# CPEN405 - ARTIFICIAL INTELLIGENCE
# COURSE PROJECT - GROUP 13

Adjei-Twum Emmanuel Sefa - 10717695
Twum Penuel Antwi - 10712173
Rebecca Seglah Sesinam - 10722022
Benjamin Powell - 10735628
Akesse Stephen King - 10732467
Osafo Kwadwo Agyeman - 10724754

June 13, 2022

# 1 INTRODUCTION

Artificial Intelligence is the ability of a digital computer or computer-controlled robot to perform tasks commonly associated with intelligent beings. The term is frequently applied to the project of developing systems endowed with the intellectual processes characteristics of humans such as the ability to reason, discover meaning, generalize or learn from the past. From the Oxford dictionary, Intelligence is defined as the ability to acquire and apply knowledge and skills. Psychologists generally do not characterize human intelligence by just one trait but by the combination of many diverse ability. Research in AI has focused chiefly on the following components of intelligence: learning, reasoning, problem solving, perception and using language. These are the definitive components that have helped shape the face of modern artificial intelligence. The aim of this project is to challenge ourselves to some problems in the fields of machine learning, natural language processing, knowledge representation, optimization and game theory.

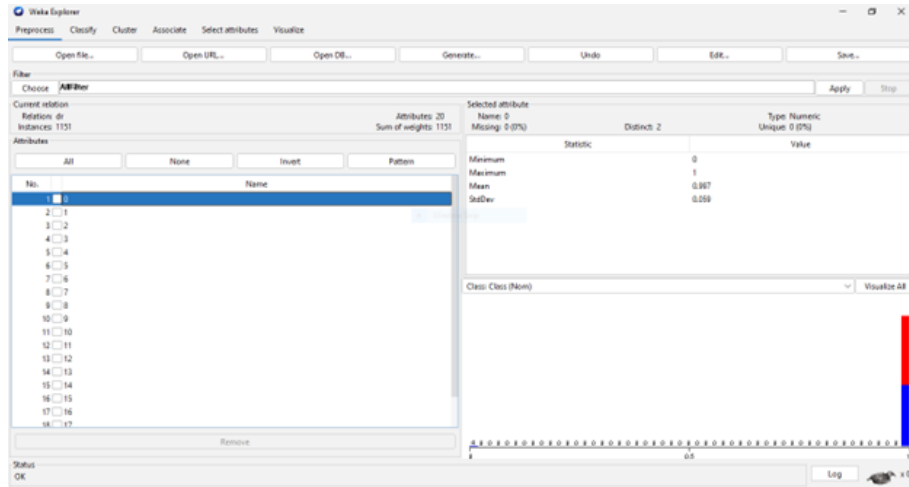# 2 PART 1 – MACHINE LEARNING USING WEKA

## 2.1 OBJECTIVES

The aim of this is to have some hands-on experience with machine learning algorithms using WEKA for solving real world data-mining problems.

## 2.2 PROBLEM FORMULATION

To evaluate and analyze the performance of at least three classification schemes on a Diabetic Retinopathy Debrecen data set. The data set was obtained from the open-source UCI Machine Learning Repository. The data set contains features extracted from the Messidor image set to predict whether an image contains signs of diabetic retinopathy or not.

## 2.3 SOLUTION APPROACH AND ALGORITHMS

The Diabetic Retinopathy Debrecen data set has 1151 instances and 20 attributes as shown in figure 1 below.

For a data set with over 1000 instances, 2/3 of the data set was used for training and 1/3 was used for testing. The classification schemes applied to the data set were J48 decision tree classification algorithm, the RepTree classification algorithm, Naives Bayes classification algorithm and KNN algorithm. A brief definition of the chosen algorithms and reasons they were chosen is given below.

### 2.3.1 J48 TREE CLASSIFICATION ALGORITHM

The J48 classification algorithm is one of the best machine learning algorithms to examine the data categorically and continuously. When it is used for instance purposes, it occupies more memory space and depletes the performance and accuracy in classifying medical data.

### 2.3.2 REDUCED ERROR PRUNING TREE CLASSIFICATION ALGORITHM

Reduced Error Pruning (RepTree) is a fast decision tree learner that builds a decision or regression tree using information gain as the splitting criterion and prunes it using reduced error pruning algorithm. As traversing over the internal nodes from the bottom to the top of a tree, the REP procedure checks for each internal node, whether replacing it with the most repeated class that does not reduce the accuracy of trees. In this case, the node is pruned. The procedure continues until any further pruning will decrease the accuracy.

### 2.3.3 NAÏVE BAYES CLASSIFICATION ALGORITHM

Naïve Bayes classification calculates explicit probabilities for superposition among the most practical approaches to definite types of learning problems. Even when Bayesian techniques are computationally difficult, they can produce standard and optimal decision making against all other methods of classification.

3

### 2.3.4 KNN CLASSIFICATION ALGORITHM

K-nearest neighbor algorithm is a simple, supervised machine learning algorithm that can be used to solve both classification and regression problems. To classify a new instance it searches the training set for one that is "most like" it. KNN algorithm decides a number k which is the nearest neighbor to the data point that is to be classified. If the value of k is 5 it will look for 5 nearest neighbors to that data point. It is easy to implement and understand but has a major drawback of becoming significantly slow as the size of data grows.

## 2.4 RESULTS AND DISCUSSION

Although having a high percentage of correctly classified instances is good, that alone cannot accurately determine if the model used is good. You can have high classification accuracy and a bad model. For instance, if you have 90 percent of instances all in one class you can say all of them belong to that class. And you're going to be right 90 percent of the time but the model is not good.

This section discusses and analyzes the statistics of the output of the models of the different algorithms that were applied and the changes in the results as the parameters were varied. In order to do that there are certain concepts we need to understand.

Kappa statistic takes into account the fact that you could randomly guess and correctly classify an instance. It is a statistic of how well you would perform if you take into consideration the fact that you can randomly guess the class of an instance

Kappa = (observed accuracy - expected accuracy)/(1 - expected accuracy

Weighted average gives a general average of how well the algorithm is doing.

True positive is the percentage of classified instances that actually belong to a particular class.

False positive is the percentage of classified instances that do not actually belong to a particular class.

Precision is the ratio of true positives to total data classified as a particular class.

Recall is ratio of true positives to total data that actually belongs to a particular class, say class 0
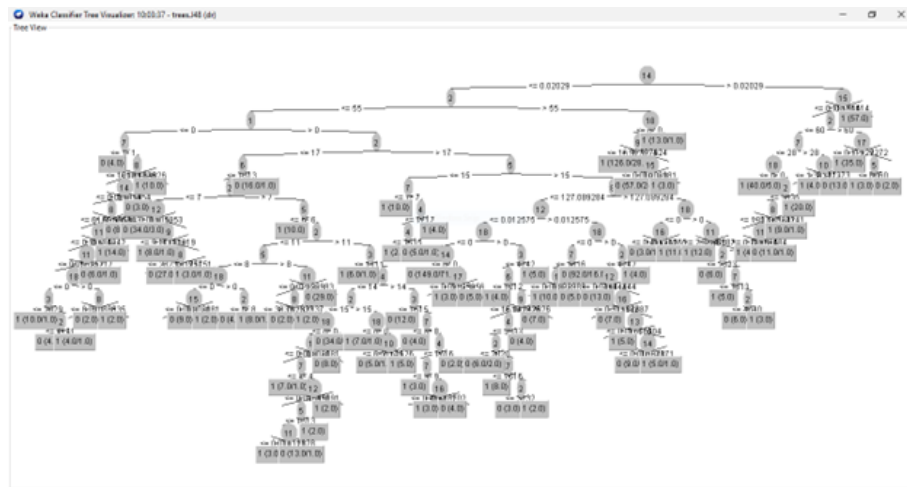
ROC area is the receiver operator characteristic area under the curve. It shows what percentage of the time you are going to correctly classify items if given one of each class.

PRC area is the precision and recall area and it does not account for true positives in the statistics.
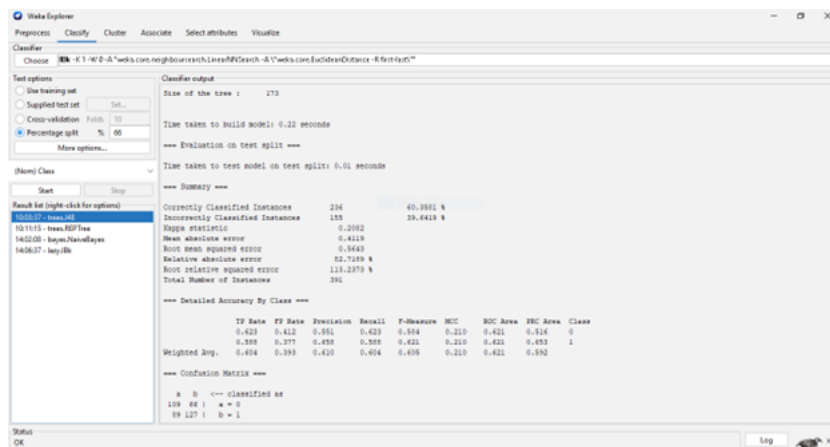
Confusion matrix shows values that are correctly classified for each class and those that are incorrectly classified.

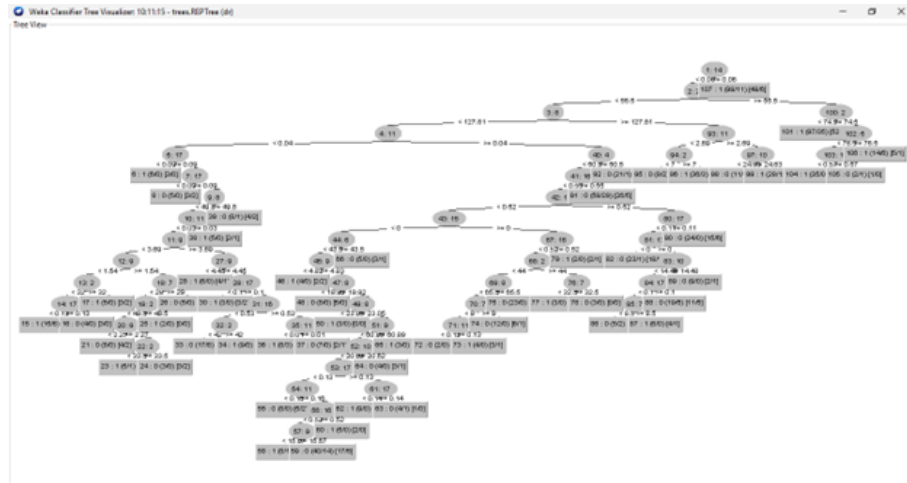### 2.4.1  J48 TREE CLASSIFICATION ALGORITHM

The visualized tree for the J48 algorithm on the dataset is given below:



The output for the classifier model is given below:

### 2.4.2 REDUCED ERROR PRUNING TREE CLASSIFICATION ALGORITHM

### 2.4.3 NAÏVE BAYES CLASSIFICATION ALGORITHM

**Weka Explorer**

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

Classifier

Choose  IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""

Test options
- Use training set
- Supplied test set    Set...
- Cross-validation  Folds  10
- Percentage split   %  66
- More options...

(Nom) Class

Start | Stop

Result list (right-click for options)
10:10:37 - trees.J48
10:11:15 - trees.REPTree
14:52:08 - bayes.NaiveBayes
14:06:37 - lazy.IBk

Classifier output

```
Time taken to build model: 0.02 seconds

=== Evaluation on test split ===

Time taken to test model on test split: 0.02 seconds

=== Summary ===

Correctly Classified Instances       204               52.1739 %
Incorrectly Classified Instances     187               47.8261 %
Kappa statistic                        0.1087
Mean absolute error                    0.475
Root mean squared error                0.6522
Relative absolute error               95.3079 %
Root relative squared error          130.5924 %
Total Number of Instances            391

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC  ROC Area  PRC Area  Class
                 0.914    0.796    0.482      0.914   0.631      0.164  0.688    0.612     0
                 0.204    0.086    0.766      0.204   0.320      0.164  0.687    0.726     1
Weighted Avg.    0.522    0.404    0.625      0.522   0.459      0.164  0.687    0.675

=== Confusion Matrix ===

   a   b   <-- classified as
 160  15 |   a = 0
 172  44 |   b = 1
```

Status
OK        Log

---

**Weka Explorer**

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

Classifier

Choose  IBk -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A \"weka.core.EuclideanDistance -R first-last\""

Test options
- Use training set
- Supplied test set    Set...
- Cross-validation  Folds  10
- Percentage split   %  66
- More options...

(Nom) Class

Start | Stop

Result list (right-click for options)
10:10:37 - trees.J48
10:11:15 - trees.REPTree
14:52:08 - bayes.NaiveBayes
14:06:37 - lazy.IBk

Classifier output

```
9
  mean          23.074   23.1
  std. dev.     19.7019  23.1387
  weight sum      540      611
  precision     0.1466   0.1466

10
  mean          0.2346   9.1221
  std. dev.     10.5544  12.3712
  weight sum      540      611
  precision     0.094    0.094

11
  mean          1.4003   2.2213
  std. dev.     2.7919   4.6671
  weight sum      540      611
  precision     0.055    0.055

12
  mean          0.1841   0.091
  std. dev.     0.5542   3.3326
  weight sum      540      611
  precision     0.0648   0.0648

13
  mean          0.0413   0.3623
  std. dev.     0.157    1.4264
  weight sum      540      611
  precision     0.0348   0.0348

14
  mean          0.0066   0.1546
```
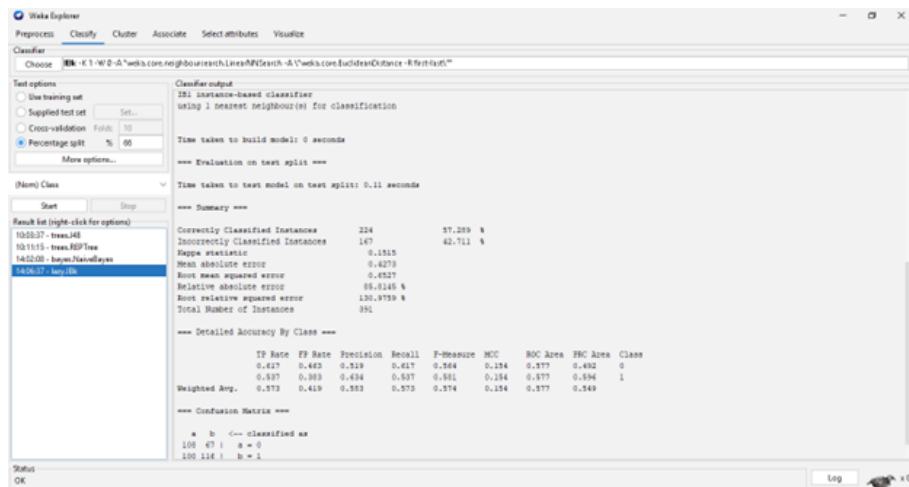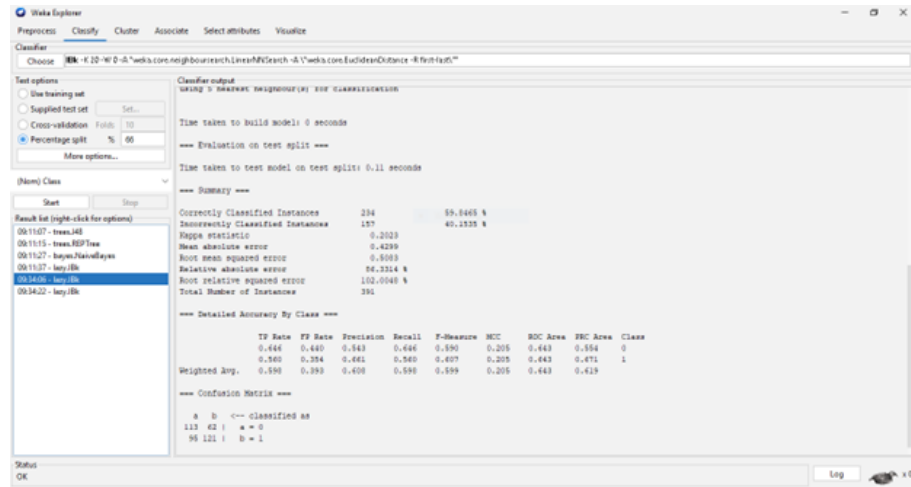
Status
OK        Log

8

### 2.4.4 KNN CLASSIFICATION ALGORITHM

The KNN algorithm has no model. When the algorithm is applied to the dataset with a KNN value of 1, the results are shown in figure below.
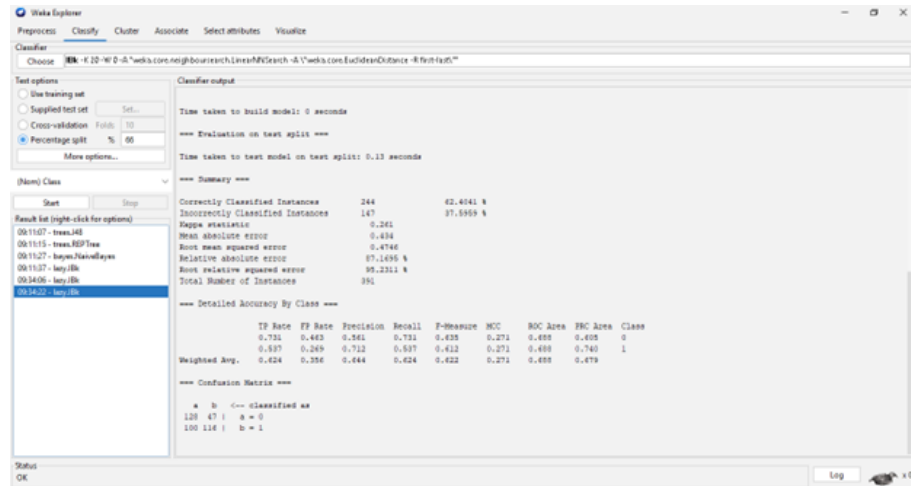


The number of correctly classified Instances were 224 at a 57.289 percent accuracy and the number of incorrectly classified instances were 167 at a 42.711percent accuracy. It has a kappa statistic of 0.1515. 0.1515 is a poor kappa statistic which is close to random guessing. The algorithm has a weighted true positive rate of 0.573 and a weighted false positive rate of 0.419. It has a weighted precision rate of 0.583 and a recall rate of 0.573. It has an F-measure of 0.573 and an ROC area of 0.577.

For the confusion matrix, 108 instances of class a were correctly classified and 67 instances of class a were incorrectly classified as class b = 1. 100 instances of class b were incorrectly classified as class a and 116 instances of class b were correctly classified. From the above statistics is can be determined that the model is not good and is close to random guessing. The KNN value was then changed to 5 which gave the following results:



The percentage of correctly classified instances improves to 59.8465 percent and the precision, recall, f-measure and ROC area rates improve as well but the model still is not good enough. The KNN value was varied to 20 and the results are shown below:

The percentage of correctly classified instances improves to 62.4041 percent and the precision, recall, f-measure and ROC area rates improve significantly. It can then be observed that for such a noisy dataset the accuracy of the model improves as the KNN value increases.

## 2.5    REFERENCES

anaN, N. S., & thri, V. G. (2018). Performance and Classification Evaluation of J48 Algorithm and Kendall's Based J48 Algorithm (KNJ48). International Journal of Computer Trends and Technology, 59(2), 73–80. https://doi.org/10.14445/22312803/ijctt-v59p112

Bishop, C. M. (2006). Bishop, C. M. (2006). Pattern Recognition and Machine Learning. (M. Jordan, J. Kleinberg, & B. Schölkopf, Eds.)Pattern Recognition (Vol. 4, p. 738). Springer. doi:10.1117/1.2819119Pattern Recognition and Machine Learning. (M. Jordan, J. Kleinberg, & B. Schölkopf, Eds.), Pattern Recognition (Vol. 4, p. 738). Springer. Retrieved from http://www.library.wisc.edu/selectedtocs/bg0137.pdf

artificial intelligence - Evolutionary computing — Britannica. (n.d.). Retrieved from www.britannica.com: https://www.britannica.com/technology/artificial-intelligence/Evolutionary-computing

Weka - classifiers. (n.d.). Retrieved from www.tutorialspoint.com: https://www.tutorialspoint.com/weka/weka

Weka Decision Tree — Build Decision Tree Using Weka. (n.d.). Retrieved from www.analyticsvidhya.com: https://www.analyticsvidhya.com/blog/2020/03/decision-tree-weka-no-coding/