



## Evaluating the performance of several algorithms using WEKA in red wine quality

Rhowel M. Delloso

Faculty Member, College of Computing, Information and Communications Technology, Isabela State University, Echague, Isabela, Philippines

### ABSTRACT

The study aimed to evaluate the performance of various regression algorithms performed at predicted values. WEKA program and a dataset from a reputable source were utilized. The study was initiated to compare the performance of several regression algorithms. The result shows that the KStar ( $r = 0.6043$ ) has the highest correlation coefficient, followed by Gaussian Processes ( $r = 0.5908$ ) and followed by M5Rules ( $r = 0.5904$ ). It also shows that the Linear Regression ( $r = 0.5872$ ) has good results as well as the SMOReg ( $r = 0.5623$ ). All of the 11 functions, rules and lazy algorithms have moderate performance except for the ZeroR which resulted in a negative correlation. Future work may consider the use of WEKA in other similar prediction analyses.

**KEYWORDS:** *Gaussian Processes, KStar, Linear regression, M5Rules, SMOReg, ZeroR*

### 1 INTRODUCTION

The Waikato Environment for Knowledge Analysis (WEKA) is a set of machine learning tools that Waikato University, located in New Zealand, has developed. It is free, open-source software that was created in Java under the terms of the GNU Public License. Linux, Mac, and Windows all support running it. It is made up of a group of machine learning algorithms used to carry out data mining activities. An environment for comparing learning approaches is provided by the graphical user interface tool, which is primarily utilized for a complete range of preprocessing tools and evaluation procedures (Vijayakamal and Narendhar, 2012).

In addition, Weka can be used to process data from agricultural fields. WEKA is a cutting-edge facility for the development and implementation of machine learning techniques for real-world data mining problems. The algorithms are directly applied to a dataset. WEKA is capable of performing a wide range of standard data mining activities which aims to pre-process, classify,

cluster, regression, visualize, and select features from the given datasets. (Buccola and Zanden, 1997).

The interest of the research is to show that the WEKA is useful in predicting and analyzing data. The red wine data from the UCI data set was used. Red wine was common to all people because it is used for occasions and gatherings. It is so interesting that red wine production is a growing market today (Smith and Margolskee, 2001).

In as much as the taste and quality of wine is a concern, the physicochemical and sensory test characteristics of the wine are to be considered. It is important to note that taste is difficult to understand when it comes to sensing it (Legin, Rudnitskaya, Lvova, Vlasov, Di Natale, and D'Amico, 2003).

This characteristic of the wine is still a challenge to wine producers on how they can accurately produce a good quality wine (Khalafyan, Temerdashev Akin'shina, and Yakuba, 2021)

The determination of wine's quality is not an easy task and needs experimentation. Experimentation includes a collection of data from the expert and the physicochemical laboratory test results of each of the wines. This process of estimation of wine items, securing and guaranteeing the quality of wines, preventing wine corruption, and controlling refreshment preparation needs financial consideration (WEKA Manual).

With so much data to be considered, there is a need to employ an intelligent machine to process machine learning functions and generate relevant data to produce patterns. With this development, WEKA is used. It has an easy-to-use graphical interface that allows for fast setup and actions. WEKA requires that the user data is in the form of a flat file or relation, which means that each data item is represented by a fixed number of characteristics, which are typical of a specified kind, alpha-numeric or numeric values. With simple options and visual interfaces, the WEKA application provides novice users with a tool for identifying hidden information from databases and file systems (Kumar, S, Agrawal K and Mandan, N, 2020).

Even though some researchers have employed machine learning approaches to evaluate wine quality,

there is still a lot of potential for improvement. An attempt to predict wine quality such as Support Vector Machine, Nave Bayes and Random Forest are used. The results are compared between the training and testing sets, with the best of the three strategies projected based on the training set results. Better results can be obtained if the best features from other techniques are extracted and combined to improve accuracy and efficiency. In the experiment which uses RStudio software, the Support Vector Machine resulted in the highest performance, with an accuracy of 67.25 percent when analyzing the red wine quality prediction, followed by the Random Forest, which has an accuracy of 65.83 percent, and the Nave Bayes algorithm, which has an accuracy of 55.91 percent (Sun, Danzer and Thiel, 1997).

Based on neural networks fed with 15 input variables, the result of (Vlassides, Ferrier, and Block, 2020). predicted six geographic wine origins. For their experiments in Germany, they utilized 170 data samples. They had a perfect prediction rate of 100 percent. A similar study by (Zaverri, and Joshi, 2017) used a neural network to classify Californian wine. Wine classification is based on grape maturity level and chemical analysis.

The study by (Jambhulkar and Baporikar, 2015) combined wireless sensor networks (WSN) and several techniques to forecast and predict heart disease from the important attributes extracted from the Cleveland dataset.

On the other hand, the study by (Beltran, Duarte-Mermoud, Soto Vicencio, Salah, and Bustos, 2008) compared several wine classification datasets. Linear discriminate analysis, support vector machine, and neural network (NN), were used to categorize Chilean wine. Three distinct types of Chilean wine were subjected to testing and analysis. The study (Gupta, 2018) used 147 bottles of rice wine to classify them into three wine groups. The abovementioned literature has been using data mining techniques to come up with a solution to data mining problems.

The objective of this study is to evaluate several regression algorithms. The UCI machine learning repository dataset of red wine quality with 11 physicochemical characteristics as inputs and 1 sensory

data as output is used. The regression algorithms were implemented in WEKA. The inputs are volatile acidity, fixed acidity, citric acid, residual sugar, chlorides, total sulfur dioxide, free sulfur dioxide, density, sulfates, pH and alcohol are the input variables based on physicochemical tests, and quality is the output variable based on sensory data (Kumar, S, Agrawal K and Mandan, N, 2020).

Furthermore, performance comparisons of the supervised algorithms are also presented to determine the best-supervised learning used in this study.

## 2 MATERIALS AND METHODS

### 2.1 Data Preparation

The data was taken from the UCI machine learning repository. The red wine data dataset contains 4898 occurrences with 11 inputs and 1 output variable. These 12 characteristics are taken into account when predicting the quality of red wine. With the help of sensory tests, the quality of red wine was determined by at least three sommeliers. The input and output features of red wine data are summarized in Table 1 (UCI Machine Learning Repository)

### 2.3. Gaussian Processes

This carries out non-linear regression using the Bayesian Gaussian process method. Users can select the kernel function as well as a regularization parameter called "noise" to regulate how well the model fits the data. Before learning the regression, they can opt to have the training data normalized or standardized. This technique is equivalent to kernel ridge regression for point estimations.

### 2.4 Linear Regression

Using either greedy backward elimination or by creating a full model from all attributes and removing terms one at a time in decreasing order of their standardized coefficients until a stopping criterion is met, Linear Regression performs least-squares multiple linear

Table 1. Summary of Inputs and Output Attribute for the Red Wine Data

Attribute	Data	Range	Description
Fixed acidity	Numeric	4.60 - 15.90	Input
Volatile acidity	Numeric	0.12 - 1.58	Input
Citric acid	Numeric	0.00 - 1.00	Input
Residual sugar	Numeric	0.90- 15.50	Input
Chlorides	Numeric	0.012 - 0.61	Input
Free sulfur dioxide	Numeric	1 -72	Input
Total sulfur dioxide	Numeric	6 - 29	Input
Density	Numeric	0.99 - 1.00	Input
pH	Numeric	2.74 - 4.01	Input
Sulfates	Numeric	0.33 - 2.00	Input
Alcohol	Numeric	8.4 - 14.9	Input
Quality	Numeric	3 - 8	Output

Table 2. Several Algorithms Tested for the Red Wine Quality Data

Algorithm	Type	Description
Gaussian Processes	Function	Implements Gaussian processes for regression without hyper parameter-tuning
Linear Regression	Function	Class for predicting using regression
Multilayer Perceptron	Function	Classifier that uses back propagation
Simple Linear Regression	Function	Learns a simple linear regression model
SMOreg	Function	Implements the support vector machine for regression
IBk	Lazy	K Nearest Neighbor Classifier
KStar	Lazy	Instance based classifier
LWL	Lazy	Locally Weighted Learning
DecisionTable	Rules	Class for building and using a simple decision table majority classifier
M5Rules	Rules	Decision using separate and conquer
ZeroR	Rules	Class for building and using 0-R classifier

regression with attribute selection. An adaptation of the AIC termination criterion is used by both techniques. It

is possible to disable attribute selection. Two additional improvements have been made to the implementation: a collinear attribute detection heuristic mechanism (which can be disabled) and a ridge parameter that stabilizes degenerate circumstances and can lessen overfitting by punishing large coefficients. Ridge regression is technically carried out via Linear Regression.

### 2.5 Multilayer Perceptron

A neural network that uses backpropagation for training is called a multilayer perceptron. It differs from the other schemes despite being included under functions since it has a unique user interface.

### 2.6 Simple Linear Regression

Based on a single attribute, Simple Linear Regression builds a linear regression model; it selects the one that produces the minimum squared error. Non-numerical characteristics and missing values are both prohibited (Witten, I, Frank, I, Hall, 2011).

### 2.7 SMOreg

According to (Witten, Frank, and Hall, 2011), the SMOreg implements the sequential minimal optimization approach for learning a support vector regression model.

### 2.8 IBk

When it comes time to classify the data, lazy learners simply store the training examples. The k-nearest-neighbor classifier, which is implemented by IBk, is the most basic lazy learner. Finding the closest neighbors can be sped up using a variety of different search strategies (Witten, Frank, and Hall, 2011).

### 2.9 LWL

A common algorithm for locally weighted learning is

LWL. It creates a classifier from the weighted instances after assigning weights using an instance-based approach. (Witten, Frank, and Hall, 2011).

### 2.10 Decision Table

A decision table classifier is created by Decision Table. It uses cross-validation to assess feature subsets and performs best-first search assessment. Instead of using the decision table's global majority based on the same set of features, an option employs the nearest-neighbor approach to identify the class for each instance that is not covered by a decision table entry (Witten, Frank, and Hall, 2011).

### 2.11 M5Rules

Regression rules are derived by M5Rules from model trees created by M5. Ridor generates the default rule, searches for exceptions with the lowest error rates using incremental reduced-error pruning, selects the best exceptions for each exception, and then iterates (Witten, Frank, and Hall, 2011).

### 2.12 Evaluative Algorithms

The most common method for evaluating numeric forecasts is the mean-squared error. The square root is occasionally employed to make it the same size as the projected value. The mean-squared error is used in different mathematical techniques, such as linear regression, because it is the easiest statistic to adjust intuitively. Another alternative is to utilize mean absolute error. Regardless of sign, the total amount of individual errors is simply averaged. Outliers are amplified by mean-squared error, but not absolute error (cases where the prediction error is larger than the others). Regardless of their magnitude, all wrong sizes are treated the same. The relative squared error, on the other hand, is a whole separate idea. The mistake is expressed as a percentage of the outcome if a simple predictor had been used instead.

The correlation coefficient is utilized to determine the

statistical link between the variables to be considered. The correlation coefficient is ranging from -1 to 1 (Witten, Frank, and Hall, 2011).

Root mean square error

$$= \sqrt{\frac{(p_1 - a_1)^2 + \dots + (p_n - a_n)^2}{n}}$$

Relative – absolute error =  $\frac{|p_1 - a_1| + \dots + |p_n - a_n|}{|a_1 - \bar{a}| + \dots + |a_n - \bar{a}|}$

Correlation Coefficient =  $\frac{S_{PA}}{\sqrt{S_P S_A}}$

where:

$$S_{pA} = \frac{\sum(p_1 - \bar{p})(a_1 - \bar{a})}{n - 1}$$

$$S_p = \frac{\sum(p_1 - \bar{p})^2}{n - 1}$$

$$S_p = \frac{\sum(a_1 - \bar{a})^2}{n - 1}$$

### 3 RESULTS AND DISCUSSION

There is a total of 12 variables of red wine collections. The variable quality rating is the dependent on predicted variable, whereas the other 11 variables are independent.

owners (39.4 %), and mainly with a monthly income range of P20,000-30,000 a month (34.6% and 24.2%, respectively).

Data file type used was the Attribute-Relation File Format (ARFF) in this experiment. The descriptive statistics of the data set are shown in Table 3. Table 3 shows the summary of Performance Output of the Quality of Red Wine. The overall performance of the model is based on the correlation coefficient. The dispersion degree with the use of standard deviation was computed to show the variability of each of the performance measures. From the result of the experiment, it shows that the KStar (r = 0.6043) ranked the highest, followed by Gaussian Processes (r = 0.5908), followed by M5Rules (r = 0.5904), followed by Linear Regression (r = 0.5872) and followed by the SMOreg (r = 0.5623). The rest of the 11 classifiers has also moderate positive correlation aside from the ZeroR which tends to have a negative to no correlation at all.

#### 3.1 Correlation Between the Output and Individual Inputs

Figure 1 shows the data visualization of the fixed acidity and quality variable taken from WEKA data visualization output. An (r = -0.39) indicates that the two variables has an inverse moderate relationship.

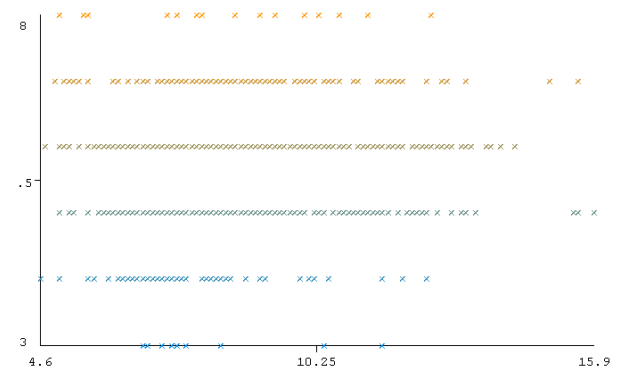


Figure 1. Data Visualization of Fixed Acidity (x) and Quality (y)

Figure 2 shows the data visualization of the volatility acidity and quality variable taken from the WEKA data visualization output. An (r = 0.23) indicates that the two variables has a positive weak relationship.

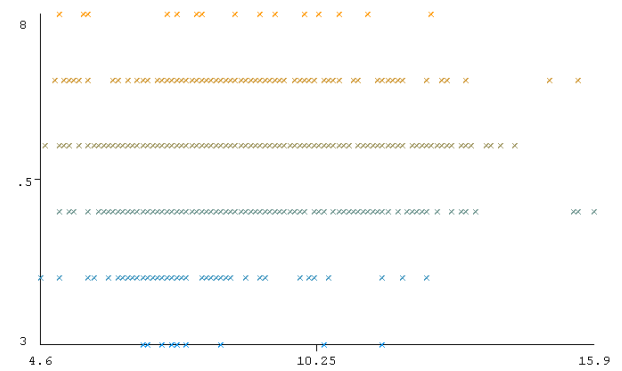


Figure 2. Data Visualization of Volatile Acidity (x) and Quality (y)

Figure 3 shows the data visualization of the citric acid and quality variable taken from the WEKA data visualization output. An (r = 0.01) indicates that the two variables has no relationship.

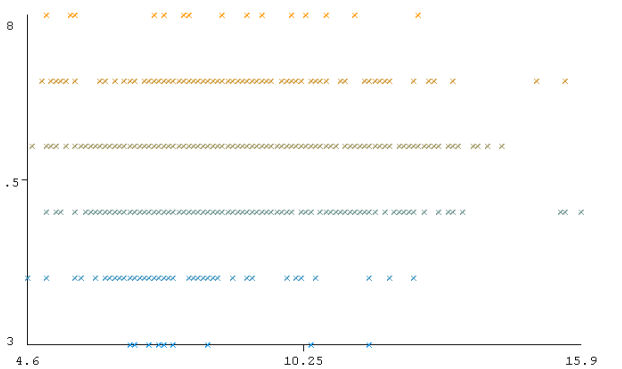


Figure 3. Data Visualization of citric acid (x) and Quality (y)

Table 3. Descriptive Statistics of the Dataset

Attributes Description	Range	Mean	Standard Deviation
Fixed acidity	4.60 - 15.90	8.32	0.53
Volatile acidity	0.12 - 1.58	0.53	0.18
Citric acid	0.00 - 1.00	0.27	0.19
Residual sugar	0.90- 15.50	2.54	1.41
Chlorides	0.012 - 0.61	0.09	0.05
Free sulfur dioxide	1 -72	15.87	10.46
Total sulfur dioxide	6 - 29	46.47	32.90
Density	0.99 - 1.00	1.00	0.002
pH	2.74 - 4.01	3.31	0.15
Sulfates	0.33 - 2.00	0.66	0.17
Alcohol	8.4 - 14.9	10.42	1.07
Quality	0 - 10	5.64	0.81

Table 4. Performance Measures of Data Set of Red Wine

Algorithm	Time taken to Build Model (Sec)	Total Instances	r	RMSE	RAE	Rank
Gaussian Processes	5.69	1599	0.5908	0.6514	73.99%	2
Linear Regression	0.02	1599	0.5872	0.6536	74.12%	4
Multilayer Perceptron	0.93	1599	0.5204	0.7131	80.38%	8
Simple Linear Regression	0.02	1599	0.4743	0.7107	82.31%	9
SMOreg	1.07	1599	0.5853	0.6574	72.94%	5
IBk	0	1599	0.5623	0.7440	60.10%	6
KStar	0	1599	0.6043	0.6773	57.38%	1
LWL	0	1599	0.4457	0.7228	82.65%	10
DecisionTable	0.16	1599	0.5271	0.6896	77.97%	7
M5Rules	0.27	1599	0.5904	0.6519	74.27%	3
ZeroR	0	1599	-0.0911	0.8081	100%	11

Figure 4 shows the data visualization of the Residual Sugar and quality variable taken from the WEKA data visualization output. An ( $r = 0.01$ ) indicates that the two variables has no relationship.

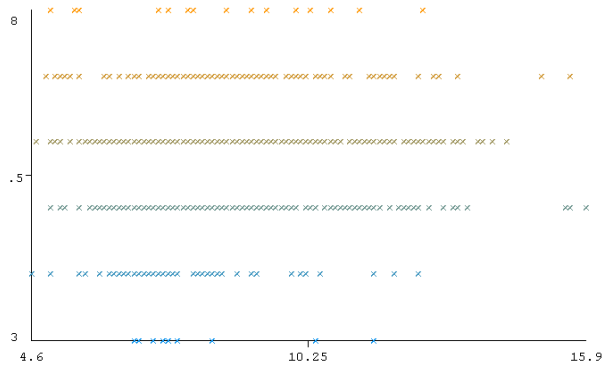


Figure 4. Data Visualization of Residual Sugar (x) and Quality (y)

Figure 5 shows the data visualization of the Chlorides and quality variable taken from the WEKA data visualization

output. An ( $r = -0.13$ ) indicates that the two variables has an inverse weak relationship.

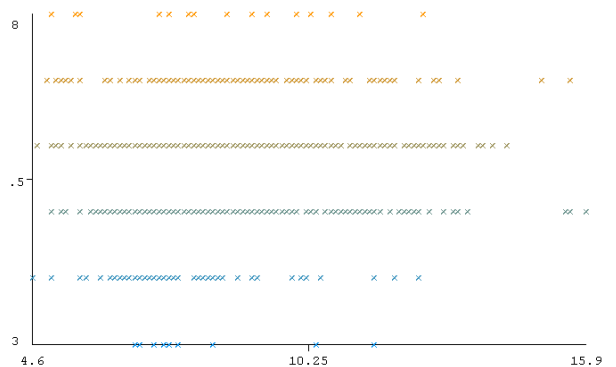


Figure 5. Data Visualization of Chlorides (x) and Quality (y)

Figure 6 shows the data visualization of the Free Sulfur Dioxide and quality variable taken from the WEKA data visualization output. An ( $r = -0.05$ ) indicates that the two variables has no relationship.

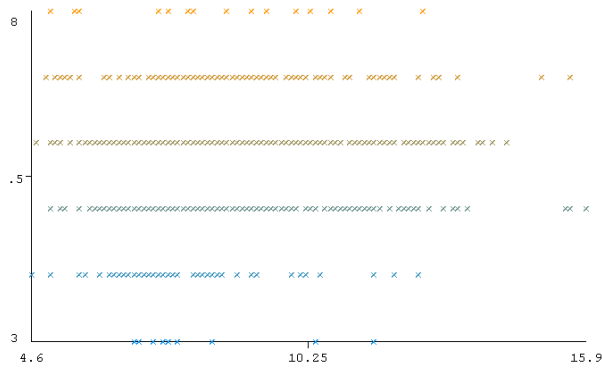


Figure 6. Data Visualization of Free Sulfur Dioxide (x) and Quality (y)

Figure 7 shows the data visualization of the total sulfur dioxide and quality variable taken from the WEKA data visualization output. An ( $r = -0.19$ ) indicates that quality and fixed acidity has an inverse weak relationship

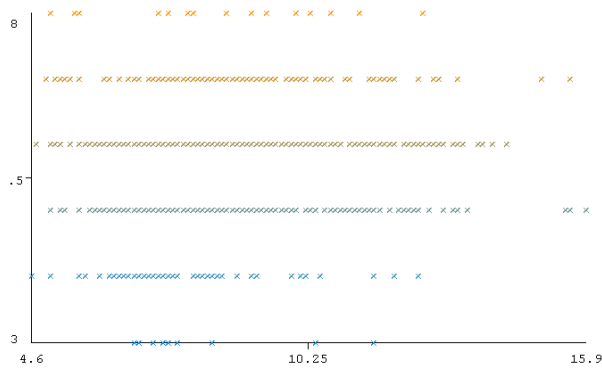


Figure 7. Data Visualization of Total Sulfur Dioxide (x) and Quality (y)

Figure 8 shows the data visualization of the density and quality variable taken from the WEKA data visualization output. An ( $r = -0.17$ ) indicates that the two variables has an inverse weak relationship.

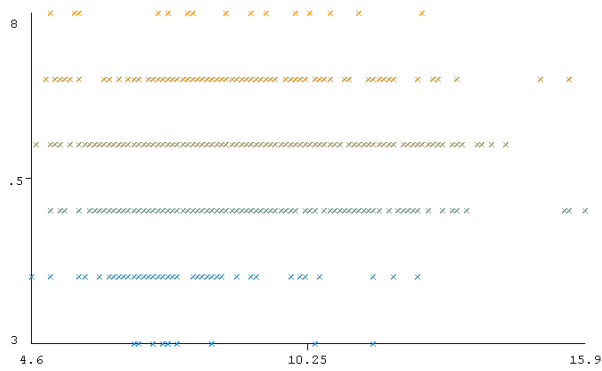


Figure 8. Data Visualization of Density (x) and Quality (y)

Figure 9 shows the data visualization of the pH and quality variable taken from the WEKA data visualization output. An ( $r = -0.06$ ) indicates that the two variables has no relationship.

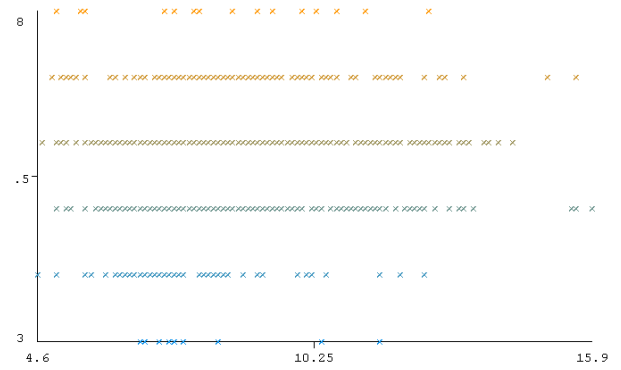


Figure 9. Data Visualization of pH (x) and Quality (y)

Figure 10 shows the data visualization of the sulphates and quality variable taken from the WEKA data visualization output. An ( $r = 0.25$ ) indicates that the two variables has a weak relationship.

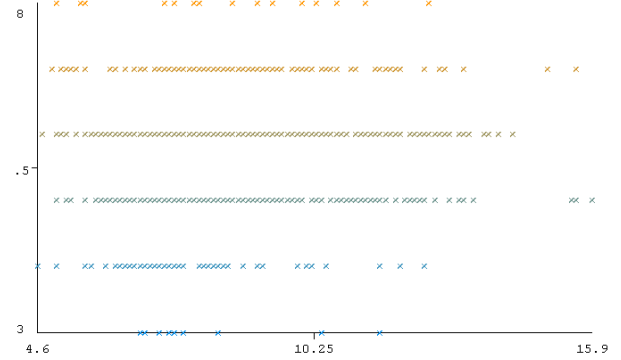


Figure 10. Data Visualization of Sulphates (x) and Quality (y)

Figure 11 shows the data visualization of the alcohol and quality variable taken from the WEKA data visualization output. An ( $r = 0.48$ ) indicates that quality and fixed acidity has a positive moderate relationship.

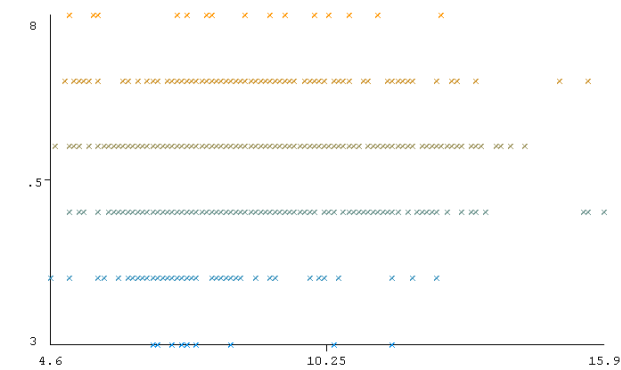


Figure 11. Data Visualization of alcohol (x) and Quality (y)

#### 4 CONCLUSIONS

Data mining is one of the available techniques that used data for investigation. The use of WEKA was able to process the data that was taken from (UCI Machine Learning Repository) that was processed to be able to predict a model and its performance. The study has already tested the utilization of machine learning techniques to predict a model and determine its accuracy to compare which among the existing supervised learning tools are the most accurate in producing a model. The (UCI Machine Learning Repository) red wine dataset was used for all experiments with the help of WEKA applications. The 1599 dataset of red wine samples was used. The red wine dataset consists of several physicochemical characteristics. 11 of which are the independent variables and one of which is the dependent variable. Independent variables are considered important factors affecting the dependent variable. The result shows that the value of the quality of the red wine can be predicted with a moderate positive correlation with the physicochemical characteristics of the wine. It shows that the KStar ( $r = 0.6043$ ) ranked the highest, followed by Gaussian Processes ( $r = 0.5908$ ), followed by M5Rules ( $r = 0.5904$ ), followed by Linear Regression ( $r = 0.5872$ ) and followed by the SMOreg ( $r = 0.5623$ ). The rest of the 11 classifiers has also moderate positive correlation aside from the ZeroR which tends to have a negative to no correlation at all. Future work may consider the use of WEKA in other datasets that will also use a similar prediction algorithm.

#### ETHICAL STATEMENT

Proper citation and credit were given to the owner of the data set and experiments are considered with prior citation and acknowledgement.

#### ACKNOWLEDGMENT

The author would like to thank the UCI Machine Learning Repository and the authors of the study entitled Modelling wine preferences by data mining from physicochemical properties for this paper.

#### REFERENCES

Barth, J., Katumullage, D., Yang, C., and Cao, J. (2020). Classification of Wines Using Principal Component Analysis. *Journal of Wine Economics*, Vol. 16, No. 1, pp. 56-

67. doi:10.1017/jwe.2020.35
- Beltran, N., Duarte-Mermoud M., Soto Vicencio, V., Salah, S., and Bustos, M. (2008). Chilean Wine Classification Using Volatile Organic Compounds Data Obtained With a Fast GC Analyzer," in *IEEE Transactions on Instrumentation and Measurement*, vol. 57, no. 11, pp. 2421-2436. doi: 10.1109/TIM.2008.925015.
- Bouckaert, R., Hall, F., Kirkby, R., Reutemann, P., Seewald, A., and Scuse, D. (2021). WEKA Manual, [https://statweb.stanford.edu/~lpekelis/13\\_datafest\\_cart/WekaManual-3-7-8.pdf](https://statweb.stanford.edu/~lpekelis/13_datafest_cart/WekaManual-3-7-8.pdf)
- Brownlee, J. (2016). How to use Regression Machine Learning Algorithms in Weka. Retrieved from <https://machinelearningmastery.com/use-regression-machine-learning-algorithms-weka/>
- Buccola, S.T and Zanden, L.V. (1997). Wine demand, price strategy, and tax policy, *Review of Agricultural Economics*, Vol. 19, No. 2, pp. 428-440, <https://www.jstor.org/stable/1349750>
- Cortez, P, Cerderia, A, Almeida, F., Matos, T., & Reis, J. (2009). Modelling wine preferences by data mining from physicochemical properties," In *Decision Support Systems*. Elsevier, Vol. 47, No. 4, pp. 547-553. ISSN: 0167-9236. 2009. <https://doi.org/10.1016/j.dss.2009.05.016>
- Data from Wine quality [dataset], UCI Machine Learning Repository [Online] (n.d). Available: <https://archive.ics.uci.edu/ml/datasets/Wine+Quality>, [Accessed: March 2021]
- Gupta, Y. (2018). Selection of important features and predicting wine quality using machine learning techniques. *Procedia Computer Science*, Vol. 125, pp 305-312, ISSN 1877-0509. <https://doi.org/10.1016/j.procs.2017.12.041>.
- Gutiérrez, J, Rubio, C, Moreno, I, González, A., Weller, D., Bencharki, N., Hardisson D., & Revert, C. (2017). Estimation of dietary intake and target hazard quotients for metals by consumption of wines from the Canary Islands," *Food and Chemical Toxicology*, Vol. 108, Part A, 2017, pp. 10-18, ISSN 0278-6915. <https://doi.org/10.1016/j.fct.2017.07.033>.
- Khalafyan, Z., Temerdashev, Z., Akin'shina, V., Yakuba, Y. (2021). Data on the sensory evaluation of the dry red and white wines quality obtained by traditional technologies from European and hybrid grape varieties in the Krasnodar Territory, Russia, *Data in Brief*, Vol 36, 2021, 106992, ISSN 2352-3409,

- <https://doi.org/10.1016/j.dib.2021.106992>.
- Kumar, S, Agrawal K and Mandan, N. (2020). White wine quality prediction using machine learning techniques. *2020 International Conference on Computer Communication and Informatics (ICCCI)*. pp. 1-6, doi: 10.1109/ICCCI48352.2020.9104095.
- Legin, A, Rudnitskaya, A, Lvova, L, Yu. Vlasov, Di Natale, C, D'Amico, A. (2003). Evaluation of Italian wine by the electronic tongue: recognition, quantitative analysis and correlation with human sensory perception. *Analytica Chimica Acta, Vol. 484*, Issue 1, pp. 33-44. ISSN 0003-2670, [https://doi.org/10.1016/S0003-2670\(03\)00301-5](https://doi.org/10.1016/S0003-2670(03)00301-5).
- Smith, D.V. and Margolskee, R.F. (2001). Making Sense of Taste. *Scientific American, Vol. 284*, pp. 32-39, <http://dx.doi.org/10.1038/scientificamerican0301-32>
- Sun, L, Danzer, K and Thiel, G. (1997). Classification of wine samples by means of artificial neural networks and discrimination analytical methods," *Fresenius J Anal Chem Vol. 359*, pp. 143-149. <https://doi.org/10.1007/s002160050551>
- Vijayakamal, M and Narendhar, M. (2012). A Novel Approach for WEKA & Study On Data Mining Tools. *International Journal of Engineering and Innovative Technology (IJEIT)*, Vol 2, Issue 2. [https://www.ijeit.com/vol%202/Issue%202/IJEIT1412201208\\_20.pdf](https://www.ijeit.com/vol%202/Issue%202/IJEIT1412201208_20.pdf)
- Vlassides, S, Ferrier J and Block, D. (2001). Using historical data for bioprocess optimization: Modeling wine characteristics using artificial neural networks and archived process information. *Biotechnol. Bioeng., Vol. 73*, pp. 55-68, 2001. [https://doi.org/10.1002/1097-0290\(20010405\)73:1<55::AID-BIT1036>3.0.CO;2-5](https://doi.org/10.1002/1097-0290(20010405)73:1<55::AID-BIT1036>3.0.CO;2-5)
- Witten F and Hall, M.(2011). *Credibility: Evaluating What's Been Learned. Data Mining: Practical Machine Learning Tools and Techniques. 3rd Edition, The Morgan Kaufmann series in data management systems*, Burlington, United States, pp 180-182, 2011
- Yu, H., Lin, H., Xu, H., Ying, Y., Li, B., and Pan, X. (2008). Prediction of Enological Parameters and Discrimination of Rice Wine Age Using Least-Squares Support Vector Machines and Near Infrared Spectroscopy. *Journal of Agricultural and Food Chemistry Vol. 56*, No. 2, pp. 307-313, <https://doi.org/10.1021/jf0725575>
- Zaveri, H and Joshi, N. (2017). Comparative Study of Data Analysis Techniques in the domain of medicative care for Disease Predication. *International Journal of Advanced Research in Computer Science. Vol. 8*, No. 3, pp. 564-566, 2017. <http://www.ijarcs.info/index.php/Ijarcs/article/view/3053/3036>