# CSE477(1)-Project Report

## Comparative analysis of different classification algorithms

Submitted to:

Dr. Mohammad Rezwanul Huq
Assistant Professor, Dept. of CSE, EWU

Submitted by:

Anamika Das Mou (2015-2-60-080)

# Introduction

In modern life, Data have increased rapidly and using those data we can measure the prediction and we also can make a proper decision to execute some works where result may use for business purpose or may be used for disease or other sectors. From the large number of dataset, we can find patterns and patterns help to make decisions. The data analysis occurs by data mining. The Data mining is the process of sorting through large data sets to identify patterns and establish relationships to solve problems through data analysis. Data mining tools allow enterprises to predict future trends.

Here, we worked with one dataset. For the python implementation, we had used dataset and for the Weka tools we had also used the same dataset.
Dataset**[who_suicide.csv]** about suicide information, here have 6 attributes. Using those attributes, we have to measure in which range age the suicide will be occurred.

For this dataset we used Naïve Bayes Classifier, Decision Tree Classifier, k Nearest-Neighbour Classifier, Random Forest tree classifier. We also used the same algorithms for Weka tools.

# Classification Algorithm

We have used four algorithms to measure the result.
1. Decision Tree
2. Naïve Bayes
3. Random forest
4. KNN

**Decision Tree:** A decision tree is a decision support tool that uses a tree-like model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. It is one way to display an algorithm that only contains conditional control statements. Decision trees
are commonly used in operations research, specifically in decision analysis, to help identify a strategy most likely to reach a goal. A decision tree is a structure that includes a root node, branches, and leaf nodes. Each internal node denotes a test on an attribute, each branch denotes the outcome of a test, and each leaf node holds a class label. The topmost node in the tree is the root node. This algorithm has the ability to handle binary and multiclass classification problem. We need two different entropies for calculating.one is for total database that means for all the attribute and the equation for that is,

$\text{Info}(F) = -\sum_{m\ i=1} p_i \log_2(p_i)$

$p_i$ is the probability that an arbitrary tuple in F belongs to class Yes or No and another is for each attribute and equation for that is,

Info$_{att}$(F) = $\sum$ |F$_j$|

$v\,j=1\,|F| \times info(F)$

After that total gain will be calculated. Here, the formula,

Gain(attribute) = Info(F) – Info$_{att}$(F)

**Naïve Bayes:** In machine learning, naive Bayes classifiers are a family of simple "probabilistic classifiers" based on applying Bayes' theorem with strong (naive) independence assumptions between the features. Naive Bayes classifiers are highly scalable, requiring a number of parameters
linear in the number of variables (features/predictors) in a learning problem. Maximum-likelihood
training can be done by evaluating a closed-form expression, which takes linear time, rather than by expensive iterative approximation as used for many other types of classifiers.

By using the Bayes theorem, posterior probability can be calculated, $P(c|x)$, from $P(c)$, $P(x)$, and $P(x|c)$. The Naive Bayes classification consider the predictor value (x) on a given category (c) is independent which is conditional independence. So, we can formulate the theorem from Bayesian.

Here,

$P(c|x_i) = p(x_1|c) \times \ldots \times P(x_n|c) \times P(c)$

$P(c|xi)$ is the posterior probability of a given predictor (attribute). $P(c)$ is the prior probability of class. $P(xi|c)$ is the likelihood. $P(xi)$ is the prior probability of the predictor.

**Random forest:**

Random forests or random decision forests are an ensemble learning method for classification, regression and other tasks that operate by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes (classification) or mean prediction (regression) of the individual trees. Random decision forests correct for decision trees' habit of overfitting to their training set. To say it in simple words: Random forest builds multiple
decision trees and merges them together to get a more accurate and stable prediction. One big advantage of random forest is, that it can be used for both classification and regression problems, which form the majority of current machine learning systems.

Random Forest has nearly the same hyper parameters as a decision tree or a bagging classifier. Fortunately, we don't have to combine a decision tree with a bagging classifier and can just easily use the classifier-class of Random Forest. Like I already said, with Random Forest, you can also
deal with Regression tasks by using the Random Forest repressor.

Random Forest adds additional randomness to the model, while growing the trees. Instead of searching for the most important feature while splitting a node, it searches for the best feature among a random subset of features. This results in a wide diversity that generally results in a better
model.

Therefore, in Random Forest, only a random subset of the features is taken into consideration by the algorithm for splitting a node. We can even make trees more random, by additionally using random thresholds for each feature rather than searching for the best possible thresholds like a normal decision tree does. This is:
1. Increasing the Predictive Power.
2. Increasing the Model's Speed.

k Nearest-Neighbour Classifier:

K-Nearest Neighbors is one of the simplest algorithms used in Machine Learning for regression and classification problem. KNN algorithms use a data and classify new data points based on a similarity measures. It works based on minimum distance from the query instance to the training samples to determine the K-nearest neighbors. The data for KNN algorithm consist of several multivariate attributes name that will be used to classify. KNN is a lazy learning algorithm.

## Dataset

Firstly, we had chosen the **"who_suicide"** dataset for the python implementation and also for Weka tools. Here have
6 attributes and it's with many null values. In this dataset, there are 43777 data. This dataset is a suicide information list. The list has briefly shows about the suicide tendency. At the end this data help to
predict what range age is particularly suicidal.

## Implementation

The configuration of my Pc is:

| System | |
| --- | --- |
| Processor: | Intel(R) Core(TM) i5-8250U CPU @ 1.60GHz 1.80 GHz |
| Installed memory (RAM): | 8.00 GB |
| System type: | 64-bit Operating System, x64-based processor |
| Pen and Touch: | No Pen or Touch Input is available for this Display |

**Weka Tools:** In below, there have all attribute figures.
In the weka tool, we have used 75% data for the training and 25% data for the testing set. So in below we attached different method measures to observe the different performance.
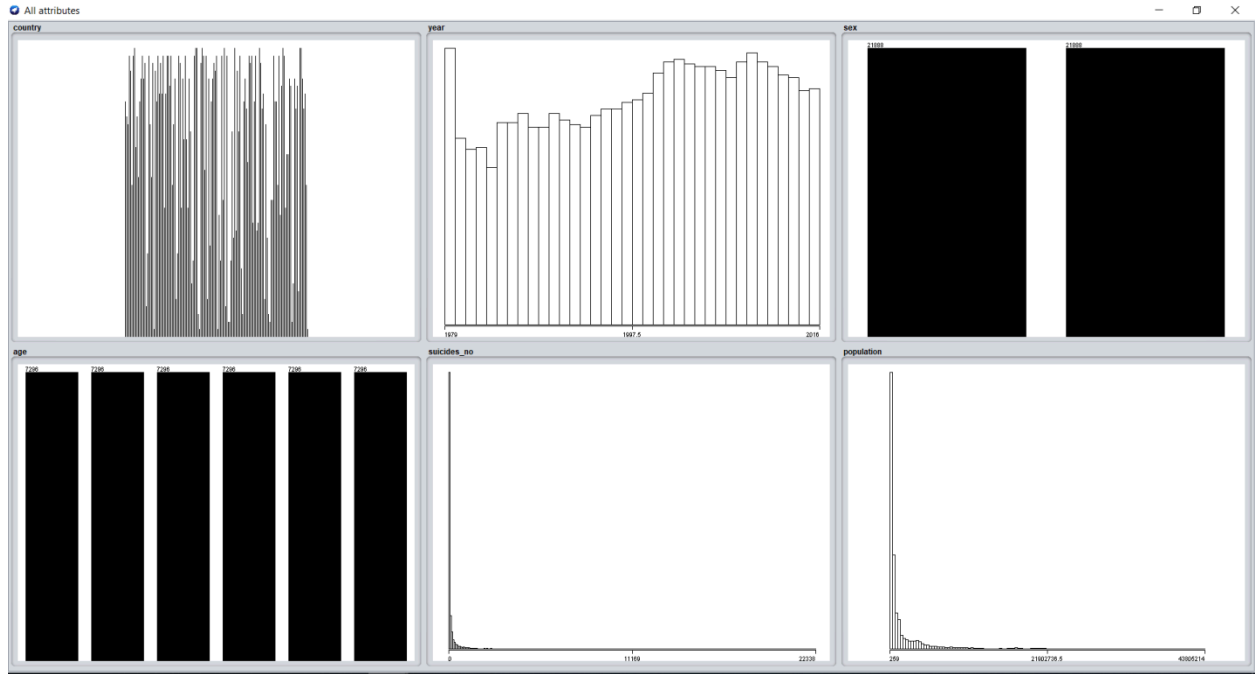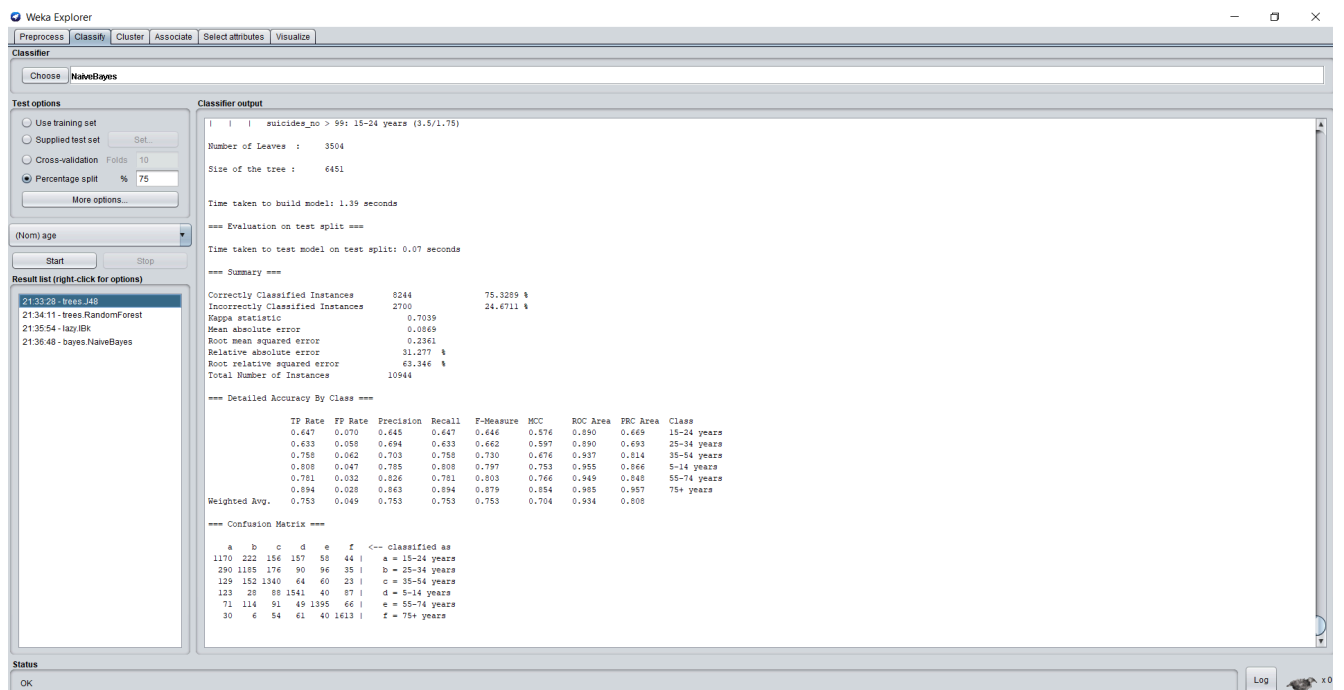
Figure1:- All attributes chart

**Decision Tree:** Correctly classified instances are 8244, about 75.31% and incorrectly classified instances are 2700, about 24.67% among the total number of 10966 instances. In Decision Tree algorithm the accuracy for that dataset is about 75%

**Naïve Bayes:** Correctly classified instances are 2615, 23.89% and incorrectly classified instances are 8329, 76.1% among the total number of 10944 instances. In Naïve Bayes algorithm the accuracy for that dataset is about 23.89%
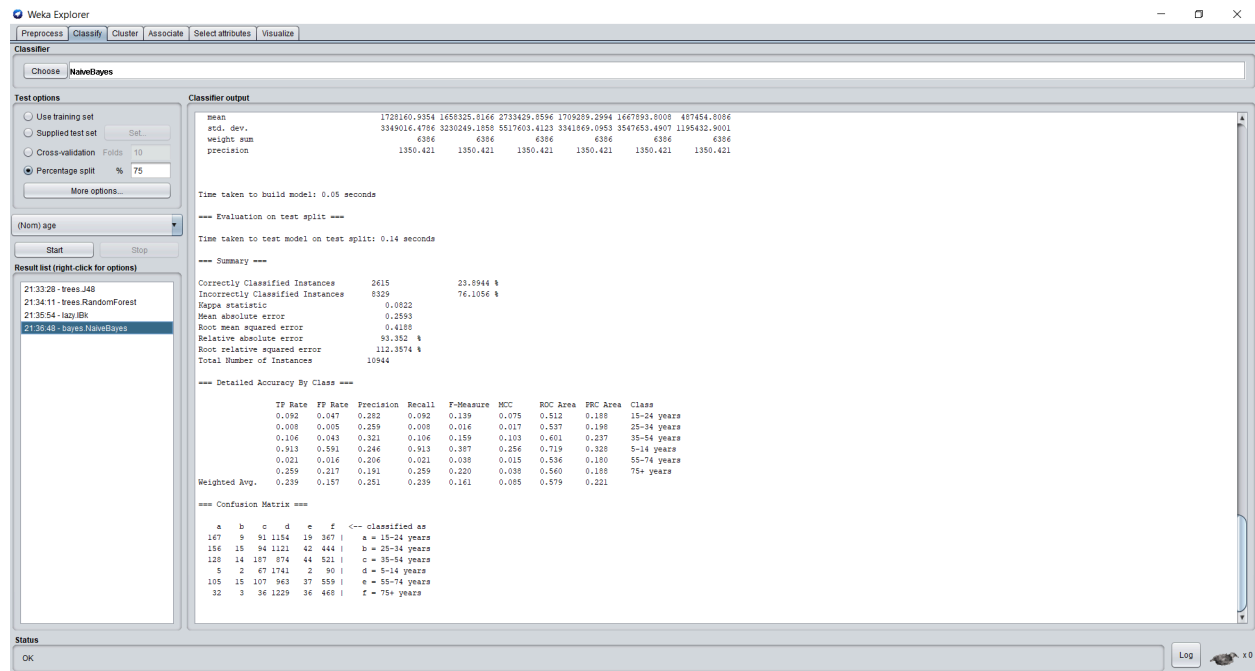


Figure3:- Naïve Bayes observation

**Random forest:** Correctly classified instances are 8442, 77.1% and incorrectly classified instances are 2502, 22.86% among the total number of 10944 instances. In Random forest algorithm
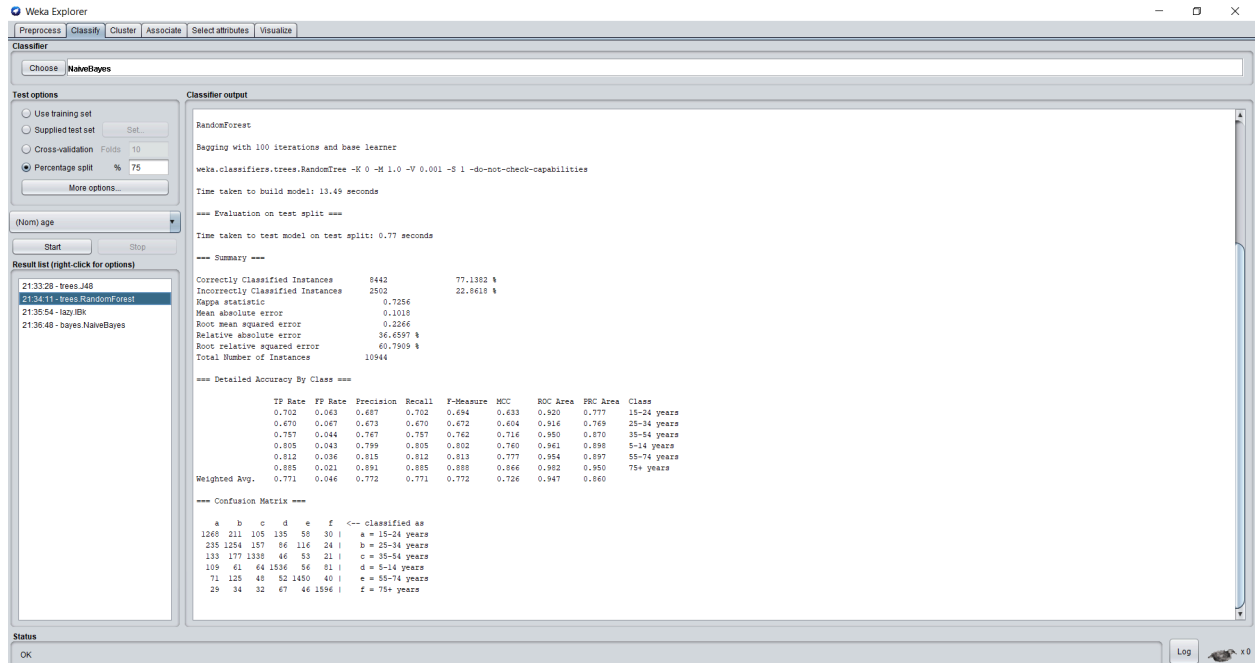the accuracy for that dataset is about 77%

Figure4:- Random forest observation

## KNN:

Correctly classified instances are 1164, 10.67% and incorrectly classified instances are 9776, 89.33% among the total number of 10944 instances. In KNN algorithm the accuracy for that dataset is about 10.97%
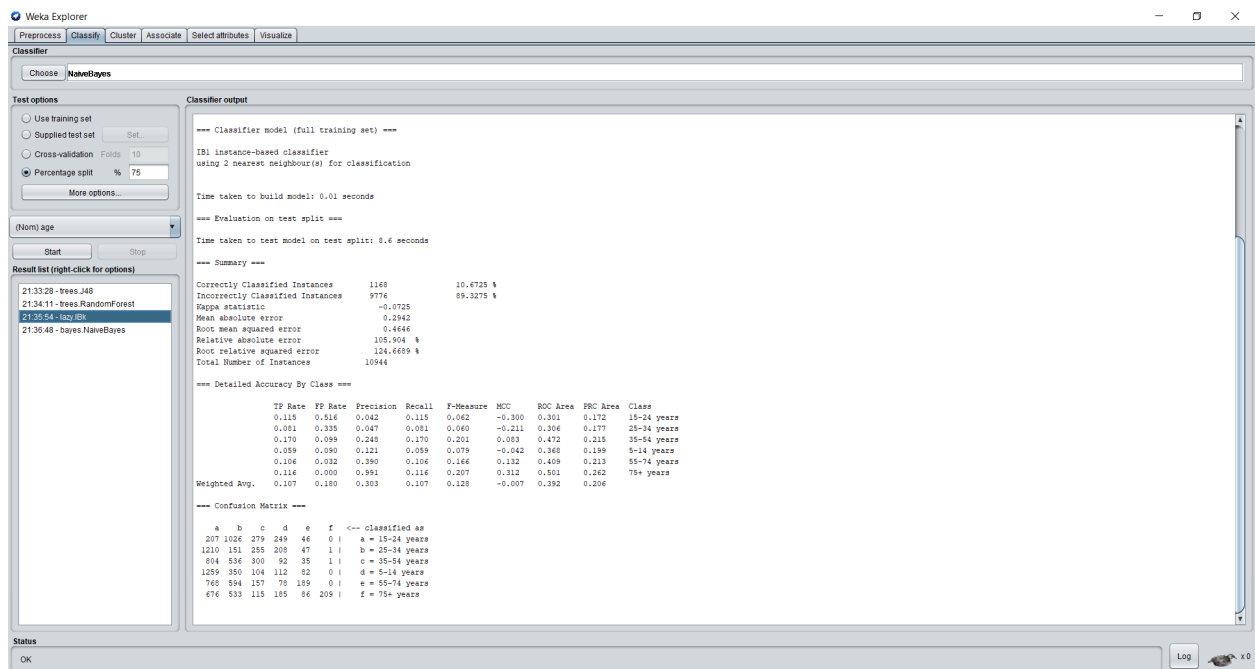
Figure4:- KNN observation

## Performance Evolution

In **table-1** in this section we show about weka tool implemented accuracy.

| Algorithm | Precision | Recall | f-score | Support | Accuracy |
|-----------|-----------|--------|---------|---------|----------|
| Decision Tree | .645 | .647 | .646 | 10944 | 75% |
| Naïve Bayes | .282 | .092 | .139 | 10944 | 23.89% |
| KNN | .042 | .115 | .062 | 10944 | 10.97% |
| Random Forest | .687 | .702 | .694 | 10944 | 77% |

**Table1**: Weka tool Implemented Confusion Matrix Report

In **table-2** we have shown python implemented accuracy

| Algorithm | Precision | Recall | f-score | Support | Accuracy |
|-----------|-----------|--------|---------|---------|----------|
| Decision Tree | 1 | 1 | 1 | 1513 | 100% |
| Naïve Bayes | .2 | .3 | .2 | 1513 | 20.31% |
| KNN | .36 | .20 | .26 | 2709 | 22.85% |
| Random Forest | 1 | 1 | 1 | 1513 | 100% |

## Discussion

In table-1 we have seen Precision, Recall, f-score, Support, Accuracy. In table-1 shows Random forest gives the best result. For the increase of data, Random forest is a good classifier. So in this project random forest gives the best data for less data.
Decision tree and Random forest both give the best result for table-2 .In table-2, Decision tree and Random Forest give the same accuracy.

## Conclusion

We learned some essentials about decisions trees, Naïve Bayes, random forests, KNN. We applied those algorithms to a dataset, and we compared the accuracy of those models. So all of the algorithms help to make decisions or predict. Using those algorithms, businessmen can predict their business and the mining might bring out the maximum profit.