FACULTY OF ENGINEERING AND TECHNOLOGY
DEPARTMENT OF COMPUTER ENGINEERING

Artificial Intelligence

ENCS3340

**Project 2 Report**
**Machine Learning for Classification**

Prepared by:

Tariq Odeh -1190699

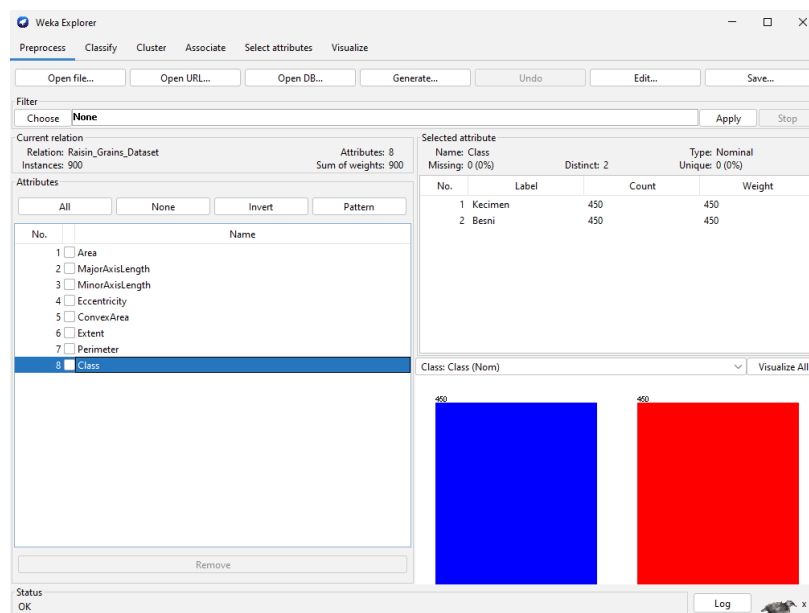Yousef Hammad - 1170625

Sec: 1

Instructor: Dr. Adnan Yahya

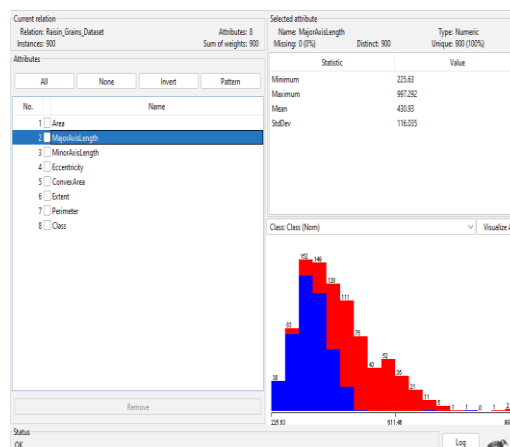Date: 11th June 2022

# Abstract

We'll learn how to assess alternative categorization methods using machine learning tools in this assignment. Additionally, we must compare multiple machine learning algorithms for a classification problem using WEKA, and the dataset we will work on is 2 according to the university id (1170625) with three models Decision.

# 1. Dataset

From the figure we can see the name of relationship Raisin_Grains _Datasetis, the instances (900) , and there are 8 attributes in the dataset, also in panel below current relation shows the name of attributes (Area, Extent, Perimenter, ...). And in the right panel, the selected attribute statistics are displayed. It exhibits the following: Name, Type, Missing, Unique and Distinct. When select class attributevwe can see that the name is Class, type is normal, missing is zero, unique is zero and the distinct are two distinct values (Kecimen and Besni). For the attribute the nominal type means it accepts no numeric values and the numeric values means it accepts numeric values. For each of the 900 instances, the count of each particular class label is provided in the count column. The output class label for the attribute will be displayed in the histogram.The class label for this dataset is either good or terrible. There are 450 Kecimen examples and 450 Besni examples.
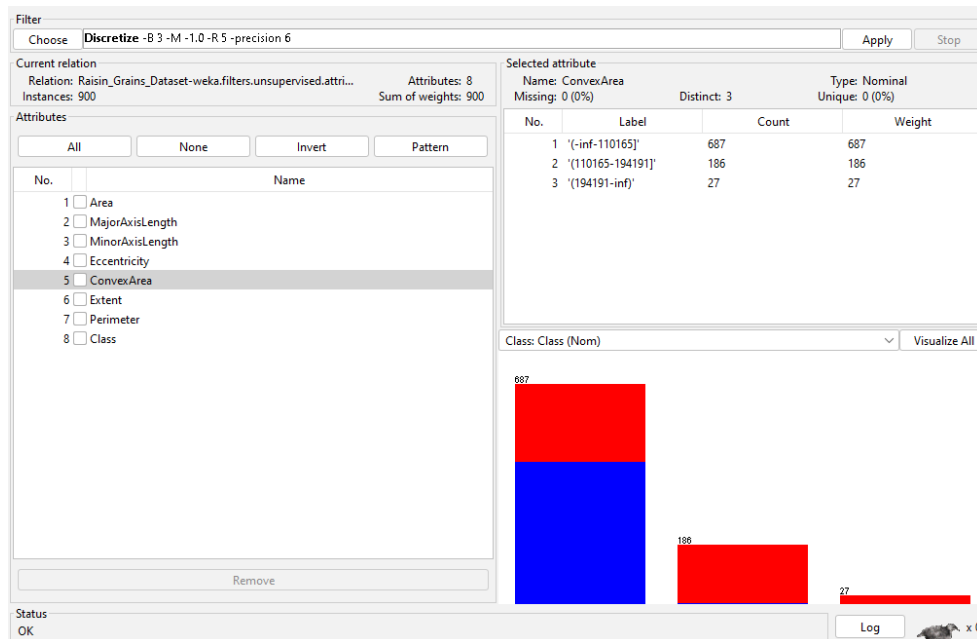


In this figure shown the name of the attribute. The type of the attribute is numeric, thus it's a number. It has 900 different values in 900 different circumstances. It's unique since there are 900 different values that don't match. Minimum monetary value for this parameter is 225. maximum value for this parameter is 997. The mean is calculated by dividing the total number of values by the number of cases. Attribute duration standard deviation and Histogram are used to calculate deviations from the mean.

# 2. Decision Tree

## Filter

Preprocess for one attribute (ConvexArea) discretization of continuous attributes with 3 pins.



## Classifer

in the confusion matrix shown below, it can be seen that the true kecimen have 397 that were classefied correctly and 53 that were classefied as Besni , the true Besni have 360 classefied correctly and 90 that were classefied as keciman.

It can be noticed that the correctly classefied have an accuracy of 84.11% and incorrectly classified have an accuracy of 15.89% .

In the Detailed Accuracy by class it shows the TP rate, FP Rate, Percision, F-Measure,MCC, ROC Area, PRC Area and the class.

As an example Tp rate of the Kecimen class will be calculated by dividing the correctly classefied Kecimans with the total true Kecimens $\frac{397}{450} = 0.882$ .

FP rate will be calculated by dividing incorrectly classefied Keciman (Kecimans classefied as Besni) on the total Kekiman values $\frac{53}{450} = 0.2$ .

Precision ca be calculated by deviding TP over $FP + TP$ $\frac{0.882}{0.2+0.882} = 0.815$ .

Recall should be equal to the TP.



F-Measure can be calculated by $\frac{2* \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$ $= \frac{2* 0.815 * 0.882}{0.815 + 0.882} = 0.847$

# Hyper-parameter

When the hyper parameter of the model is changed to reduced error pruning as shown in figure below it can be noted that the changes in the values displayed on the screen such as correctly classified instances became higher and equal 85.33% and also the incorrectly classified instances changed and it became less and equal 14.66%. In addition, the effects are obvious for the confusion Matrix, TP Rate, FP Rate, Precision, Recall and F- Measure.



# Visualize the trees
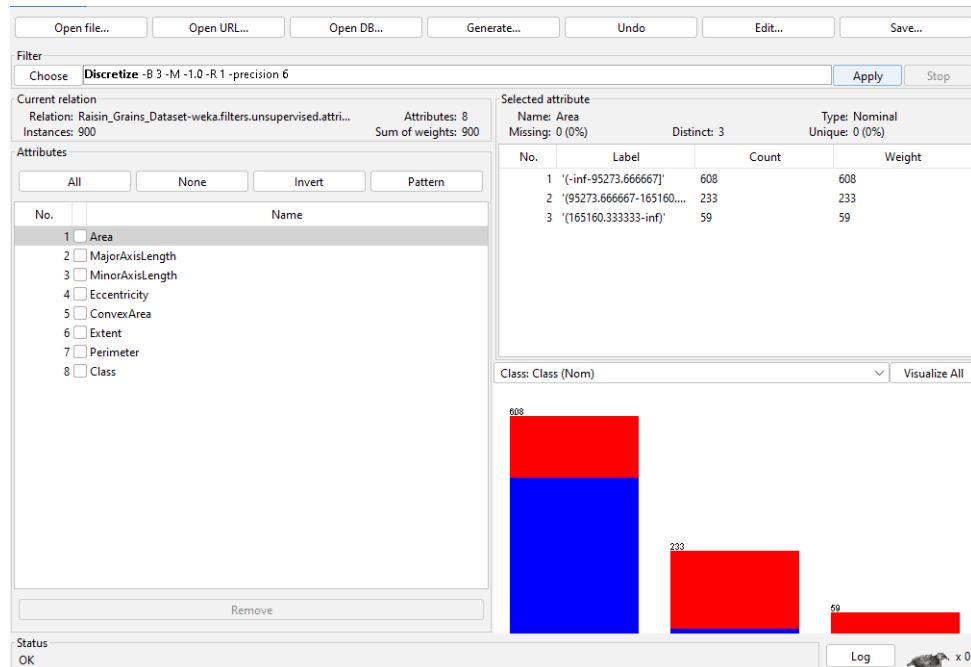
| Default Tree | With reduced error pruning set to true |

# 3. Naïve Bayes

## Filter

Preprocess for one attribute (Area) discretization of continuous attributes with 3 pins.



## Classifer

in the confusion matrix shown below, it can be seen that the true kecimen have 426 that were classefied correctly and 24 that were classefied as Besni , the true Besni have 323 classefied correctly and 127 that were classefied as keciman.

It can be noticed that the correctly classefied have an accuracy of 83.22% and incorrectly classified have an accuracy of 16.78% .

In the Detailed Accuracy by class it shows the TP rate, FP Rate, Percision, F-Measure,MCC, ROC Area, PRC Area and the class.
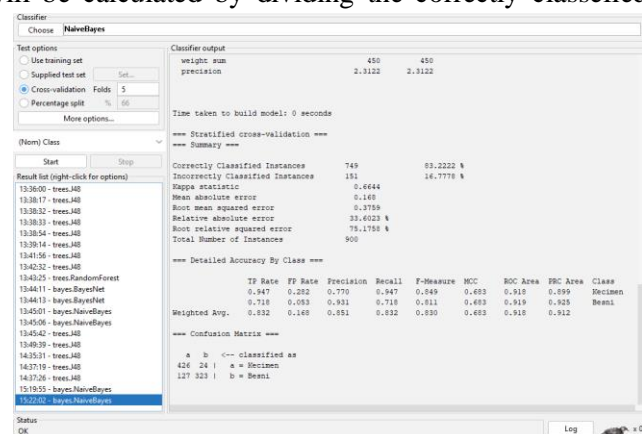
As an example Tp rate of the Kecimen class will be calculated by dividing the correctly classefied Kecimans with the total true Kecimens $\frac{426}{450} = 0.947$ .

FP rate will be calculated by dividing incorrectly classefied Keciman (Kecimans classefied as Besni) on the total Kekiman values $\frac{24}{450} = 0.282$ .

Precision ca be calculated by deviding TP over FP + TP $\frac{0.947}{0.282+0.947} = 0.77$ .
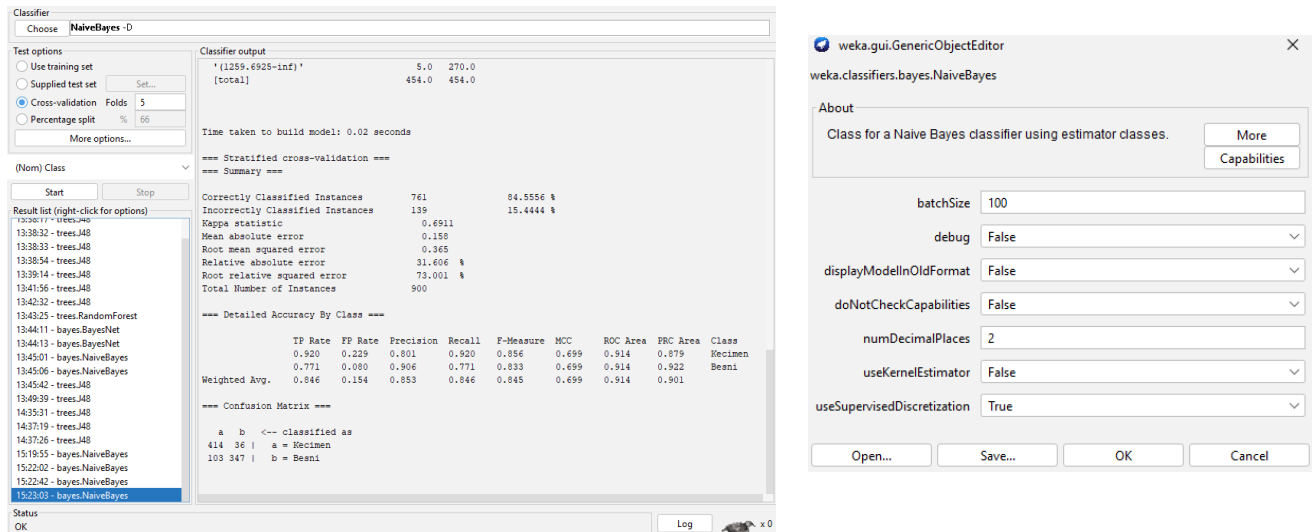
Recall should be equal to the TP.



F-Measure can be calculated by $\frac{2*\text{Precision}*\text{Recall}}{\text{Precision}+\text{Recall}} = \frac{2*0.77*0.947}{0.77+0.947} = 0.849$
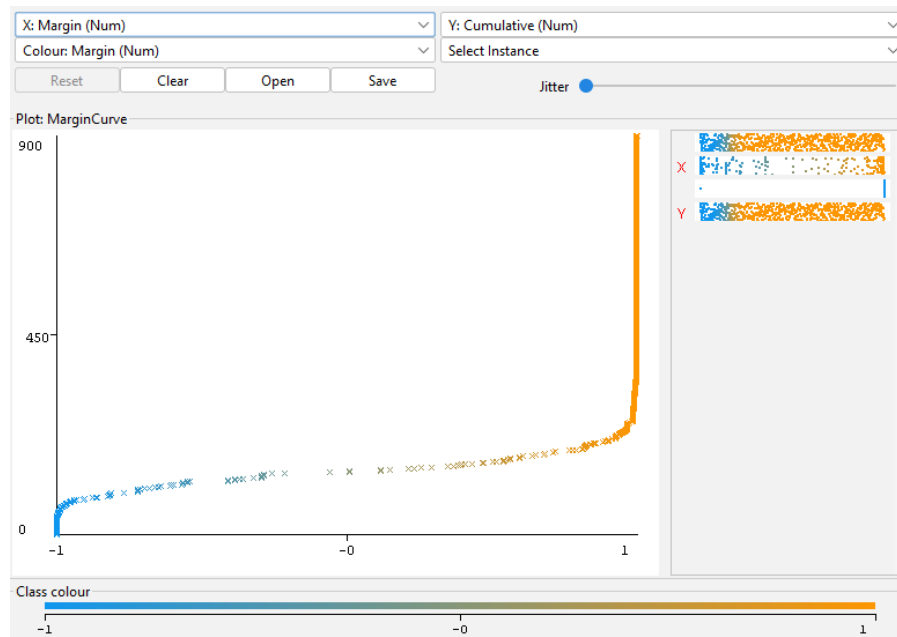
## Hyper-parameter

When the hyper parameter of the model is changed to supervised discretization as shown in figure below it can be noted that the changes in the values displayed on the screen such as correctly classified instances became higher and equal 84.56% and also the incorrectly classified instances changed and it became less and equal 15.44%. In addition, the effects are obvious for the confusion Matrix, TP Rate, FP Rate, Precision, Recall and F- Measure.
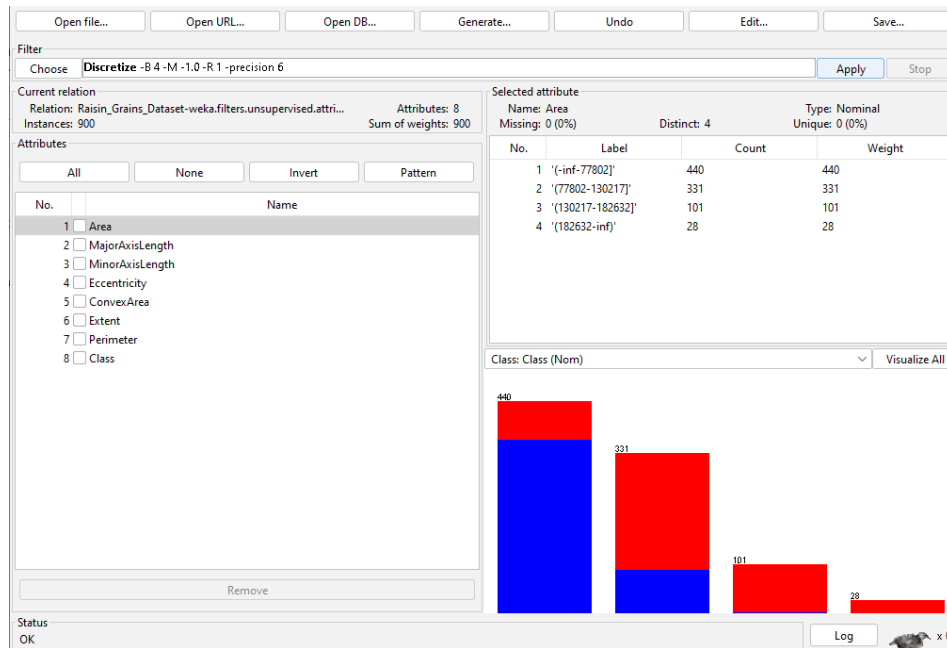


## Visualize the curve

### Use Supervised discretization

# 4. Lazy IBK

## Filter

Preprocess for one attribute (Area) discretization of continuous attributes with 4 pins.



## Classifer

In the confusion matrix shown below, it can be seen that the true kecimen have 368 that were classefied correctly and 82 that were classefied as Besni , the true Besni have 366 classefied correctly and 84 that were classefied as keciman.

It can be noticed that the correctly classefied have an accuracy of 81.56% and incorrectly classified have an accuracy of 18.44% .

In the Detailed Accuracy by class it shows the TP rate, FP Rate, Percision, F-Measure,MCC, ROC Area, PRC Area and the class.
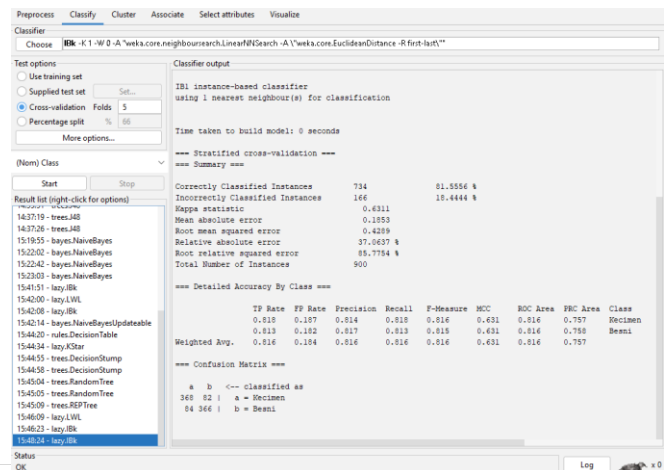
As an example Tp rate of the Kecimen class will be calculated by dividing the correctly classefied Kecimans with the total true Kecimens $\frac{368}{450} = 0.818$ .

FP rate will be calculated by dividing incorrectly classefied Keciman (Kecimans classefied as Besni) on the total Kekiman values $\frac{82}{450} = 0.187$ .

Precision ca be calculated by deviding TP over FP + TP $\frac{0.818}{0.187+0.818} = 0.814$ .
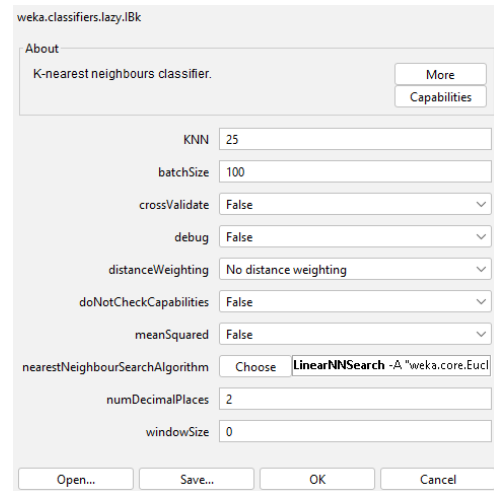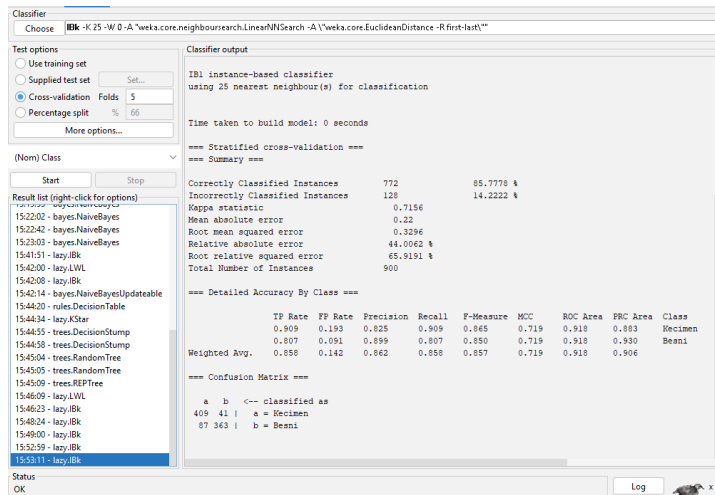
Recall should be equal to the TP.

F-Measure can be calculated by $\frac{2 * \text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} =$

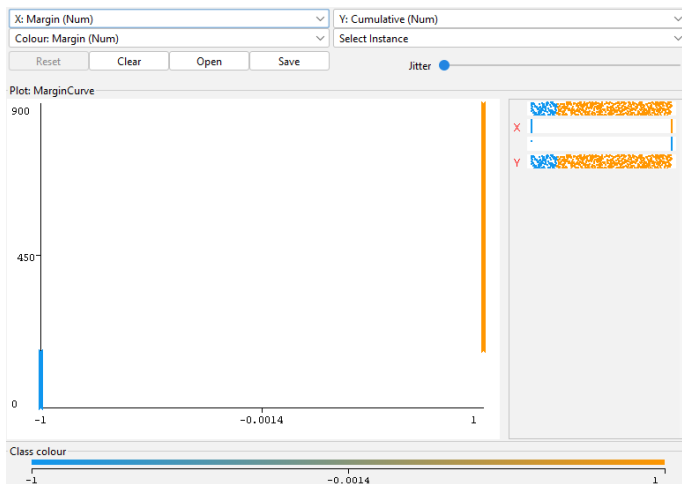$\frac{2 * 0.814 * 0.818}{0.814 + 0.818} = 0.816$

## Hyper-parameter

When the hyper parameter of the model is changed to KNN as shown in figure below it can be noted that the changes in the values displayed on the screen such as correctly classified instances became higher and equal 85.77% and also the incorrectly classified instances changed and it became less and equal 14.22%. In addition, the effects are obvious for the confusion Matrix, TP Rate, FP Rate, Precision, Recall and F-Measure.



## Visualize the curves

### Default



### Use KNN