



RED WINE QUALITY

Machine Learning



Name	ID
Rubaiyad Noor Shahriar	19-39541-1
Muhammad, Akib	19-39372-1
Md. Adnan Shakib	18-38015-2
Md Amir Habib	19-41429-3



JUNE 14, 2022

AMERICAN INTERNATIONAL UNIVERSITY BANGLADESH
408/1, Kuratoli, Khilkhet, Dhaka 1229, Bangladesh

Contents

List of Figures	2
List of Tables	2
Notations.....	2
Introduction	3
Background study	3
Problem Statement.....	3
Objective	4
Dataset	4
Data Information.....	4
Model Development	5
Naïve Bayes.....	5
K-Nearest Neighbor (KNN) Algorithm	5
Decision Trees	6
Graphical Representations.....	6
Variable value	6
Evaluating Training set.....	10
Naïve Bayes Algorithm Implementation.....	11
Decision Tree Implementation	12
K-Nearest Neighbor (KNN) Algorithm Implementation	18
Discussion & Conclusion	19
Naïve Bayes	19
K-Nearest neighbor Algorithm	19
Decision Tree Algorithm.....	19
Comparative analysis	20
References	20

List of Figures

1. Naïve bayes classification algorithm results from Weka tool
2. Decision Tree classification algorithm results from Weka tool
3. K-Nearest Neighbor classification algorithm from Weka tool
4. Variable values ranges
5. Confusion matrices of the 3 implemented algorithms

List of Tables

1. Accuracy comparison table

Notations

- $P(c/x)$ = the posterior probability of *class* (c , *target*) given *predictor* (x , *attributes*).
- $P(c)$ = the prior probability of *class*.
- $P(x/c)$ = the likelihood which is the probability of *predictor* given *class*.
- $P(x)$ = the prior probability of *predictor*.

Introduction

Background study

Wine is an alcoholic drink typically made from fermented grapes. Yeast consumes the sugar in the grapes and converts it to ethanol and carbon dioxide, releasing heat in the process. Different varieties of grapes and strains of yeasts are major factors in different styles of wine. These differences result from the complex interactions between the biochemical development of the grape, the reactions involved in fermentation, the grape's growing environment (terroir), and the wine production process. Many countries enact legal appellations intended to define styles and qualities of wine. These typically restrict the geographical origin and permitted varieties of grapes, as well as other aspects of wine production. Wines not made from grapes involve fermentation of other crops including rice wine and other fruit wines such as plum, cherry, pomegranate, currant and elderberry.

Wine has been produced for thousands of years. The earliest evidence of wine is from ancient China (c. 7000 BC) Georgia (6000 BC), Persia (5000 BC), and Italy (4000 BC). New World wine has some connection to alcoholic beverages made by the indigenous peoples of the Americas, but is mainly connected to later Spanish traditions in New Spain. Later, as Old-World wine further developed viticulture techniques, Europe would encompass three of the largest wine-producing regions. Today, the five countries with the largest wine-producing regions are in Italy, Spain, France, the United States, and China.

Wine has long played an important role in religion. Red wine was associated with blood by the ancient Egyptians and was used by both the Greek cult of Dionysus and the Romans in their Bacchanalia; Judaism also incorporates it in the Kiddush, and Christianity in the Eucharist. Egyptian, Greek, Roman, and Israeli wine cultures are still connected to these ancient roots. Similarly, the largest wine regions in Italy, Spain, and France have heritages in connection to sacramental wine, likewise, viticulture traditions in the Southwestern United States started within New Spain as Catholic friars and monks first produced wines in New Mexico and California.

Problem Statement

The two datasets are related to red and white variants of the Portuguese "Vinho Verde" wine. For more details, consult the reference [Cortez et al., 2009]. Due to privacy and logistic issues, only physicochemical (inputs) and sensory (the output) variables are available (e.g., there is no data about grape types, wine brand, wine selling price, etc.).

These datasets can be viewed as classification or regression tasks. The classes are ordered and not balanced (e.g., there are much more normal wines than excellent or poor ones).

We are going to classify this dataset according to the value of the quality

Objective

Our objective of this project is to classify the dataset for red wine quality. We will be using Naïve bayes algorithm, K-Nearest Neighbor algorithm, Decision tree algorithm to classify the dataset. The class mentioned as “*quality*”. There are 5 classes which are numerical but determines the quality of the wine. By using the three-classification algorithm we will provide a comparative discussion on which classification algorithm provides the most efficient result.

Dataset

Data Information

fixed acidity: most acids involved with wine or fixed or nonvolatile (do not evaporate readily).

volatile acidity: the amount of acetic acid in wine, which at too high of levels can lead to an unpleasant, vinegar taste.

citric acid: found in small quantities, citric acid can add 'freshness' and flavor to wines.

residual sugar: the amount of sugar remaining after fermentation stops, it's rare to find wines with less than 1 gram/liter and wines with greater than 45 grams/liter are considered sweet.

chlorides: the amount of salt in the wine.

free sulfur dioxide: the free form of SO₂ exists in equilibrium between molecular SO₂ (as a dissolved gas) and bisulfite ion; it prevents microbial growth and the oxidation of wine.

total sulfur dioxide: amount of free and bound forms of SO₂; in low concentrations, SO₂ is mostly undetectable in wine, but at free SO₂ concentrations over 50 ppm, SO₂ becomes evident in the nose and taste of wine.

density: the density of wine is close to that of water depending on the percent alcohol and sugar content.

pH: describes how acidic or basic a wine is on a scale from 0 (very acidic) to 14 (very basic); most wines are between 3-4 on the pH scale.

sulphates: a wine additive which can contribute to sulfur dioxide gas (SO₂) levels, which acts as an antimicrobial and antioxidant.

alcohol: the percent alcohol content of the wine.

quality: output variable (based on sensory data, score between 0 and 10)

Model Development

Naïve Bayes

It is a classification technique based on Bayes' Theorem with an assumption of independence among predictors. In simple terms, a Naive Bayes classifier assumes that the presence of a particular feature in a class is unrelated to the presence of any other feature.

For example, a fruit may be considered to be an apple if it is red, round, and about 3 inches in diameter. Even if these features depend on each other or upon the existence of the other features, all of these properties independently contribute to the probability that this fruit is an apple and that is why it is known as 'Naive'.

Naive Bayes model is easy to build and particularly useful for very large data sets. Along with simplicity, Naive Bayes is known to outperform even highly sophisticated classification methods.

Bayes theorem provides a way of calculating posterior probability $P(c|x)$ from $P(c)$, $P(x)$ and $P(x|c)$. Look at the equation below:

$$P(c|x) = \frac{P(x|c) P(c)}{P(x)}$$

Above,

- $P(c|x)$ is the posterior probability of *class (c, target)* given *predictor (x, attributes)*.
- $P(c)$ is the prior probability of *class*.
- $P(x|c)$ is the likelihood which is the probability of *predictor* given *class*.
- $P(x)$ is the prior probability of *predictor*.

K-Nearest Neighbor (KNN) Algorithm

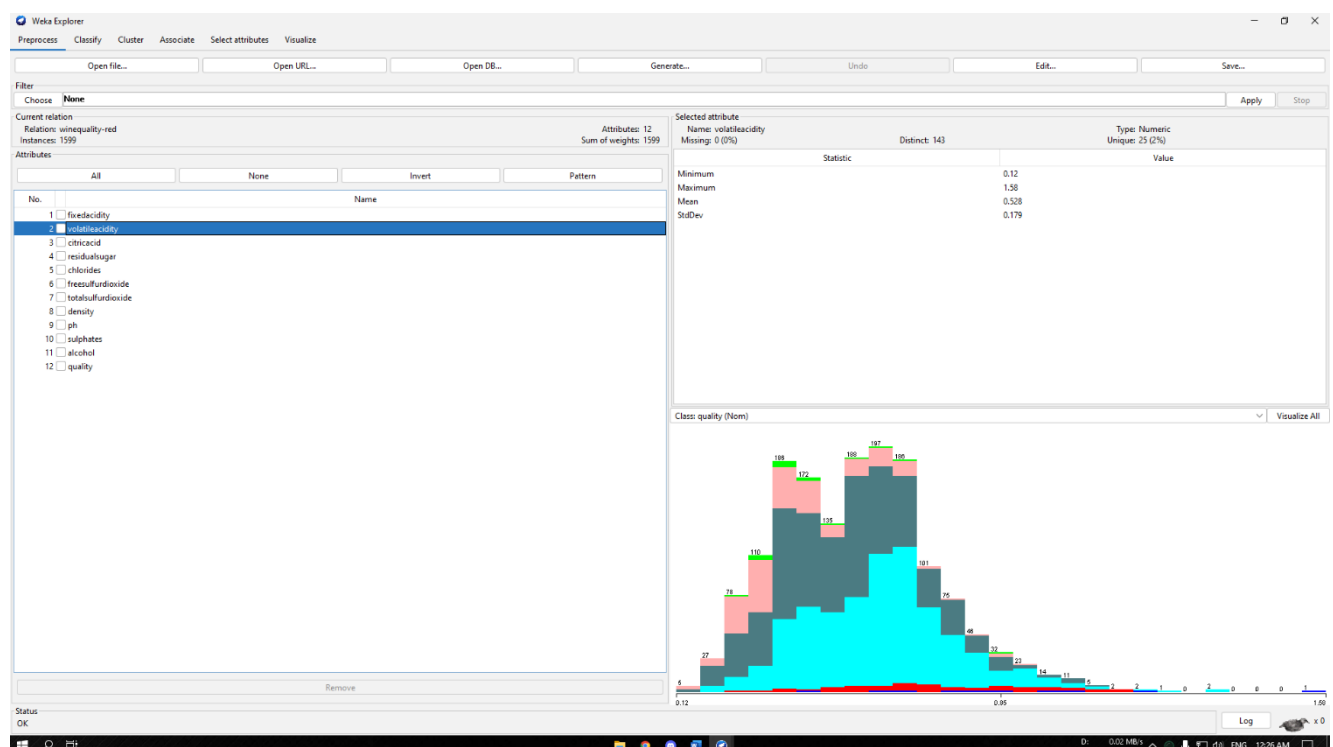
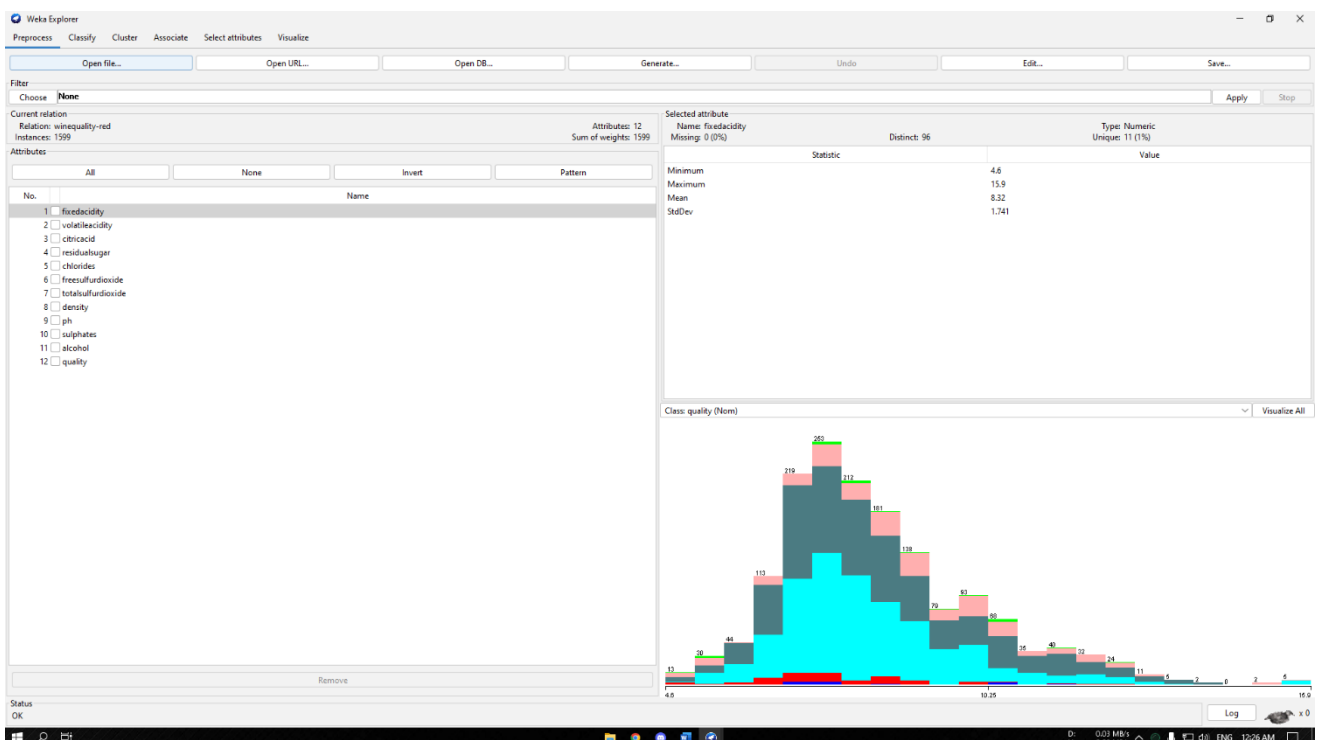
K-Nearest Neighbor is one of the simplest Machine Learning algorithms based on Supervised Learning technique. K-NN algorithm assumes the similarity between the new case/data and available cases and put the new case into the category that is most similar to the available categories. K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well suite category by using K- NN algorithm. K-NN algorithm can be used for Regression as well as for Classification but mostly it is used for the Classification problems. K-NN is a non-parametric algorithm, which means it does not make any assumption on underlying data. It is also called a lazy learner algorithm because it does not learn from the training set immediately instead it stores the dataset and at the time of classification, it performs an action on the dataset. KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to the new data.

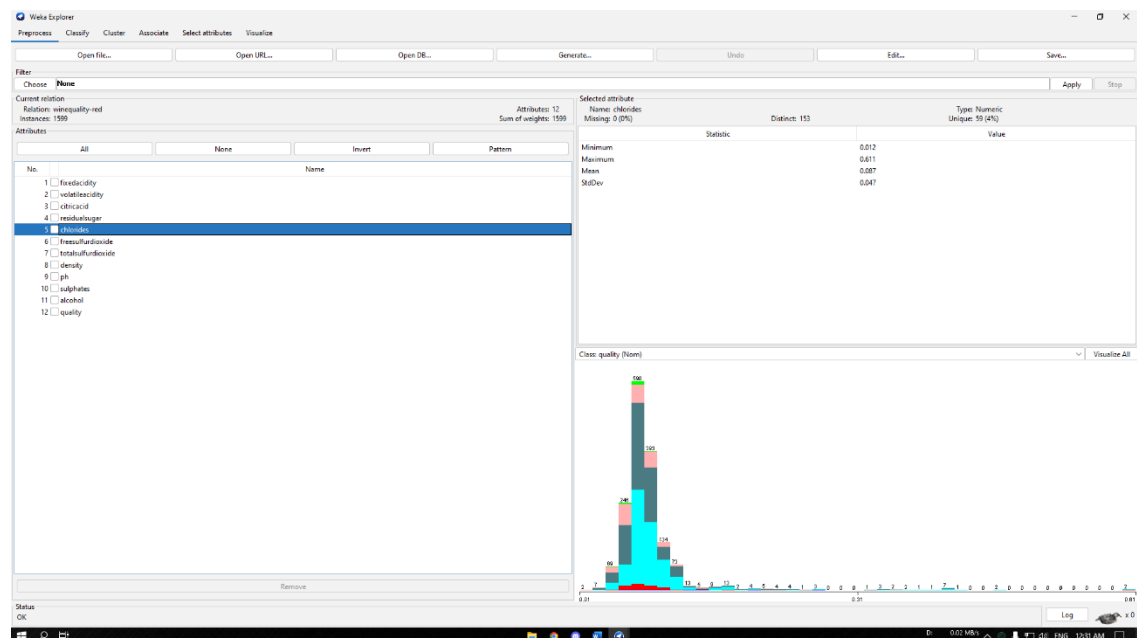
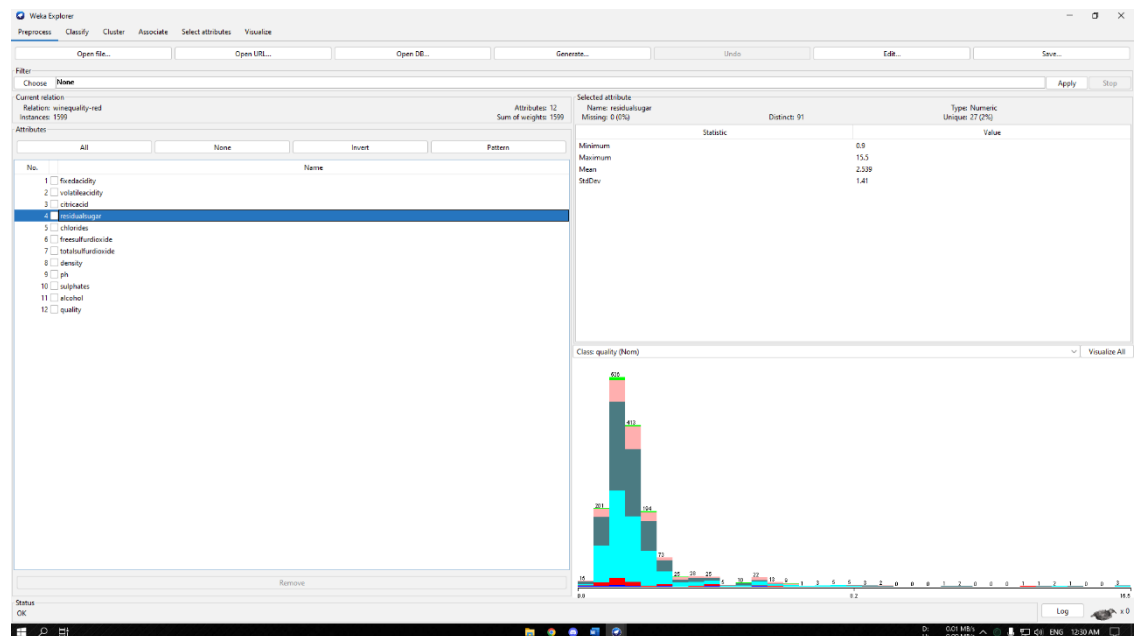
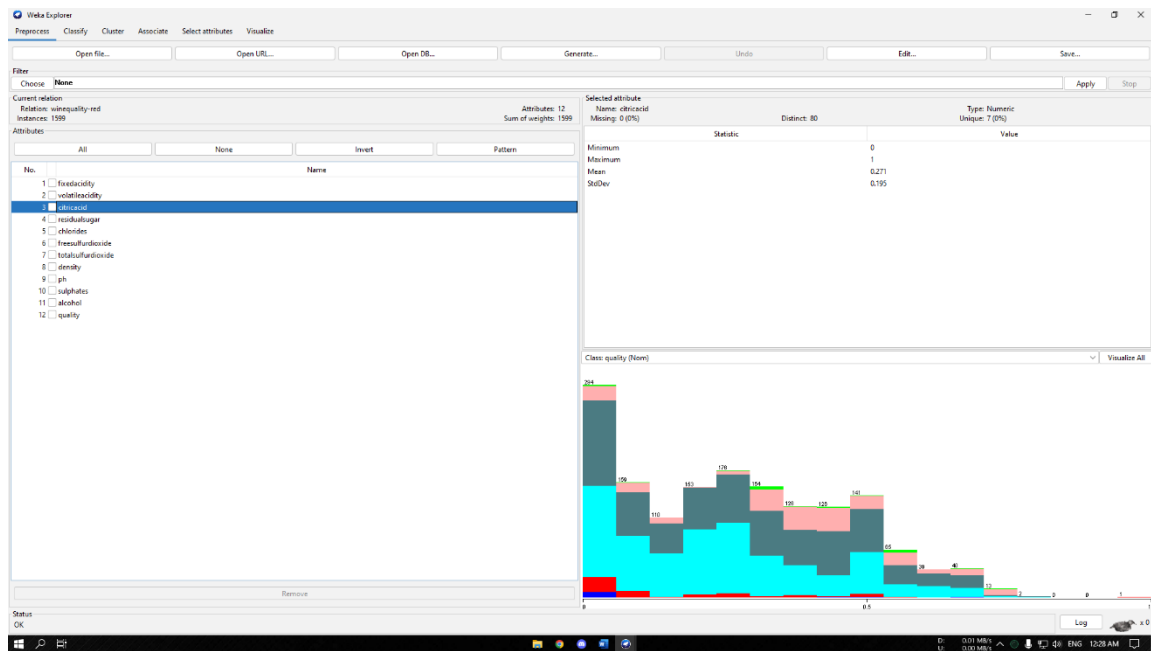
Decision Trees

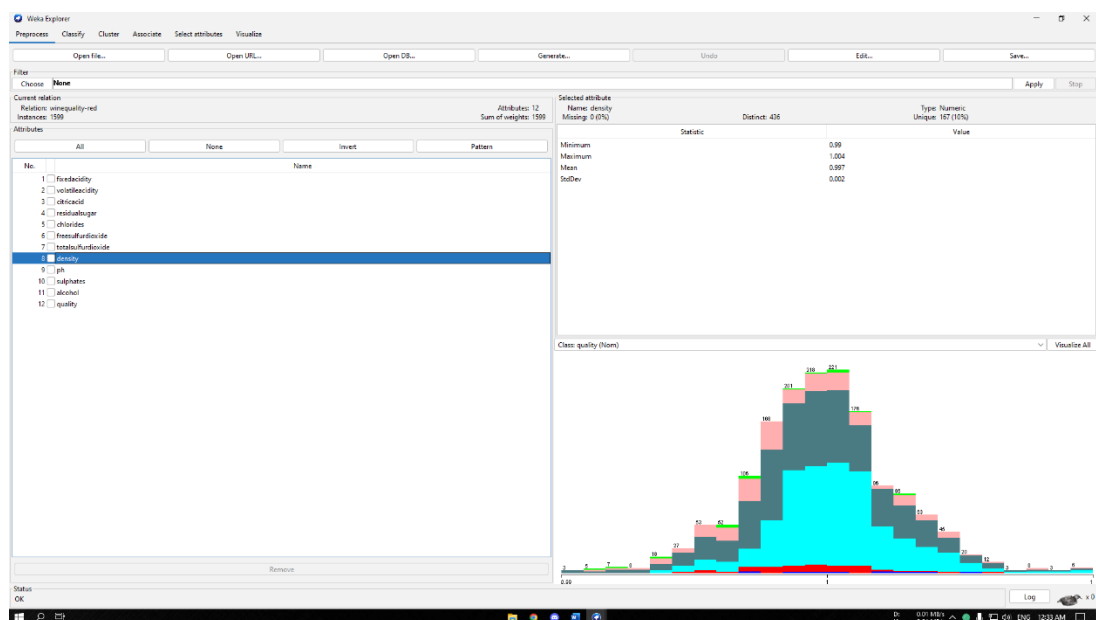
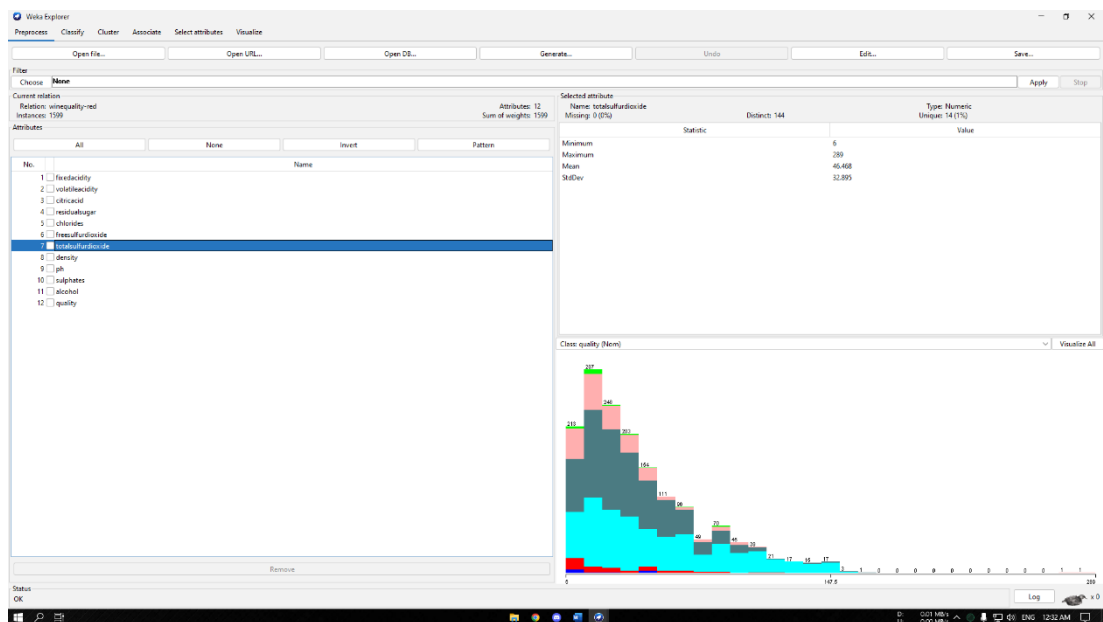
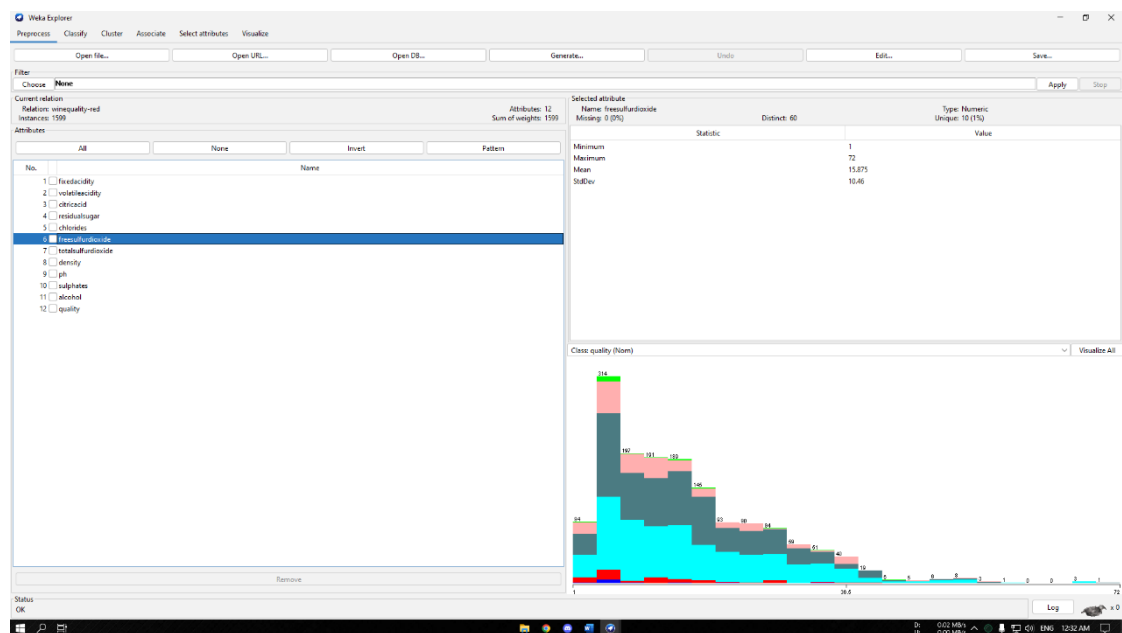
Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from the data features. A tree can be seen as a piecewise constant approximation.

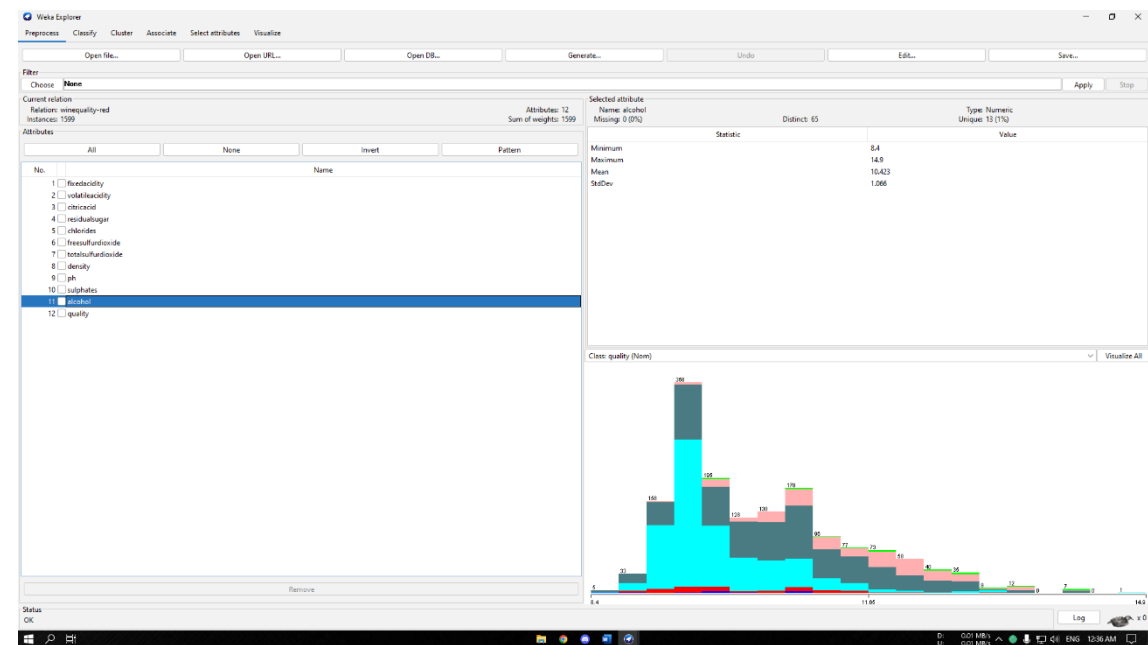
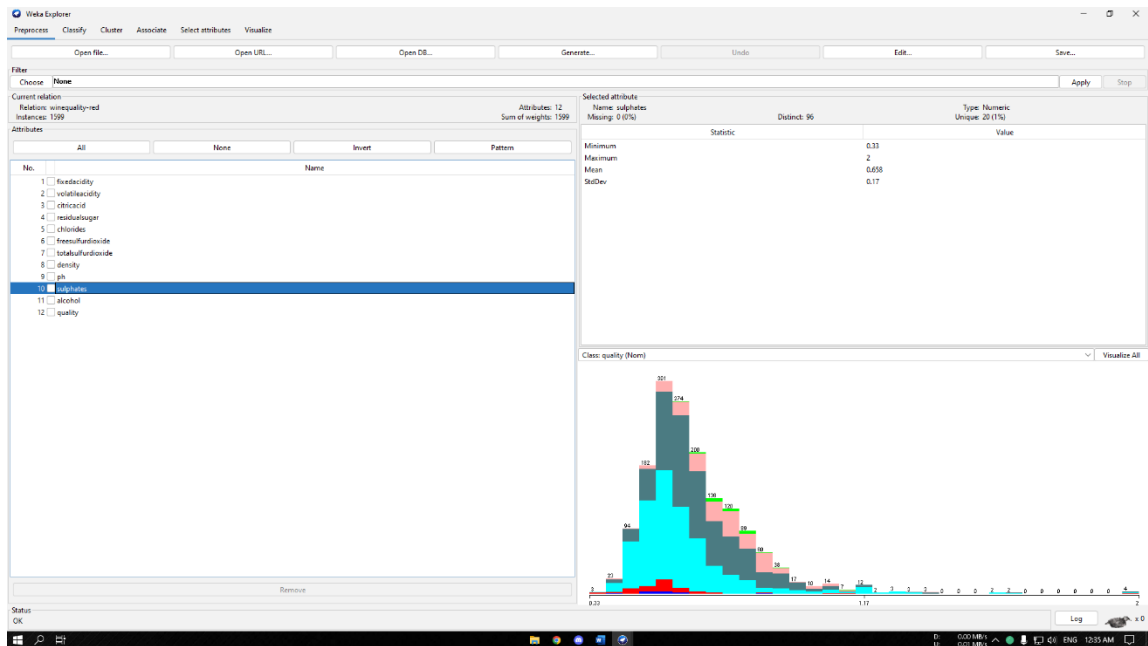
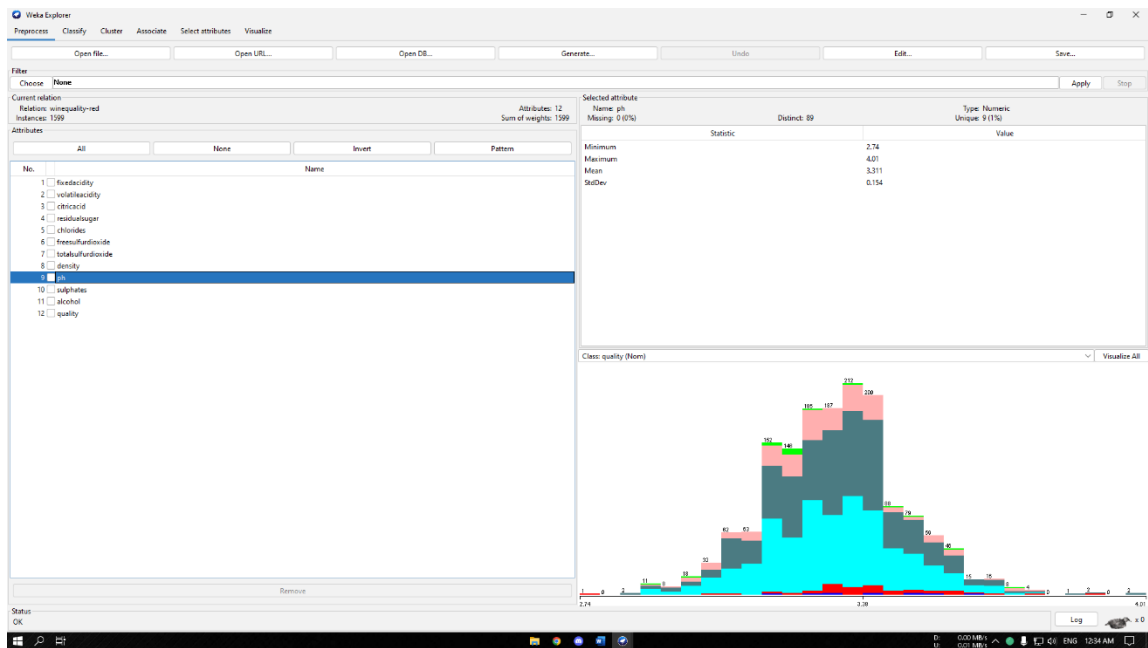
Graphical Representations

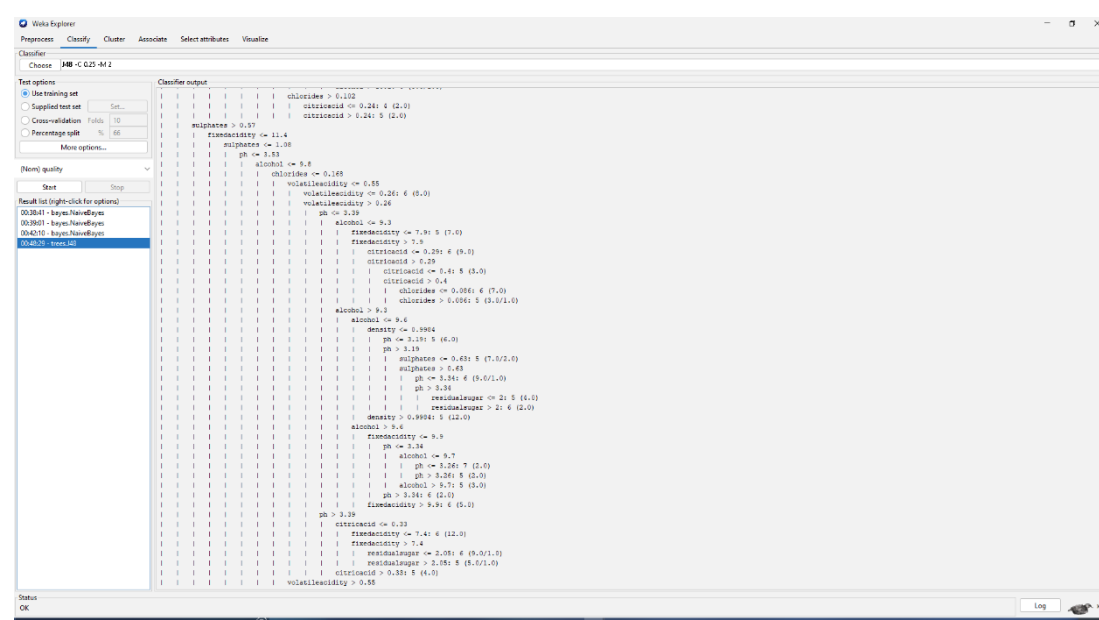
Variable value

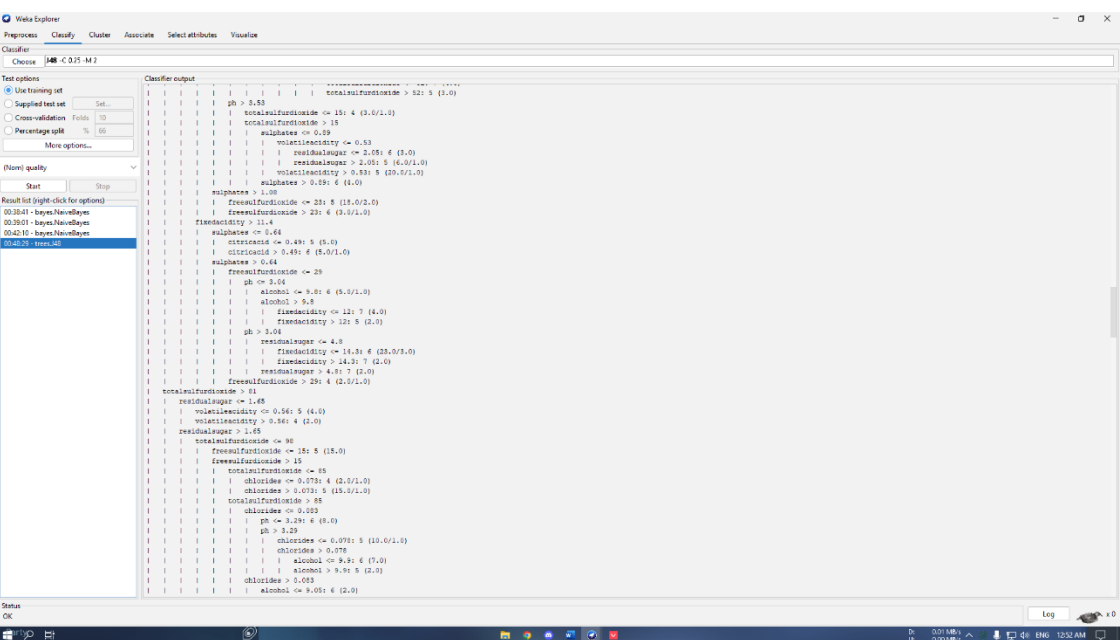












Weka Explorer

Preprocesses Classify Cluster Associate Select attributes Visualize

Classifier

Choice **AH** C 0.25 M 2

Test options

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

(Nom) quality

Start Stop

Result list (right-click for options)

Class	Count	%	Mean	StdDev	Min	Max	InfoGain	ReliefF	GainRatio	MSE	MAE	RMSE	MAPE	Kappa	F1	Precision	Recall	Accuracy
00:04:1 haveNameDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
00:00:0 haveNameDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNameDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
00:45:0 haveNoDays	1	1%	0.0	0.0	0.0													

[illegible]

K-Nearest Neighbor (KNN) Algorithm Implementation

KNN is conceptually simple and has the advantage of being nonparametric. That is, the method can be used even when the variables are categorical—though if we are using numeric variables in the mix, it is best to standardize them to eliminate differences in scale. The challenge is that when the number of data points is very large special methods must be employed to rapidly search the space and find the “most similar” items.

The screenshot shows the Weka Explorer interface with the K-Nearest Neighbor (KNN) classifier selected. The classifier output is displayed, showing the model's performance metrics and a detailed accuracy table.

Classifier output

```
=== Run information ===
Scheme:   weka.classifiers.lazy.IBK -K 1 -W 0 -A "weka.core.neighboursearch.LinearNNSearch -A "weka.core.EuclideanDistance -B first-last""
Relation: winequality-red
Instances: 1599
Attributes: 12
  fructose
  volatileacidity
  citricacid
  residualalugar
  chlorides
  freesulfur dioxide
  totalsulfur dioxide
  density
  ph
  sulphates
  alcohol
  quality

Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===
IBK instance-based classifier
using 1 nearest neighbour(s) for classification

Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      1036      64.7905 %
Incorrectly Classified Instances    563      35.2095 %
Kappa statistic                    0.4507
Mean absolute error                 0.118
Root mean squared error             0.3419
Relative absolute error             55.0092 %
Root relative squared error        104.4764 %
Total Number of Instances          1599

=== Detailed Accuracy By Class ===
```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.000	0.005	0.000	0.000	0.000	0.000	-0.006	0.565	0.007	3
0.075	0.028	0.085	0.075	0.080	0.081	0.534	0.038	4	
0.746	0.211	0.724	0.746	0.735	0.833	0.764	0.647	5	
0.636	0.223	0.655	0.636	0.645	0.416	0.700	0.569	6	
0.583	0.066	0.555	0.583	0.569	0.506	0.766	0.377	7	
0.111	0.007	0.154	0.111	0.129	0.122	0.500	0.027	8	
Weighted Avg.	0.648	0.188	0.643	0.648	0.645	0.459	0.727	0.551	

The screenshot shows the Weka Explorer interface with the K-Nearest Neighbor (KNN) classifier selected. The classifier output is displayed, showing the model's performance metrics and a detailed accuracy table, including a confusion matrix.

Classifier output

```
=== Classifier model (full training set) ===
IBK instance-based classifier
using 1 nearest neighbour(s) for classification

Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      1036      64.7905 %
Incorrectly Classified Instances    563      35.2095 %
Kappa statistic                    0.4507
Mean absolute error                 0.118
Root mean squared error             0.3419
Relative absolute error             55.0092 %
Root relative squared error        104.4764 %
Total Number of Instances          1599

=== Detailed Accuracy By Class ===
```

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
0.000	0.005	0.000	0.000	0.000	0.000	-0.006	0.565	0.007	3
0.075	0.028	0.085	0.075	0.080	0.081	0.534	0.038	4	
0.746	0.211	0.724	0.746	0.735	0.833	0.764	0.647	5	
0.636	0.223	0.655	0.636	0.645	0.416	0.700	0.569	6	
0.583	0.066	0.555	0.583	0.569	0.506	0.766	0.377	7	
0.111	0.007	0.154	0.111	0.129	0.122	0.500	0.027	8	
Weighted Avg.	0.648	0.188	0.643	0.648	0.645	0.459	0.727	0.551	

Confusion Matrix

```

a b c d e f <-- classified as
0 4 5 1 0 0 | a = 3
3 4 26 17 3 0 | b = 4
4 17 508 133 19 0 | c = 5
1 20 142 406 64 5 | d = 6
0 2 19 56 116 6 | e = 7
0 0 2 7 7 2 | f = 8
```

Discussion & Conclusion

The analysis of the dataset using the WEKA tool comparison among data mining classification algorithms (Decision tree, KNN, Naive Bayes), shows that all KNN algorithms are more accurate and they have less error rate and they are easier algorithms as compared to the Decision tree and Naive Bayes. The result of implementation in WEKA on the same dataset showed that the Decision Tree outperforms and Bayesian classification are less than the accuracy of KNN. The comparative study has shown that each algorithm has its own set of advantages and disadvantages as well as its own area of implementation. None of the algorithms can satisfy all constraints and criteria. Depending on the application and requirements, a specific algorithm can be chosen. We think KNN will be the right choice for this dataset according to the WEKA result.

Naïve Bayes

```
=== Confusion Matrix ===

  a   b   c   d   e   f   <-- classified as
1   4   4   1   0   0 |   a = 3
0   8  29  14   1   1 |   b = 4
6  26 455 168  26   0 |   c = 5
1  20 189 308 113   7 |   d = 6
0   2  14  73 107   3 |   e = 7
0   0   0   4  13   1 |   f = 8
```

K-Nearest neighbor Algorithm

```
=== Confusion Matrix ===

  a   b   c   d   e   f   <-- classified as
0   4   5   1   0   0 |   a = 3
3   4  26  17   3   0 |   b = 4
4  17 508 133  19   0 |   c = 5
1  20 142 406  64   5 |   d = 6
0   2  19  56 116   6 |   e = 7
0   0   2   7   7   2 |   f = 8
```

Decision Tree Algorithm

```
=== Confusion Matrix ===

  a   b   c   d   e   f   <-- classified as
1   3   2   3   1   0 |   a = 3
4   7  23  16   3   0 |   b = 4
2  24 478 159  18   0 |   c = 5
3  19 168 393  50   5 |   d = 6
1   4  22  65 100   7 |   e = 7
0   0   0  11   7   0 |   f = 8
```

Comparative analysis

Total number of instances **1599**

Perspective	Algorithm	Naïve Bayes	K-Nearest Neighbor	Decision Tree
Correctly Classified Instances		880	1036	979
Incorrectly Classified Instances		719	563	620
Relative absolute Error		82.1845%	55.0092%	63.6062%
Root relative squared Error		97.7148%	104.4764%	101.8383%

Table: Accuracy comparison table

References

1. <https://www.kaggle.com/datasets/uciml/red-wine-quality-cortez-et-al-2009>
2. https://datauab.github.io/red_wine_quality/
3. <https://drive.google.com/file/d/1CiRI6IqOM77SdjKvNGAJCdAzUCDBaXWG/view?usp=sharing>