Nishi Gupta
Gu1978

# NAME: NISHI GUPTA
# CSC 5800

# PROJECT TOPIC:

## Predict Student Performance in Secondary Education

## (high school)

# <u>INDEX</u>

Nishi Gupta
Gu1978

# OVERVIEW

Educational system is one of the important parts for development of any country. Most of the developed countries have their own educational system and evaluation criteria. In this work, we will analyze real world data from two Portuguese secondary schools. The problem statement is based on paper published at 5[th] future public technology conference held at Portugal during the year 2008.[1]

Modeling student performance is an important tool for both educators and students, since it can help a better understanding of this phenomenon and ultimately improve it. Although the educational level of the Portuguese population has improved in the last decades. In particular, lack of success in the core classes of Mathematics and the Portuguese language is extremely serious. As mentioned in the paper[1], they have modeled under three DM goals:

i)       Binary classification (pass/fail)

ii)      Classification with five levels (from I very good or excellent to V - insufficient)

iii)     Regression, with a numeric output that ranges between zero (0%) and twenty (100%)

Here, in this project we have classified students into three categories, "Good", "Fair", and "Poor", on the basis of their final exam performance and created various DM models to predict students' final performance and compared it.

# DATASET OVERVIEW -

## 1.  Dataset Information :-

We have used the Student Performance Dataset published by UCI Machine Learning Repository. The dataset has 395 rows and 33 features. We have used only Mathematics dataset for analysis in this project. Each sample consists of student's background information and grades. Background information consists of some personal information of student (i.e. smoking, study time, etc.), family information of student (i.e. father-mother education, job, etc.). There are three columns of grades given in dataset; G1, G2 and G3. G1 and G2 are periodic grades throughout the year and G3 is the final grades for Mathematics.
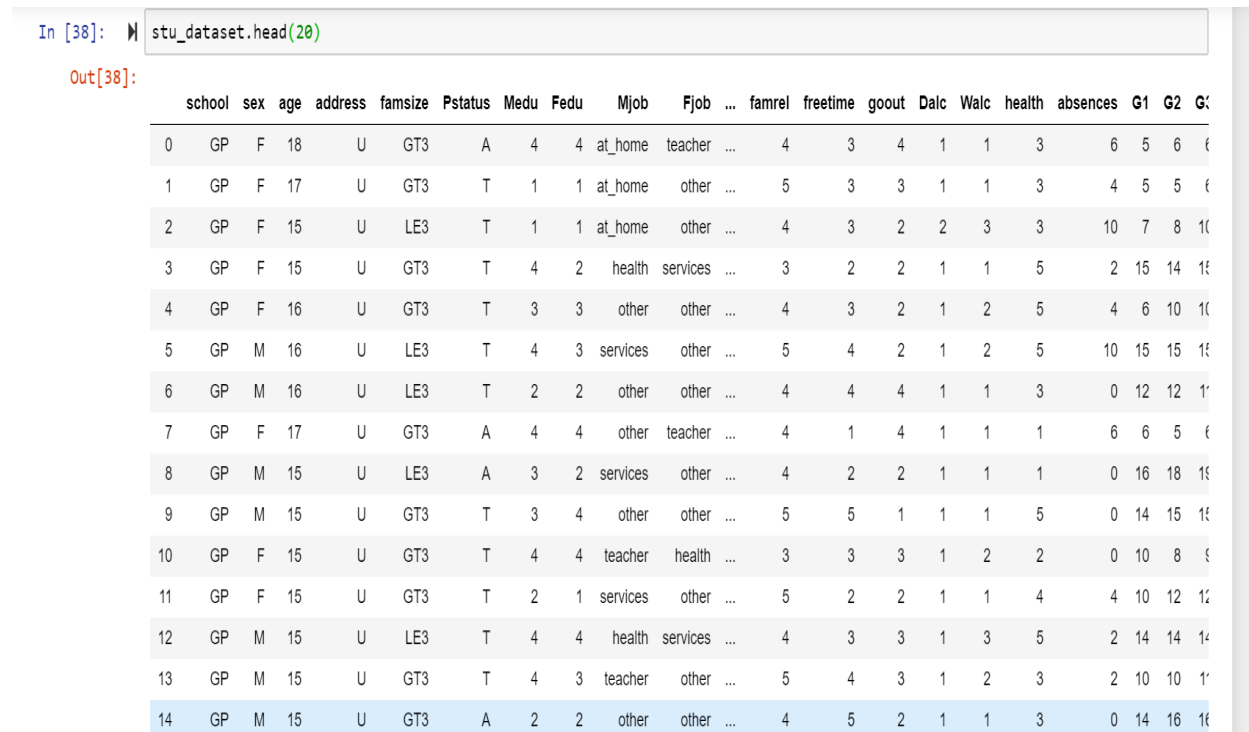
File Name: student-mat.csv

File Size: 57KB

Nishi Gupta
Gu1978

Table 1 : Description of the student attribute taken into build the dataset[1]

| Attribute | Description (Domain) |
|---|---|
| sex | student's sex (binary: female or male) |
| age | student's age (numeric: from 15 to 22) |
| school | student's school (binary: Gabriel Pereira or Mousinho da Silveira) |
| address | student's home address type (binary: urban or rural) |
| Pstatus | parent's cohabitation status (binary: living together or apart) |
| Medu | mother's education (numeric: from 0 to 4a) |
| Mjob | mother's job (nominalb) |
| Fedu | father's education (numeric: from 0 to 4a) |
| Fjob | father's job (nominalb) guardian student's guardian (nominal: mother, father or other) |
| famsize | family size (binary: $\leq 3$ or $> 3$) |
| famrel | quality of family relationships (numeric: from 1 – very bad to 5 – excellent) |
| reason | reason to choose this school (nominal: close to home, school reputation, course preference or other) |
| traveltime | home to school travel time (numeric: $1 - < 15$ min., $2 - 15$ to 30 min., $3 - 30$ min. to 1 hour or $4 - > 1$ hour). |
| studytime | weekly study time (numeric: $1 - < 2$ hours, $2 - 2$ to 5 hours, $3 - 5$ to 10 hours or $4 > 10$ hours) |
| failures | number of past class failures (numeric: n if $1 \leq n < 3$, else 4) |
| schoolsup | extra educational school support (binary: yes or no) |
| famsup | family educational support (binary: yes or no) |
| activities | extra-curricular activities (binary: yes or no) |
| paidclass | extra paid classes (binary: yes or no) |
| internet | Internet access at home (binary: yes or no) |
| nursery | attended nursery school (binary: yes or no) |
| higher | wants to take higher education (binary: yes or no) |
| romantic | with a romantic relationship (binary: yes or no) |
| freetime | free time after school (numeric: from 1 – very low to 5 – very high) |
| goout | going out with friends (numeric: from 1 – very low to 5 – very high) |
| Walc | weekend alcohol consumption (numeric: from 1 – very low to 5 – very high) |
| Dalc | workday alcohol consumption (numeric: from 1 – very low to 5 – very high) |
| health | current health status (numeric: from 1 – very bad to 5 – very good) |
| absences | number of school absences (numeric: from 0 to 93) |
| G1 | first period grade (numeric: from 0 to 20) |
| G2 | second period grade (numeric: from 0 to 20) |
| G3 | final grade (numeric: from 0 to 20) |

Nishi Gupta
Gu1978

## 2. Dataset Exploration :-

We have done dataset exploration in IPython Notebook as well as the attribute information given above provided valuable information regarding the dataset. Some of the important discoveries are found as follows:



Figure 2.1: Overview of Dataset

- There are 395 instances and 33 attributes in this dataset.
- All the 32 attributes except G3 are independent and support in the prediction of G3 which is output label.
- Grades G1, G2 & G3 are given on scale of 1-20.
- The dataset is mixture of numeric and nominal attribute.
- There are **no null values or missing values** for all the attributes.
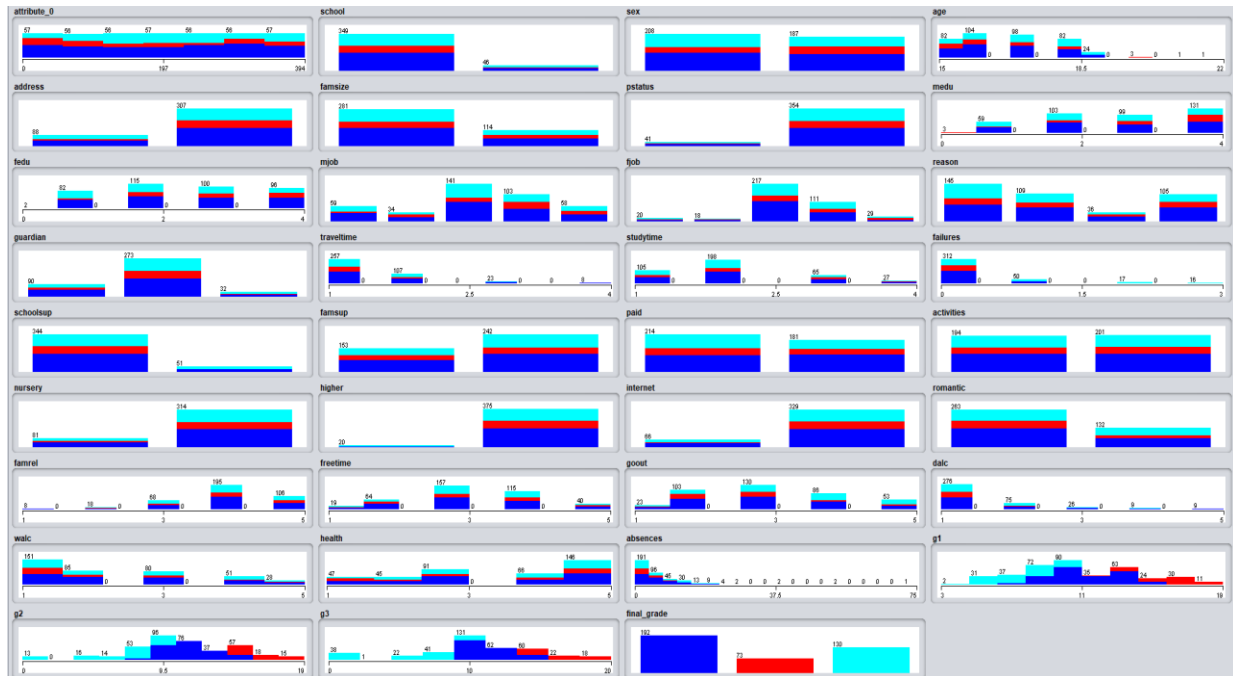
Nishi Gupta
Gu1978

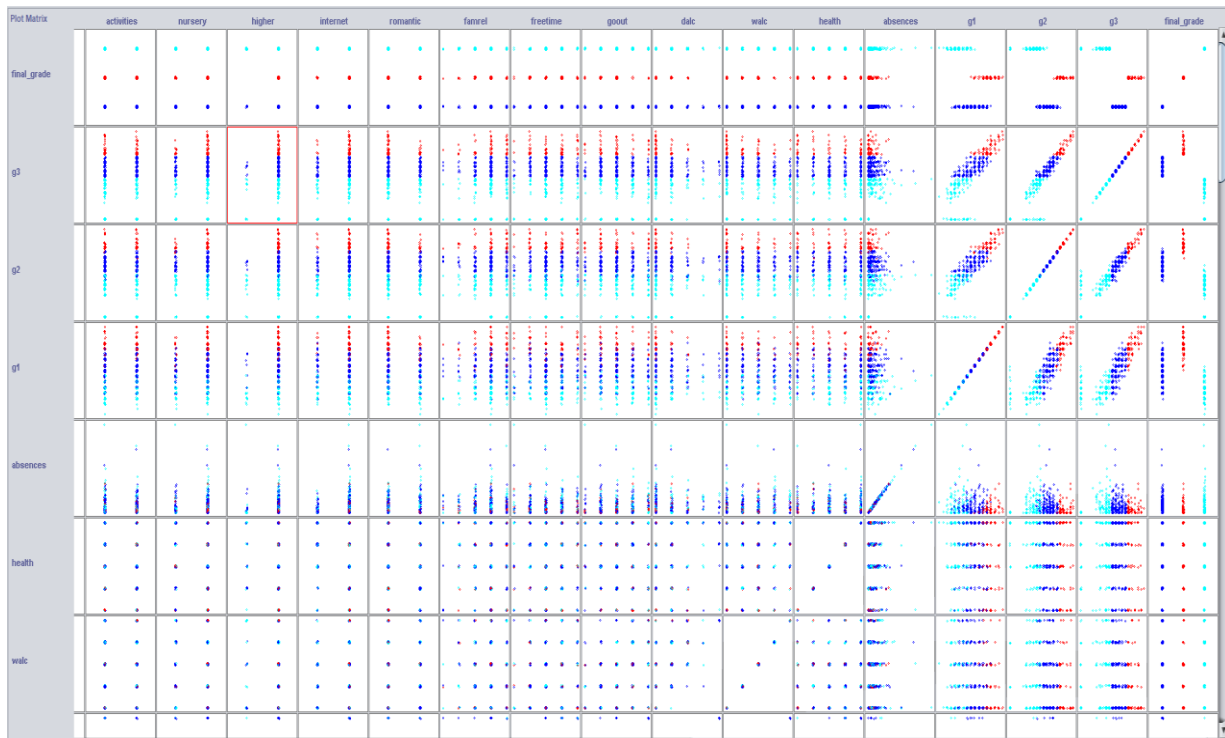Figure 2.2: Histogram Plot for some of the features



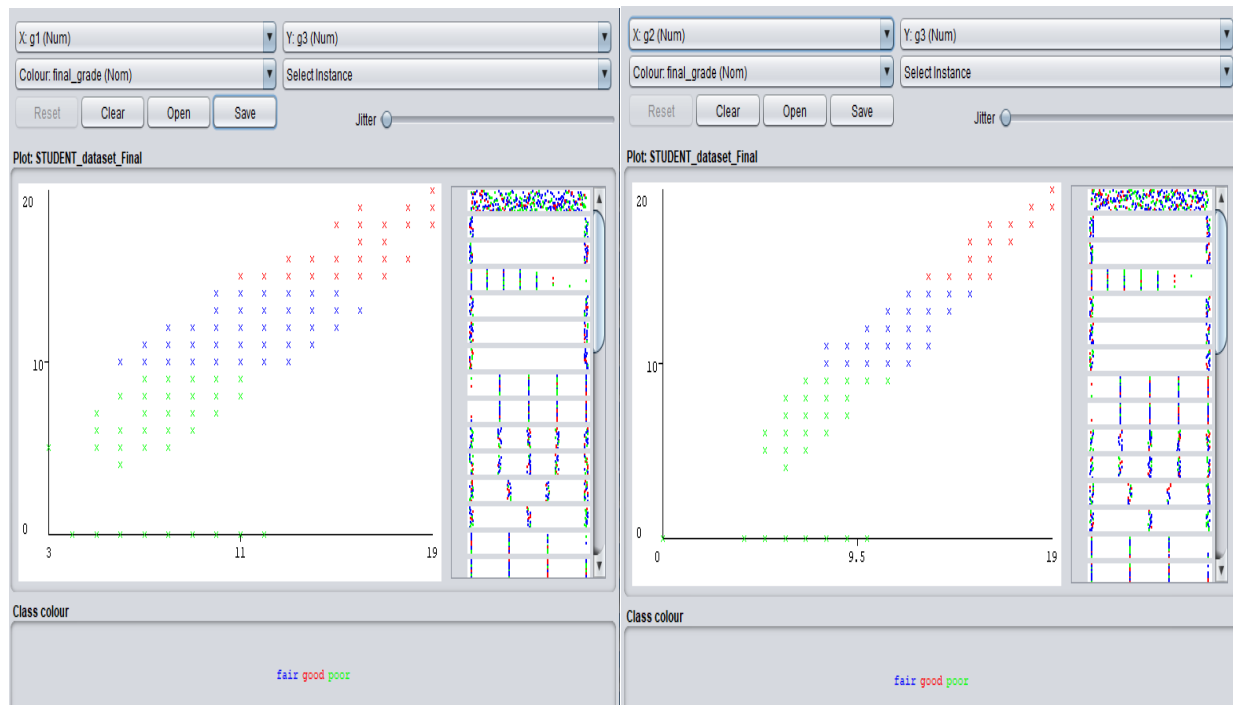Figure 2.3: Plotting of some features

Figure 2.4: Snapshot of Correlation between G1 & G3, G2 & G3 respectively

## 3. Data Mining Software Tools:

- **IPython Notebook** – Python in IPython Notebook for data exploration and preprocessing.
- **WEKA** 3.8.2 - Waikato Environment for Knowledge Analysis (WEKA) tool for classification models.

# DATA PRE-PROCESSING –

Data preprocessing is a data mining technique that involves transforming raw data that is often incomplete, inconsistent, lacking in certain trends, into an understandable format, so that we can further process the data to apply Data mining algorithms.

Steps Performed:-

- Performed pre-processing with *PYTHON in IPython notebook* by using numpy and pandas library. we have checked in each column for missing values. But in given dataset, there are no missing values available. So, there is no need of data cleaning in given dataset.
- Since, most of the columns are categorical attributes, if we apply variance filter, then it might remove some important features. So, we are not applying variance filter.

- Initially, the target output class (G3) ranges from 0 to 20 then mapped into Final_grades attribute class by applying categorical value-

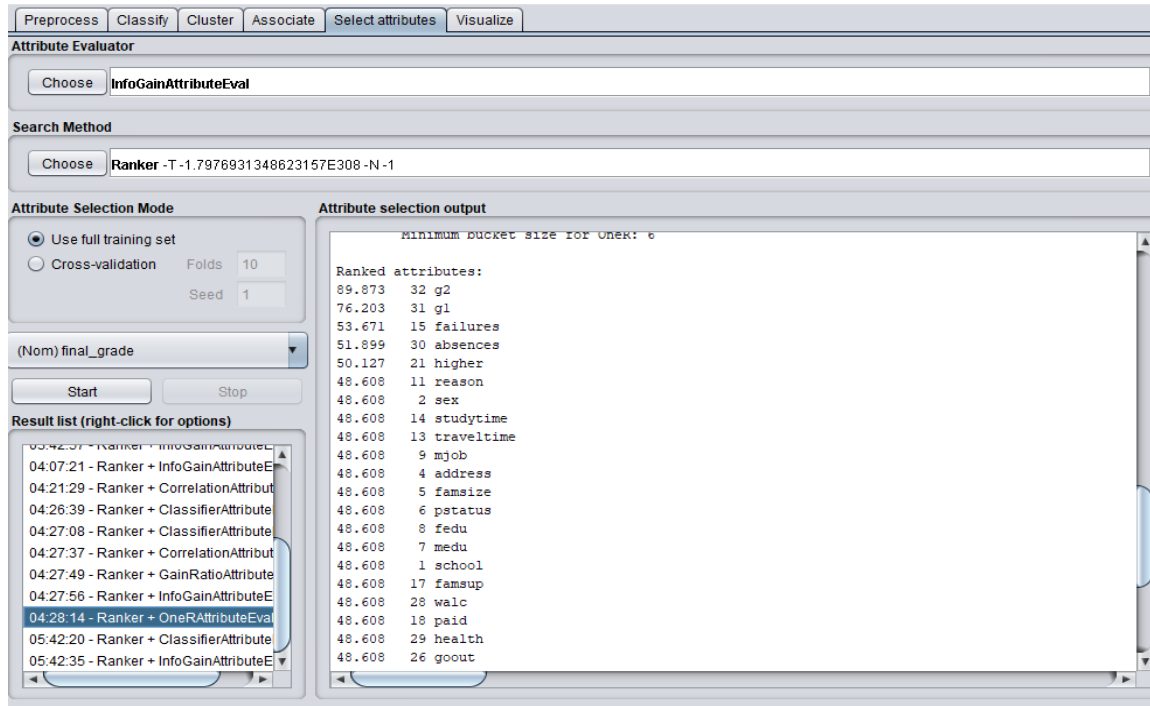| | |
|---|---|
| Good | 15-20 |
| Fair | 10-14 |
| Poor | 0-9 |

- Label Encoder method is applied on Final_grades attribute class which encode target labels; 0,1 & 2 are encoded label.
- Created data frame named dfd for classification; where G3 column is dropped for future prediction.
- As the downloaded data came in csv format. Thus, in order to use the dataset in *WEKA*, we converted the CSV file into ARFF.

# **ALGORITHM SELECTION-**

As mentioned before, our goal is mainly focused on classification. we have created various data mining models to predict student final performance and compared model performance based on sample accuracy score.

After going through the paper[4] as mentioned that they have performed classification with five class and eight random attributes and their maximum accuracy was 76.73%. So, my goal was to increase the accuracy and thought to remove redundant variables.

We are selecting the best attribute on the basis of attribute's evaluator method namely CorrelationAttributeEval, GainRatioAttributeEval and InfoGainAttributeEval with Ranker Search Method in WEKA. After finding the result of the implementing attribute evaluator algorithm, we are selected 5 different attribute which mostly affect our prediction result.

We have implemented total eight models out of which we have selected only top five model on the basis of accuracy. The lowest accuracy models are lazy IBK, Kstar, ZeroR having accuracy in between 45-70% accuracy. Here are the selected model for our dataset-

    I.    Decision Tree (J48) classifier

    II.    Naïve Bayes Classifier

    III.    Random Forest Classifier

    IV.    JRip Classifier

    V.    Multilayer Perceptron

NOTE :- For implementation of all the classification models, we have tested on 10-fold cross validation.

## 1. DECISION TREE (J48):

➢ Decision Tree is one of the most popular Data Mining algorithm. Decision Tree is used for both classification and regression problems. Here, we used it as a classification problem

➢ J48 showed accuracy 87.3418 % when experimented with all 32 attributes.

➤ When we experimented with randomly 5 selected attributes and set minNumObj to six and numFolds to five, then the accuracy is increased to 89.8734 %, which is in increment of 2.5316%.
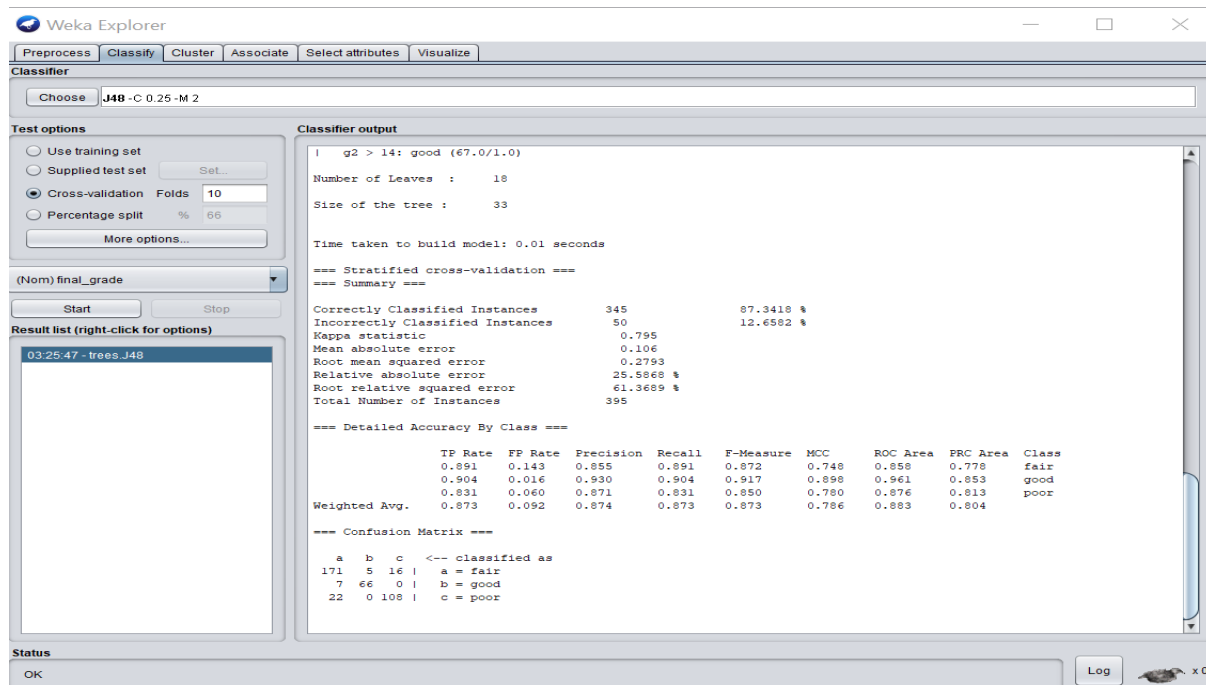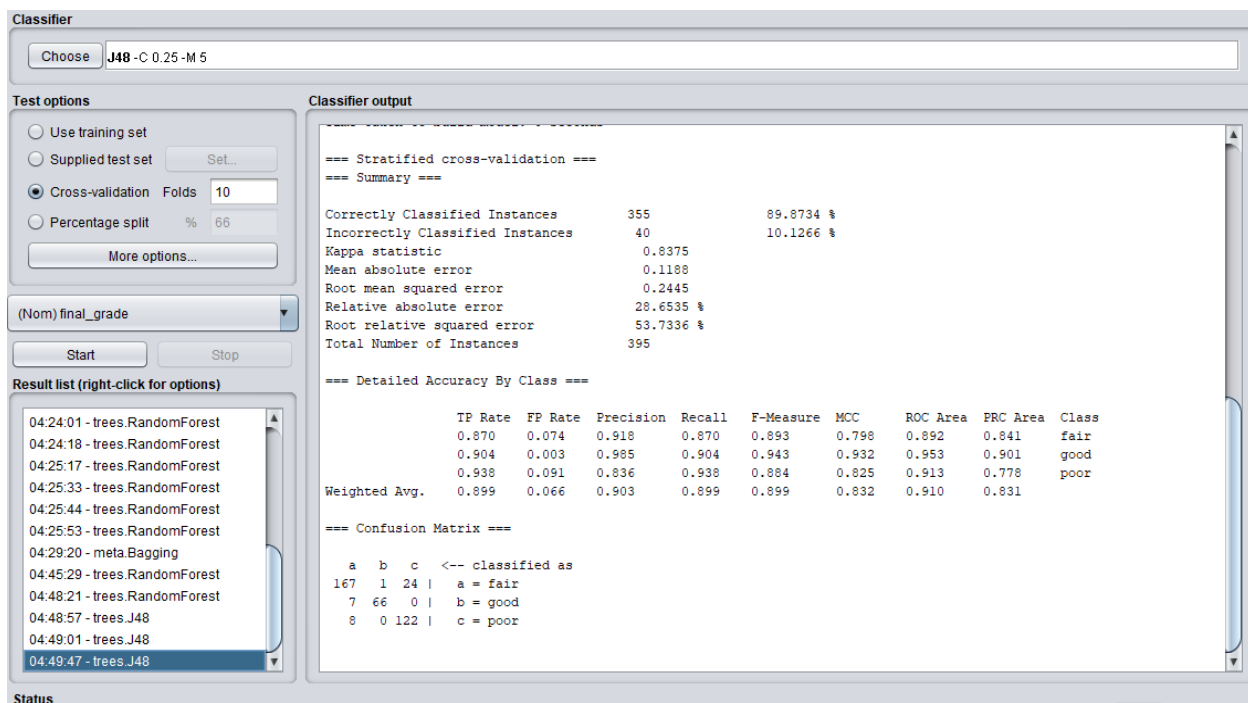


Figure1.1: J48 with 32 attributes



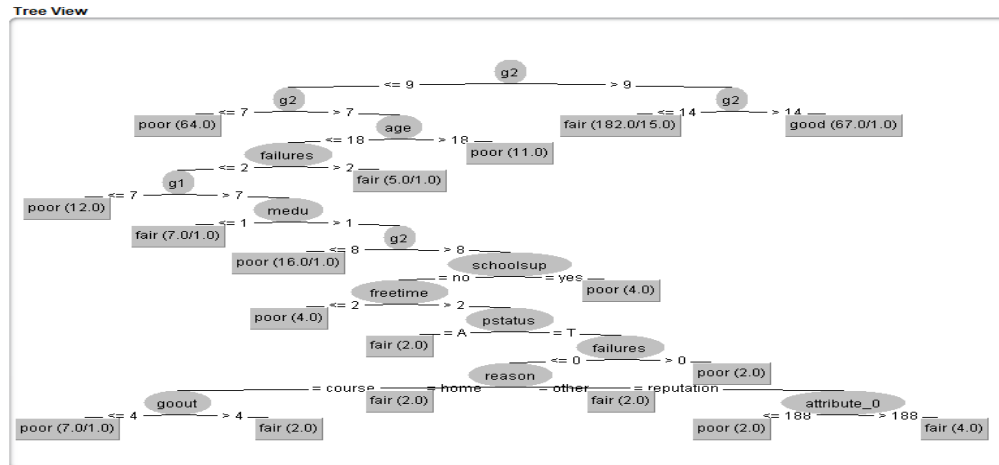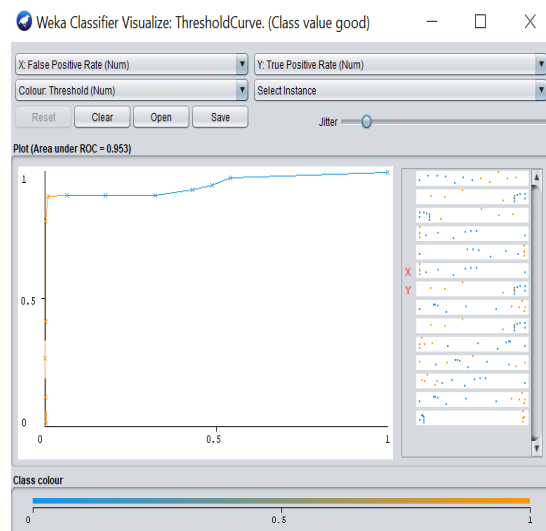Figure 1.2: J48 with random 5 attributes

Figure 1.3: Tree



Figure 1.4: ROC curve for J48(class value good)

## 2. NAIVEBAYES CLASSIFICATION :

➢ Naïve Bayes classifier is based on Bayes theorem to classify objects.

➢ When we experimented with all 32 attributes, Naïve Bayes showed the accuracy of 81.77%.

➢ With the randomly selected five attributes, accuracy increased to 83.038%, this shows accuracy of Naive Bayes is increased to 1.268%.

➢ For better performance of this model, we set kernelEstimator to True and then model performance was quite good which is 85.3165%.
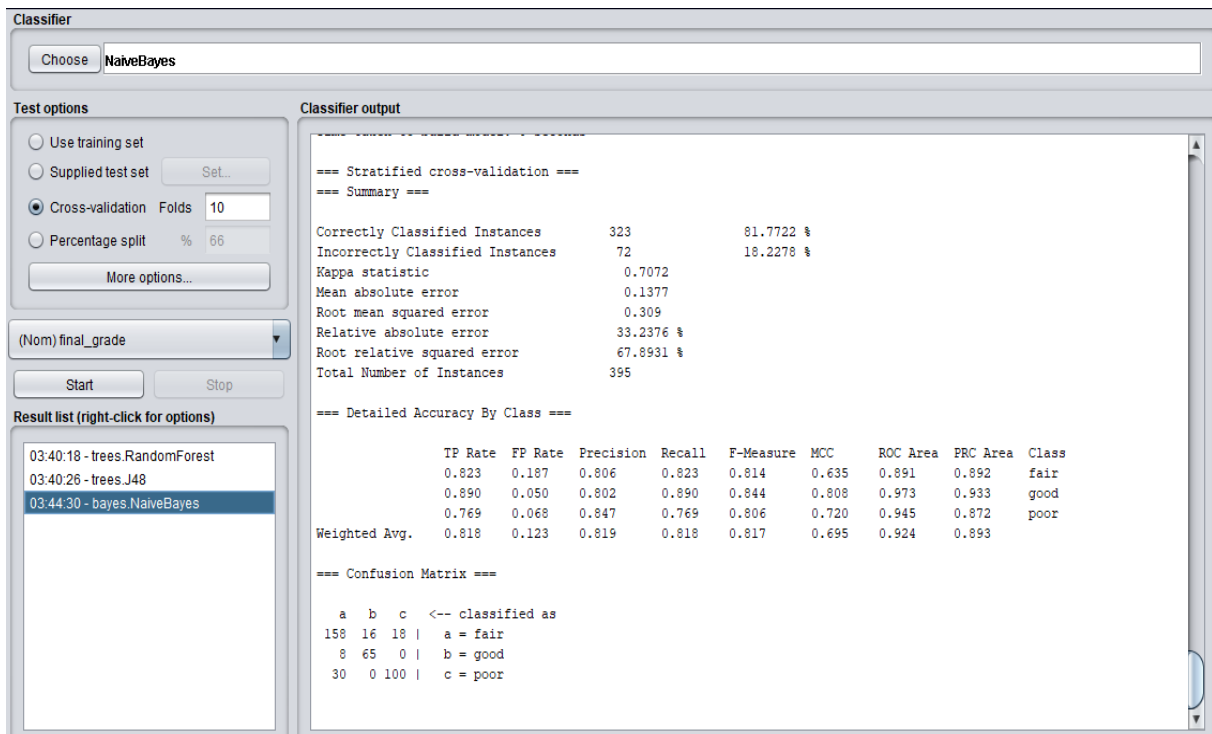
Nishi Gupta
Gu1978


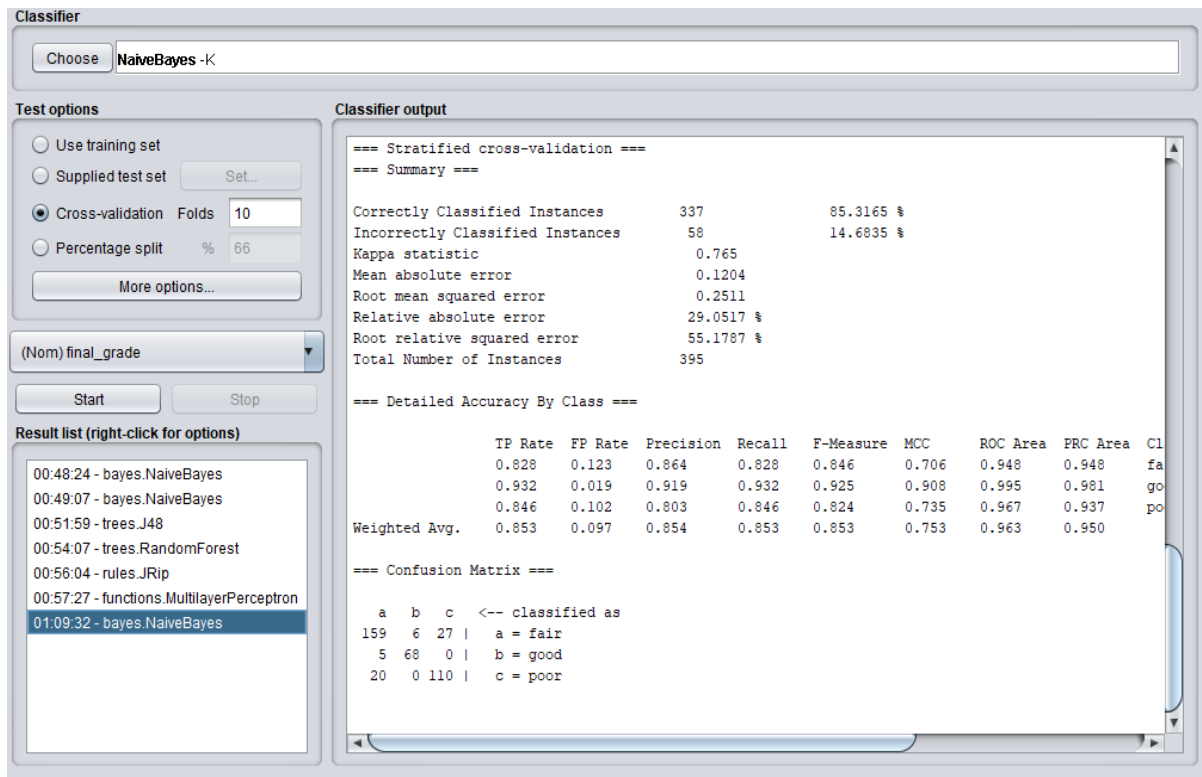
Figure 2.1: NaiveBayes with 32 attributes



Figure 2.2: NaiveBayes with 5 random attributes and KernelEstimator set TRUE
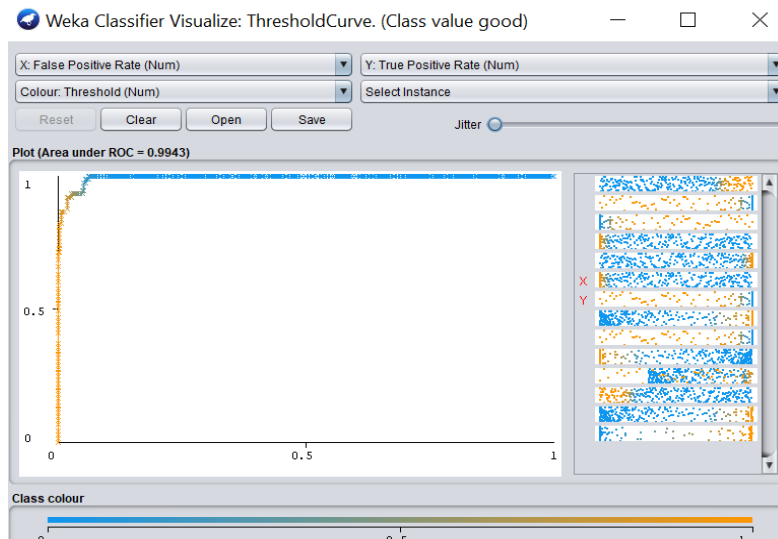
Figure 2.3: ROC curve for Naïve Bayes(class value good)

## 3. RANDOM FOREST :

➢ Random Forest algorithm is used as a supervised learning and known for its powerfulness. In addition to powerful decision tree representation, it is capable of generalizing well.

➢ It showed accuracy with all 32 attributes is 88.6076 %.

➢ When we randomly select 5 attributes and change the maxDepth to 3, accuracy become 89.8734 % which is slight increase in the accuracy i.e 1.2658%.
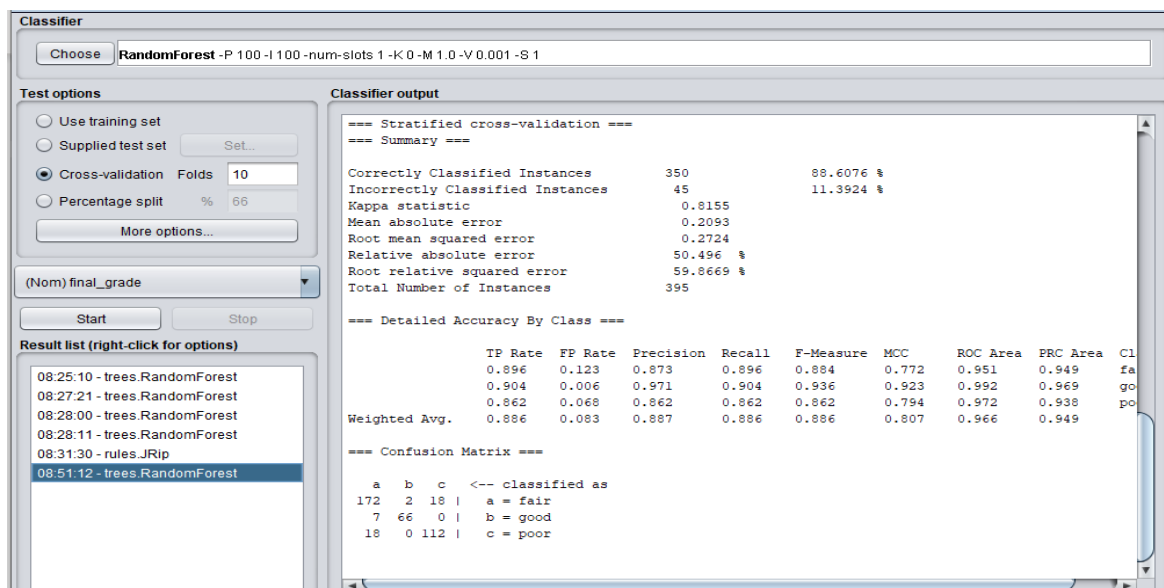


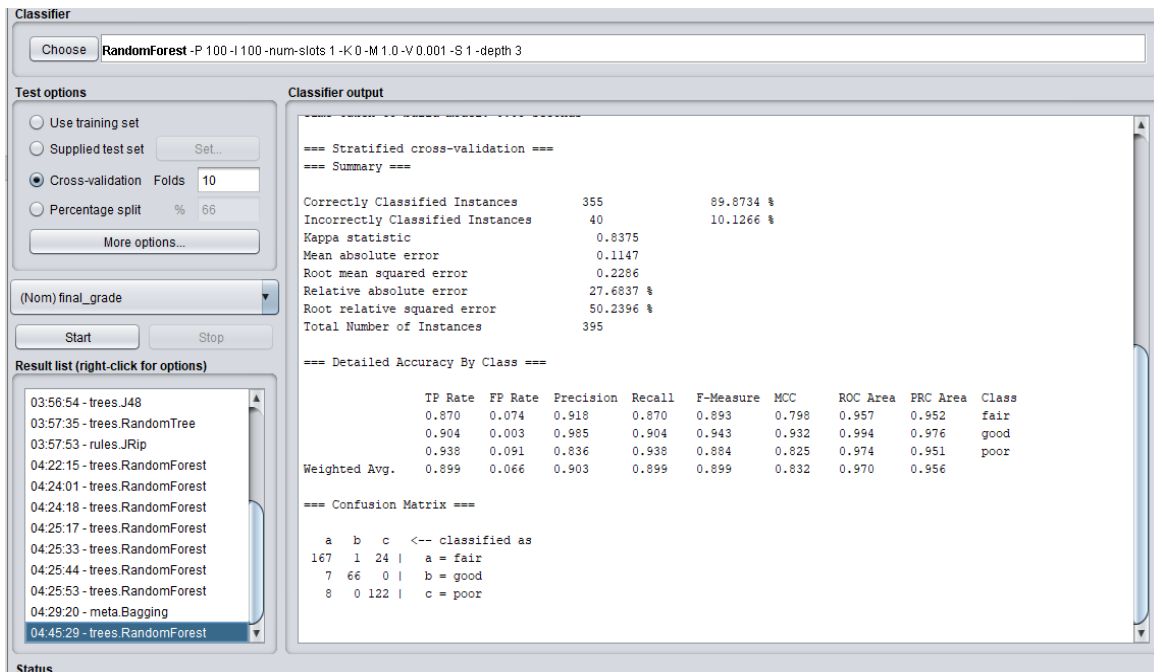Figure 3.1: Random Forest with 32 attributes

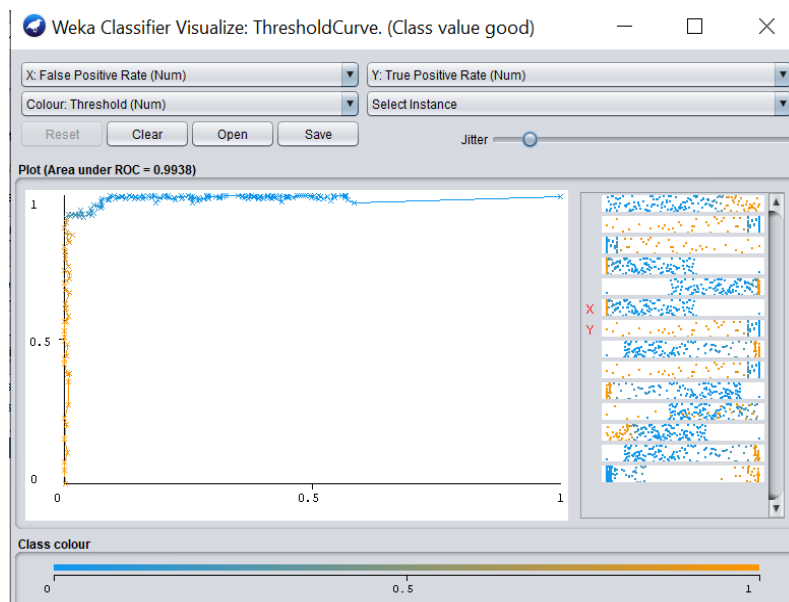Figure 3.2: Random Forest with 5 random attributes



Figure 3.3: ROC curve for RandomForest(class value good)

## 4. JRIP CLASSIFICATION ALGORITHM:

➢ JRip is based on association rules with reduced error pruning (REP), a very common and effective technique found in decision tree algorithms.

➢ When we experimented with all 32 attributes, JRip Classification Algorithm showed 89.1139% accuracy.

➢ When randomly five attributes selected and set folds to 5, seed value is set to 5, numDecimalPlaces value is set to 4, its accuracy increased to 89.8734% which is slight increase in the accuracy i.e. 0.7595% improve the prediction accuracy.
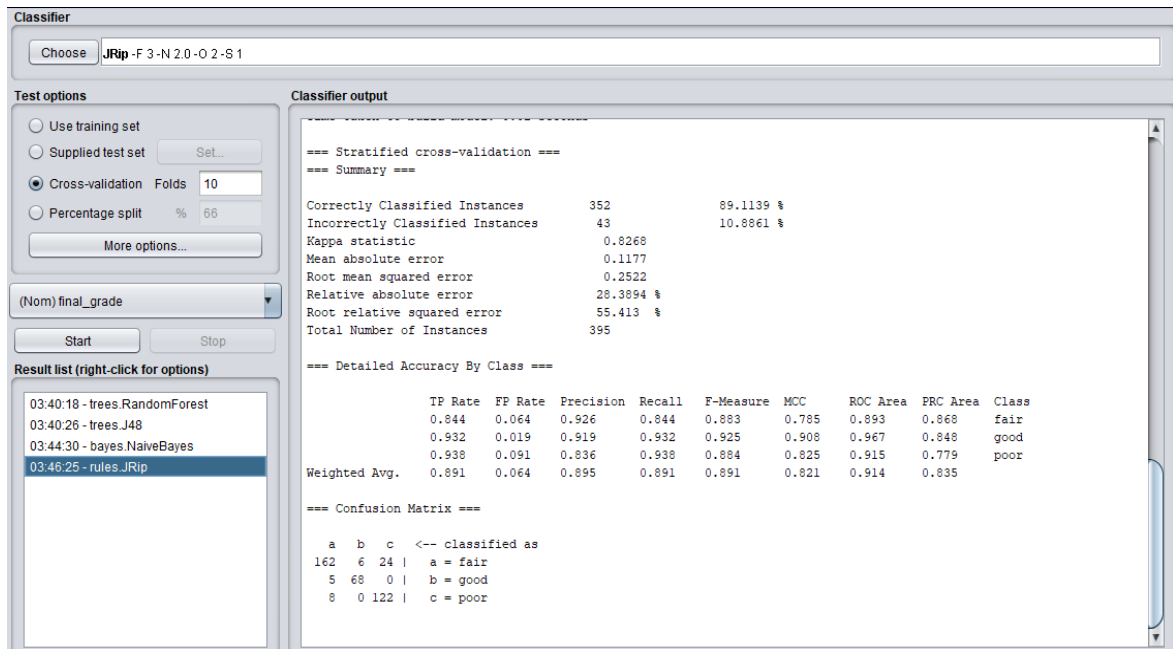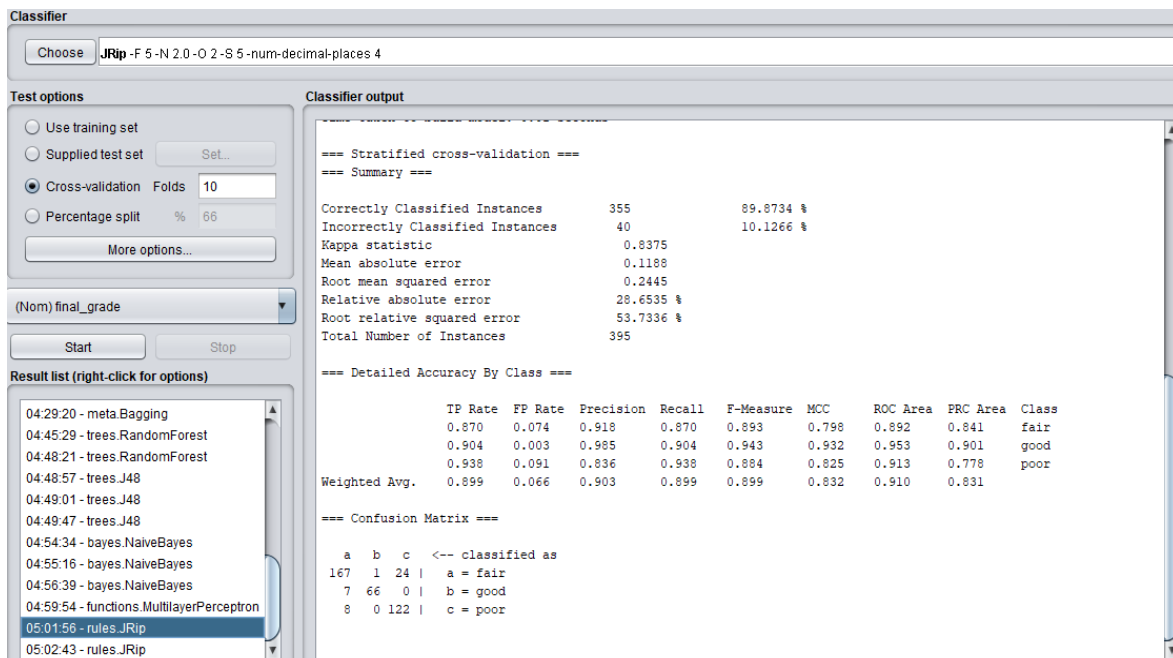


Figure4.1: JRip with 32 attributes



Figure 4.2: JRip with random 5 attributes

Figure 4.3: ROC curve for JRip(class value good)

## 5. MULTILAYER PERCEPTRON :

➢ Multilayer Perceptron is a feed forward artificial neural network model trained with the standard back propagation algorithm that maps sets of input data onto a collection of acceptable output.[3]

➢ It showed accuracy with all 32 attributes is 86.5823%.

➢ When randomly five attributes were selected it showed the increased accuracy of 88.1013 %.



Figure 5.1: Multilayer Perceptron with 32 attributes

Figure 5.2: Multilayer Perceptron with 5 random attributes



Figure 5.3: ROC curve for Multilayer Perceptron(class value good)
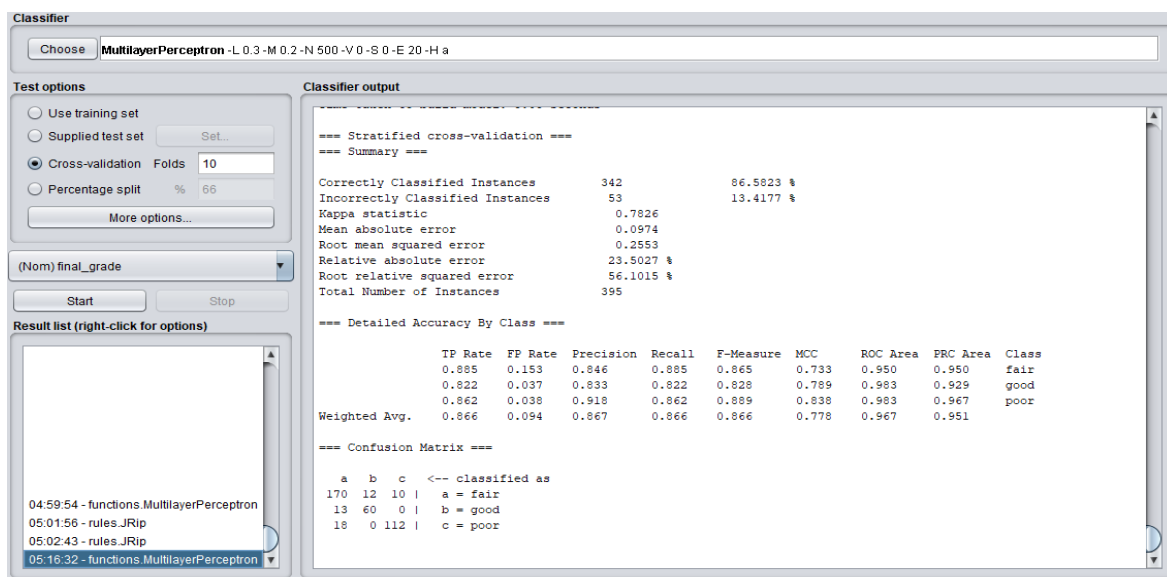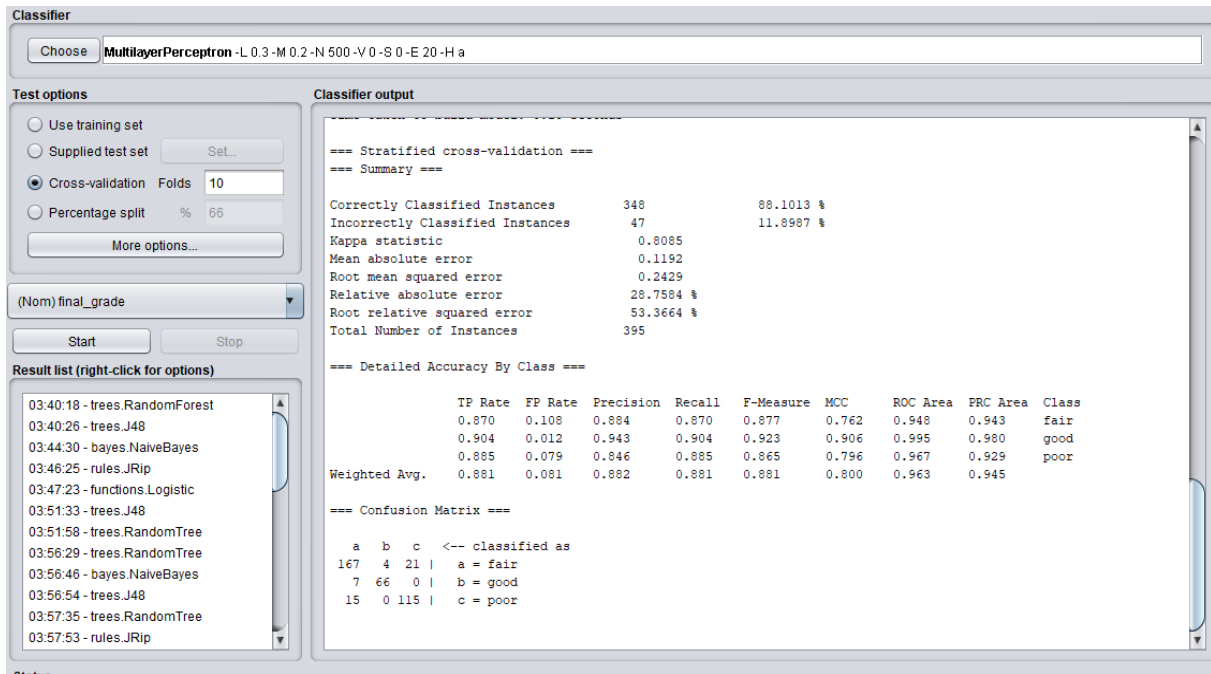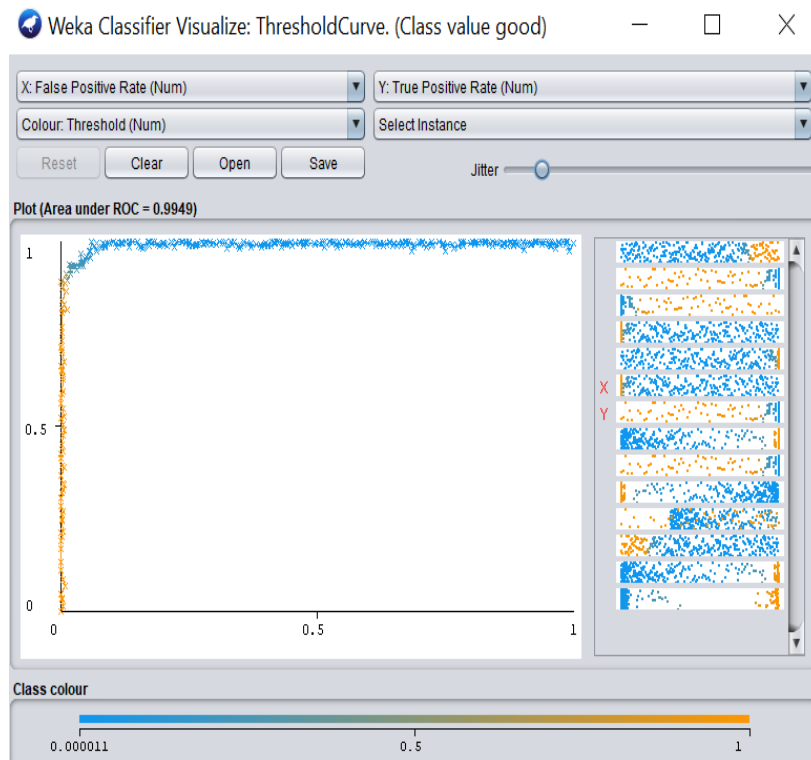
Nishi Gupta
Gu1978

# ANALYSIS RESULTS AND COMPARISON:

After implementing all five models, I get different measures like accuracy, precision, recall, F1, TP rate, FP rate, Error rate, etc. and compared all model on the basis of these measures.

| Measures | Decision Tree (J48) | Naïve Bayes | Random Forest | JRip Classifier | Multilayer Perceptron |
|---|---|---|---|---|---|
| Accuracy | 89.8734 %, | 85.3165%. | 89.8734% | 89.8734% | 88.101% |
| TP Rate | 0.899 | 0.830 | 0.899 | 0.899 | 0.881 |
| FP Rate | 0.066 | 0.112 | 0.066 | 0.066 | 0.081 |
| Precision | 0.903 | 0.830 | 0.903 | 0.903 | 0.882 |
| Recall | 0.899 | 0.830 | 0.899 | 0.899 | 0.881 |
| F-measure | 0.899 | 0.830 | 0.899 | 0.899 | 0.881 |
| ROC Curve | 0.910 | 0.954 | 0.910 | 0.910 | 0.963 |
| Error Rate | 10.126% | 16.962% | 10.126% | 10.126% | 11.898% |

By comparing, we got the three models having same highest accuracy with 89.87%, which is J48, Random Forest and JRip for predicting student result and they perform equally well among all five classification models.

All models in the descending order of their accuracy :-

➢ J48 : 89.87%
➢ Random Forest : 89.87%
➢ JRip : 89.87%
➢ Multilayer Perceptron : 88.10%
➢ Naïve Bayes : 85.31%

Nishi Gupta
Gu1978

# CONCLUSION:

After implementing 5 different classification algorithms on our dataset. The obtained result shows that we can achieve high accuracy, provided that the first and second period grades are known. From the decision tree classification figure1.3, it is clear that G2 is the most significant factor as it appear on the top of decision tree for the prediction of G3 (final grades).

Overall, from our implementation of the models we can conclude that J48, Random Forest and JRip are performed equally well with more than 89% accuracy and surpass other classification algorithm.

# CHALLENGES

1. After dataset preprocessing, the difficulty which faced is in converting CSV file to ARFF in WEKA.
2. Choosing perfect classification for a data mining problem like this is always a challenge as we have to perform lot of trials in selection of correct attribute which would lead to highest accuracy.
3. While using WEKA, we have chosen cross-validation for 10 folds. Having 10 folds means 90% of full data is used for training (10% for testing). However, the current model also shows good accuracy result. But if we can proceed implementation in R or python then we can check overfitting of data.

# FUTURE ANALYSIS

To develop the same model using conventional programing language so that we can understand the data and handle the data manually. However, the current model also shows the good accuracy but if we can check implementation in R or python then we can check data is overfitted.

Nishi Gupta
Gu1978

# **REFERENCE**

I read the following papers and then implement my project. So, here are the reference:

1. P. Cortez and A. Silva. Using Data Mining to Predict Secondary School Student Performance. In A. Brito and J. Teixeira Eds., Proceedings of 5th FUture BUsiness TEChnology Conference (FUBUTEC 2008) pp. 5-12, Porto, Portugal, April, 2008, EUROSIS, ISBN.
2. Comparative Analysis of Models for Student Performance with Data Mining Tools A. K. Shrivas[1], Pragya Tiwari[2]
3. Prediction and Analysis of Student Performance by Data Mining in WEKA. www.rcciit.org/students_projects/projects/it/2018/GR4.pdf
4. Educational Data Mining: Student Performance Prediction in Academic Y. K. Salal, S. M. Abdullaev, Mukesh Kumar