.554    0.607    0.643    0.592    0.122    0.444    0.548

=== Confusion Matrix ===
a b <-- classified as
8 1 | a = yes
4 1 | b = no
```

Conclusion:

What performance metrics were used to evaluate the Naïve Bayes classifier in WEKA?

The Naïve Bayes classifier in WEKA was evaluated using several key performance metrics. The classifier correctly classified 64.29% of instances, while 35.71% were incorrectly classified. The Kappa statistic, which measures the agreement between predicted and actual classifications beyond chance, was 0.1026. Additionally, the mean absolute error, a measure of the average magnitude of errors in predictions, was 0.4649, and the root mean squared error was 0.543, indicating the overall prediction error. The relative absolute error was 97.63%, while the root relative squared error was 110.05%. These error metrics provide insights into how well the model generalizes across the dataset. A breakdown of detailed accuracy by class shows that for the "yes" class, the true positive (TP) rate was 0.889, and for the "no" class, it was 0.200. Precision, which measures the proportion of correct positive predictions, was 0.667 for "yes" and 0.500 for "no." Recall, the model's ability to correctly identify actual positives, was 0.889 for "yes" and 0.200 for "no." The weighted averages of precision, recall, and other metrics across both classes were also provided. These metrics offer a comprehensive view of the model's classification performance.