# Predicting Wine Quality by Statistical Machine Learning

**[Abstract]** *In the essay, a series of models are proposed to predict wine quality, providing support for wine expert and advice for customers. Efficient and comprehensive as possible, accurate result will help wine society. The essay first utilizes non-parametric models to provide an intuitive classification method with excellent results from cutting-edge techniques. Next the essay carries parametric part by 1) a polynomial model with important product item; 2) utilize natural cubic spline and smoothing spline to check every variable. The essay finally makes comparison between red wine and white wine characteristics, discovering significant differences between the two differing kinds.*
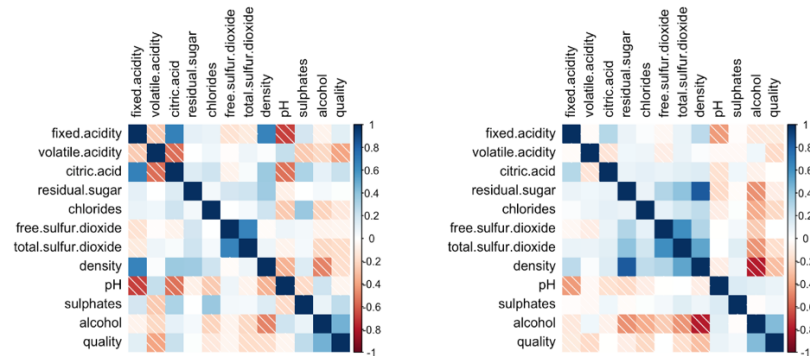
## 1. Data Exploration



*Figure 1 Red Wine Correlation Matrix (left), White Wine Correlation Matrix (right)*

The main target of the essay is to provide excellent model to predict the quality of wine and seek into the wine characteristics.

The essay has two different datasets, i.e. red wine and white wine (Cortez et al., 2009). As Figure 1, both correlation matrices have similar correlation tendency, but neither is too significant. Both datasets confront non-uniform distribution problem. Only 0.5% of data is Class 3 (the lowest quality), and nearly no wine data hits the highest class that only happens in white wine, as Figure 2. Thus, we divide the wine into Good wine and Bad wine, in which have Class 3, 4, 5 and Class 6, 7, 8, 9 respectively, by which the target value is now more uniform for the robustness of further modelling.
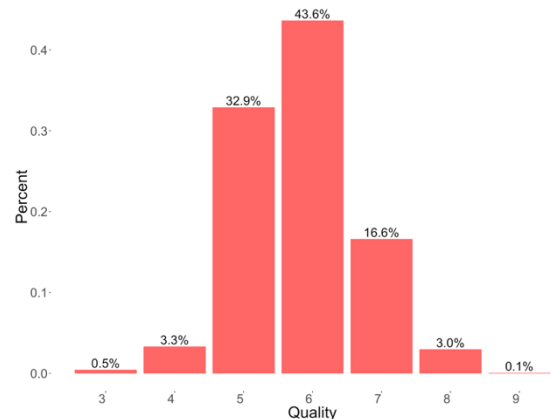


*Figure 2 Quality Distribution*

## 2. Model Solution

### 2.2 Parametric approach

#### 2.2.1 Polynomial Classification
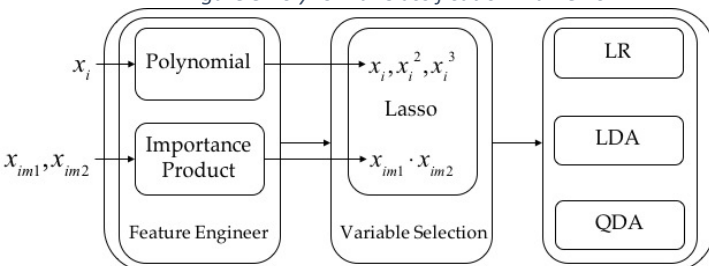


*Figure 3 Polynomial Classification Framework*

In polynomial models, the essay extends attributes by adding polynomial of degree up to three for each attribute (Keating & Cherry., 2011). The essay also picks the most two influencing attributes, making their product as a new predictor, by visualizing the importance of each attributes using random forest in previous work (Jaccard, 2003). For selection methods, best subset selection does not fit our massive data amount; therefore, the paper chooses Lasso to determine the best suitable model.

*Table 1 Polynomial Classification*

|  | LR | LDA | QDA |
| --- | --- | --- | --- |
| red wine | 0.77 | 0.77 | 0.78 |
| white wine | 0.78 | 0.79 | 0.78 |

*Table 2 Spline Classification*

|  | Natural Spline | Smoothing Spline |
| --- | --- | --- |
| red wine | 0.77 | 0.77 |
| white wine | 0.77 | 0.79 |

Most of the parametric classification approaches work well and have little differences between red and white as can be showed by Table 1. Since the paper add polynomial items and product variable into models, the distribution is not strict normal, so that it cannot expect an even better result by them.

### 2.2.2 Spline Classification

Spline classification goes from the original predictors. We set natural cubic spline transform and smoothing spline transform to all variables respectively to acquire different classification result. The results are showed in Table 2. Same as previous work, wine quality does have high degree correlations with predictors, and smoothing spline with penalty on twice differentiable functions concerns more on fluctuations and thus lead to better result with white wine, who holds more data to fit non-linear part.

### 2.1 Non-parametric approach

Parametric approach is not as good as non-parametric method, where the paper proposes seven models and get the accuracy results as following table shows.

<center>Table 3 Non-parametric Approach</center>

|  | Tree | Bagging | Random Forest | Boosting | SVM | Neural Network | **XGBoosting** |
|---|---|---|---|---|---|---|---|
| Red Wine | 0.76 | 0.84 | 0.85 | 0.81 | 0.81 | 0.78 | **0.86** |
| White Wine | 0.73 | 0.84 | 0.86 | 0.84 | 0.81 | 0.77 | **0.85** |

Both datasets can be explained well in these models, among which Random Forest and XGBoosting keep the relative advantage among all the models. Bagging and Boosting also show good ability to classify the wine, in which Boosting performing better may result from larger white wine set that strengthens boosting. Random Forest is only slightly better than Bagging, due to the fact that these two methods have almost the same intuition, but the former provides random selection in each split, reducing overfitting. Support vector machine performance is also fairly good under radial basis kernel, a similarity measure. It concerns interaction between observations.

Here the essay introduces two cutting-edge classifiers, i.e. Neural Network, XGBoosting. The essay hereby implements a 2-hidden-layer fully connected backpropagation neural network. In hidden layer, it has 5 and 3 neurons, respectively and uses sigmoid function as activation function, whose graph can be checked in appendix. Thus, the non-linear relationship among variables will be covered comprehensively (Cortez et al., 2009). Nonetheless, it performs only a little bit better than simple tree. It is normal when data is insufficient for neural network, where effects between response and predictors cannot be fully explained.

The XGBoosting (Chen & Guestrin, 2016) has strong improvements towards Boosting, which boosts trees with subjectively chosen depth. Boosting is trying to gather the functions among all trees, which are acquired by fitting the residuals from previous tree. Intuitively, a tree should move only a small path from its predecessor, so the whole tree set can get the optimum appropriately at last. However, Subjective depth without complexity concern may prevent from getting global optimum. XGBoosting can, in contrast, decide the depth by minimizing a target function below. It concerns not only loss indicating how model fits training data, but also model complexity including the number of leaves and L2 term of residuals.

$$Obj = \sum_{i=1}^{n} l\left( y_i, \hat{y}_i \right) + \sum_{K=1}^{K} \left( \gamma T + \frac{1}{2} \lambda \|w\|^2 \right)$$

$l(y_i, \hat{y}_i)$ is the fitting loss function, $T$ is the number of leaves and $\|w\|^2$ indicates the L2 term of residuals, whereas $\gamma$ and $\lambda$ are corresponding penalty coefficients. From that, boosting depth is controlled with respect to both fitting and complexity that reduces overfitting as possible.

The essay tests other XGBoosting result, i.e. 1) precision: 0.85 for how accurate out of those predicted positive; 2) recall 0.87 for how accurate out of real positive data, and 3) F1: 0.86 for combining both of them. All of these show robustness and incredible prediction power of XGBoosting. The ROC curve is presented in Appendix, where it is easy to find the model is even approaching the top.

## 3. Comparison Analysis

The essay also provides analysis into wine characteristics. By empirical logistic regression comparison to all original variables (see below in Table 4), both red wine and white wine take emphasis on *alcohol*, *sulphates, volatile acidity* and *free.sulfur.dioxide*. The first influences white more negatively, while other three more positively affect red wine. The *total.sulfur.dioxide* takes an significant negative impact only on red, as well as *citric acid* and *chlorides*, with less significance, while *residual.sugar*, *density* and *pH* significantly influence white wine. Unlike *volatile.acidity*, *fixed.acidity* does not show any impact with both kinds.

| | Red Wine | White Wine |
|---|---|---|
| *alcohol* | **0.867** *** | **0.7429** *** |
| *sulphates* | **2.7951** *** | **1.797** ** |
| *volatile.acidity* | **-3.282** *** | **-6.459** *** |
| *free.sulfur.dioxide* | **0.022** ** | **0.010** *** |
| | | |
| *total.sulfur.dioxide* | **-0.017** *** | -0.001 |
| *citric.acid* | **-1.274** * | 0.116 |
| *chlorides* | **-3.916** * | 0.8852 |
| | | |
| *residual.sugar* | 0.055 | **0.170** *** |
| *density* | -50.932 | **-270.9** *** |
| *pH* | -0.381 | **1.09** * |
| | | |
| *fixed.acidity* | 0.136 | 0.036 |

*Table 1 Wine Comparison*

*** (99.9%) ** (99%) * (95%) ' (90%)

*Figure 4 Important Variables, Red Wine (Red Upper), White Wine (Green Lower)*

Obviously, *alcohol* is essential to all wine, but differing kinds of wine have differing the second concern. For red wine it is *sulphates*, but *volatile.acidity* for white wine. Red wine focuses on the *sulfur.dioxide*s, whereas *sugar* and *density* may be vital concern of white wine, which implies we care more how sweet white wine is.

We also find similar result XGBoosting importance detection. In Figure 4 for points, the higher represents the more contribution to the model. (Gain means the contribution of the feature by its splits). The righter indicates the number of observations related to this feature. And a larger point means more nodes splits this feature. The *alcohol* which is higher contribute mostly for red and white wine. In upper figure, *total.sulfur.dioxide* and *chlorides* are right to the end, whereas in lower figure, *volatile.acidity*, which is more important to white wine showed in the table, dominates others except alcohol. Other orders are also presented above. It is interesting to see although white wine view *alcohol* as a very important predictor, it is hardly to see the times it splits the data, as the point is so small. It can even show that this variable is rather essential, with higher efficiency.

## 4. Conclusion

The essay focuses on binary classification, which is due to the small amount of data. The data is split by threshold 6. That is quality Larger than and equal to 6 is regarded as "Good", "Bad" on the other side.
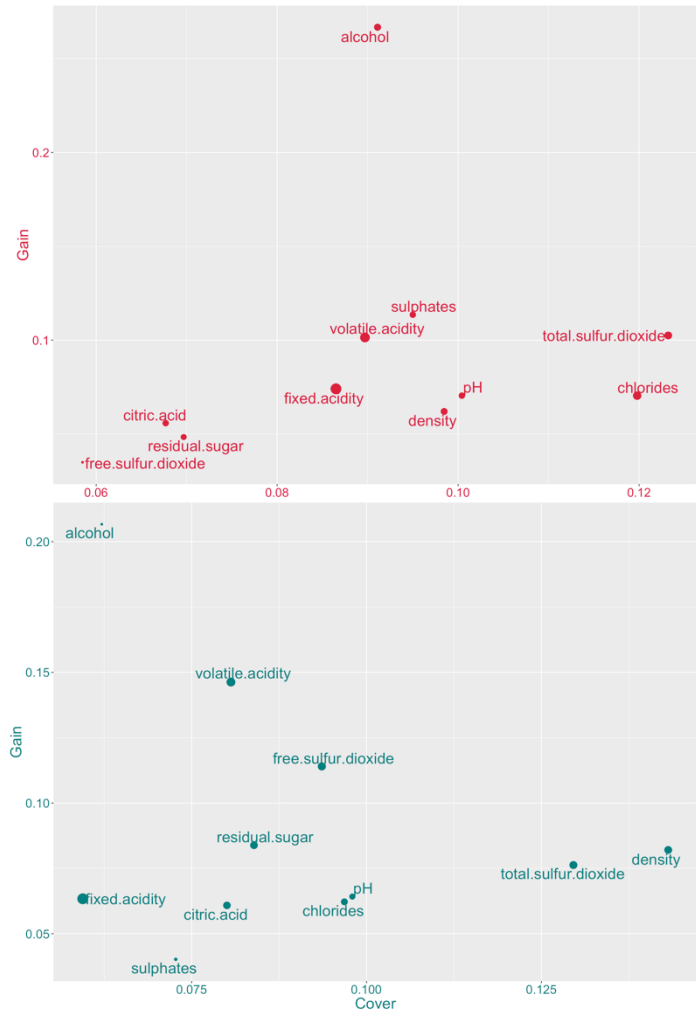
Non-parametric models with stronger techniques, i.e. XGBoosting overwhelm parametric models, which, in terms of accuracy particularly, is the most suitable model for both wines. On the other hand, it is easy to understand the result, especially for features importance comparison. Parametric models, on the contrary, does not provide as good result as the others. Nonetheless, coefficients of Logistic Regression, by which we compare issues of two differing wine, and describe characteristics of them, is still meaningful.

The essay hereby introduces XGBoosting as a supreme model for wine quality prediction. The essay recommends its use in wine recommendation or initially classification. Experts can refer to our result of wines or double check their subjective predictions.

As for wines quality, the essay believes degree of alcohol is vital to wine quality. The more sulphates will contribute to quality for both kinds, but mainly for red wine. The higher volatile acidity will lead to bad quality, but especially white wine. Free sulfur dioxide will enhance wines but the total should be controlled, red wine particularly. For better red wine, producer should also take care of citric acid and chlorides. And white wine tasters also care residual sugar, the density of wine and its pH value.

Due to the small amount of data, especially in the extreme cases. exceeding unbalance classes are severely influence model robustness if modeling multiclassification. It also prevents better result from Neural Network implemented before. Dealing with more data will be further wish to achieve even greater results. At the same time, wine recommendation towards individuals can be also interesting to explore.
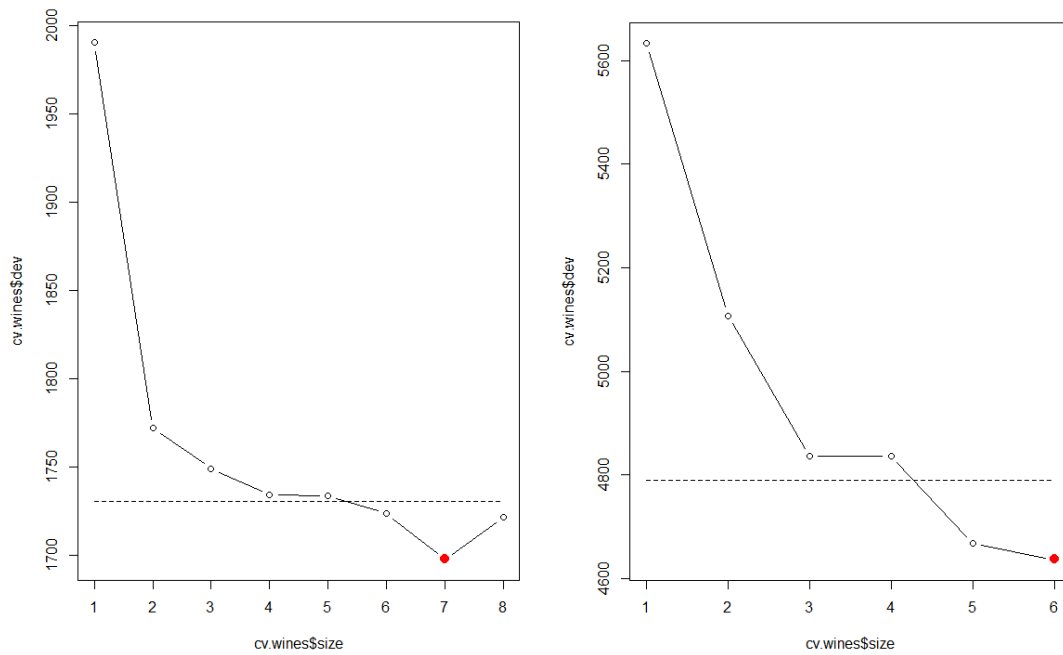
# [APPENDIX I: Visualization Graphs]



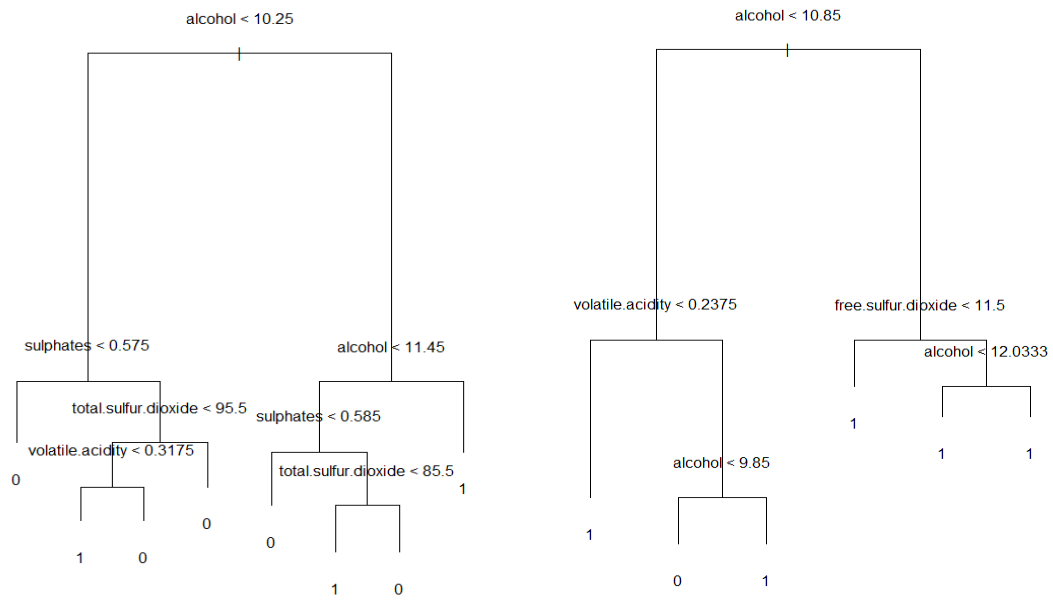Figure 1 Cross-validation results for tree size, Red Wine (left), White Wine (right)



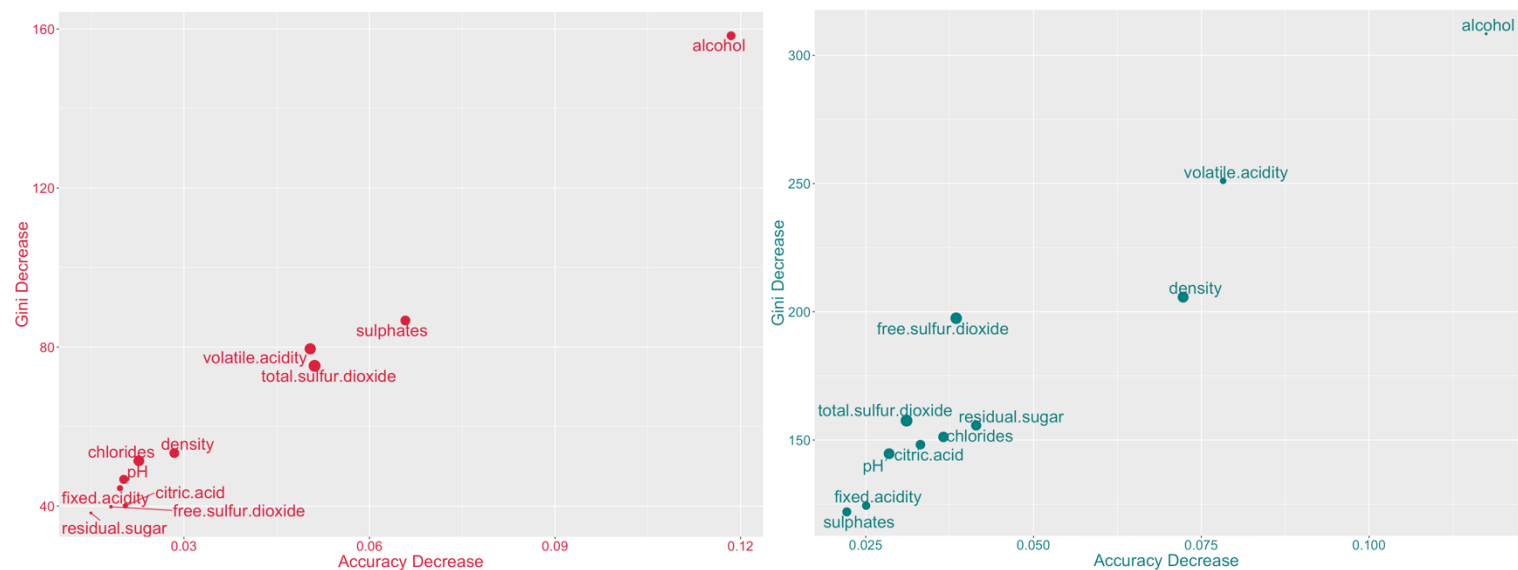Figure 2 Pruned tree diagram, Red Wine (left), White Wine (right)

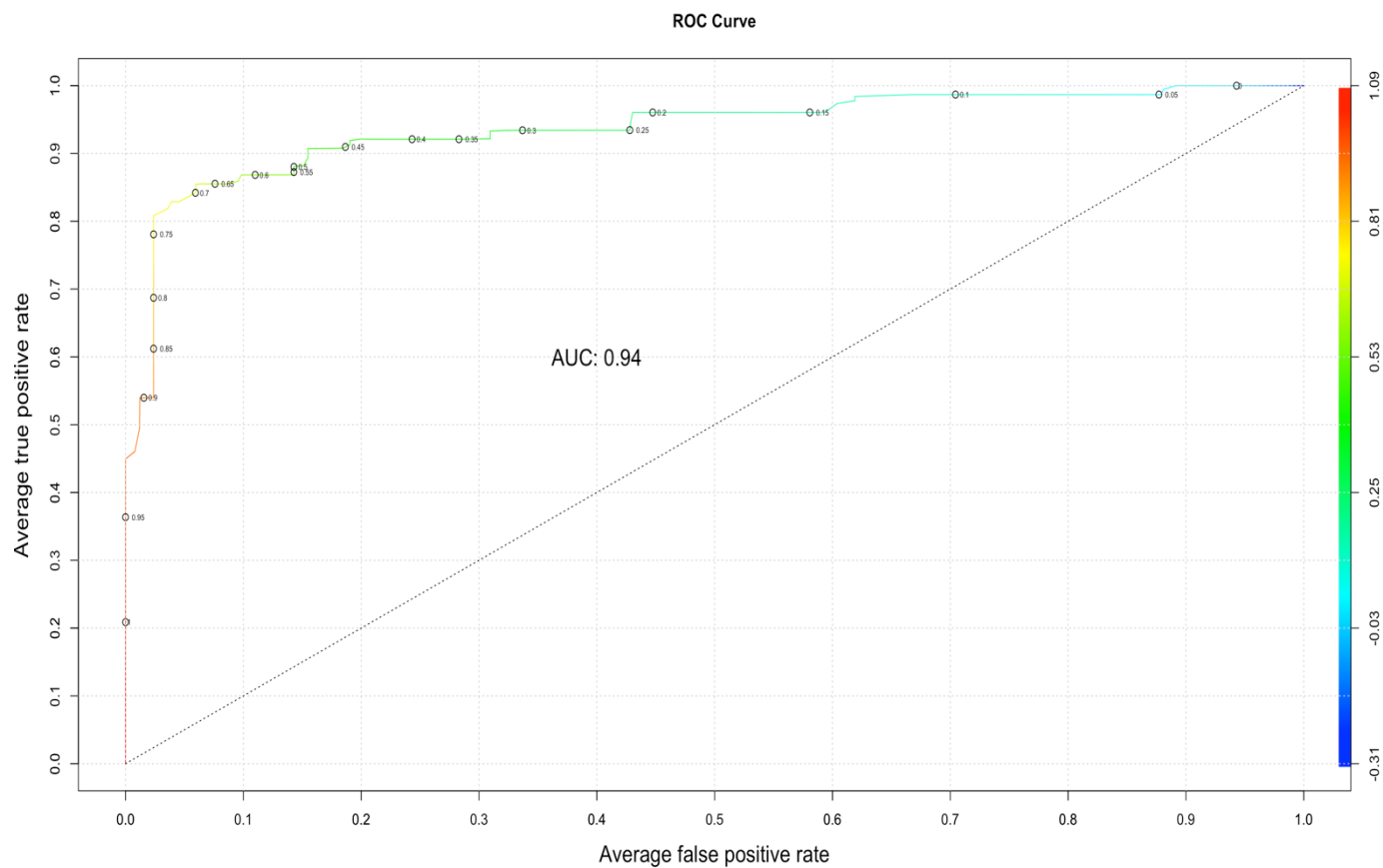*Figure 3 Important Variables by Random Forest, Red Wine (Red Left), White Wine (Green Right)*
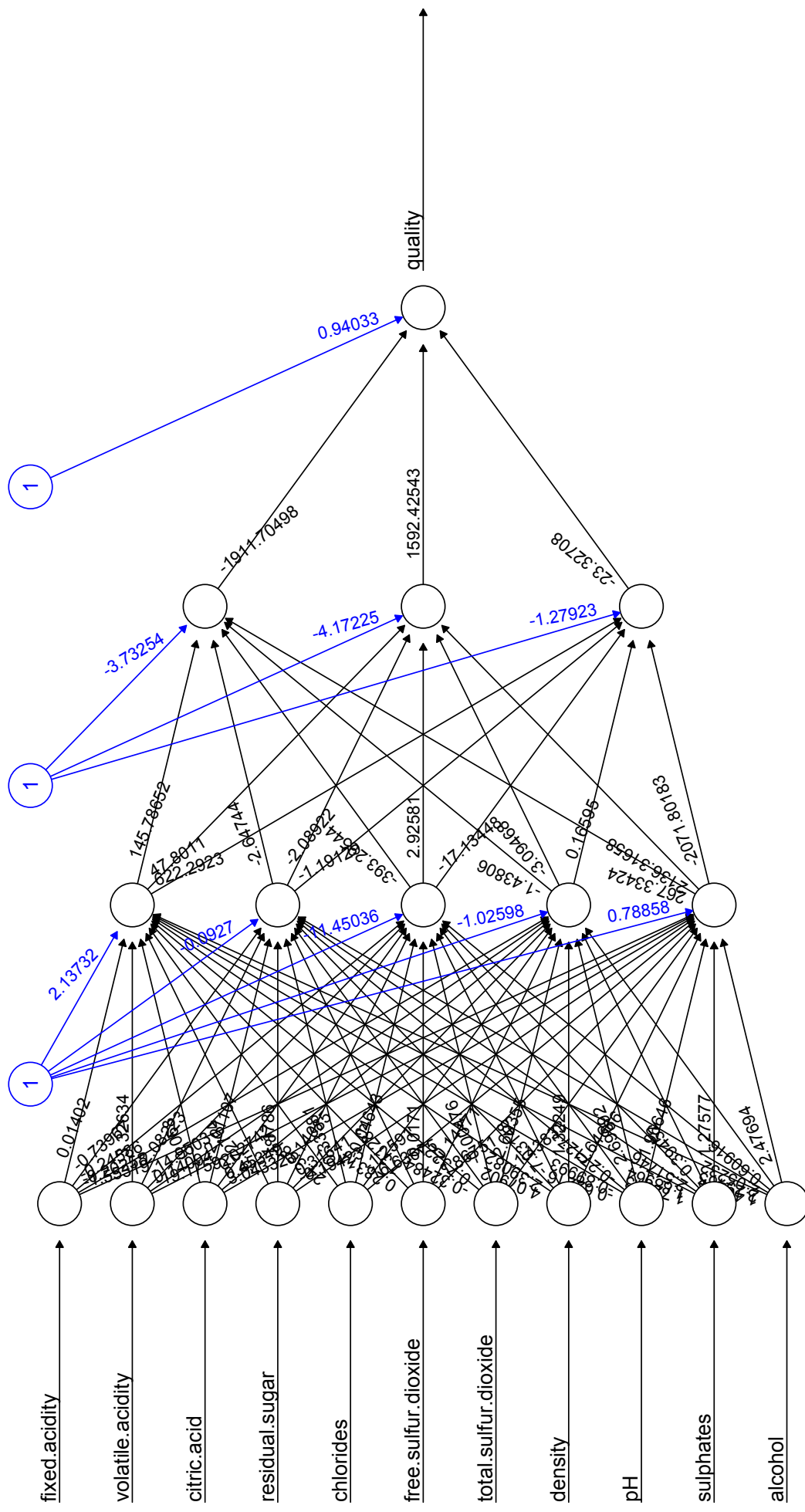


*Figure 4 ROC for XGBoosting*
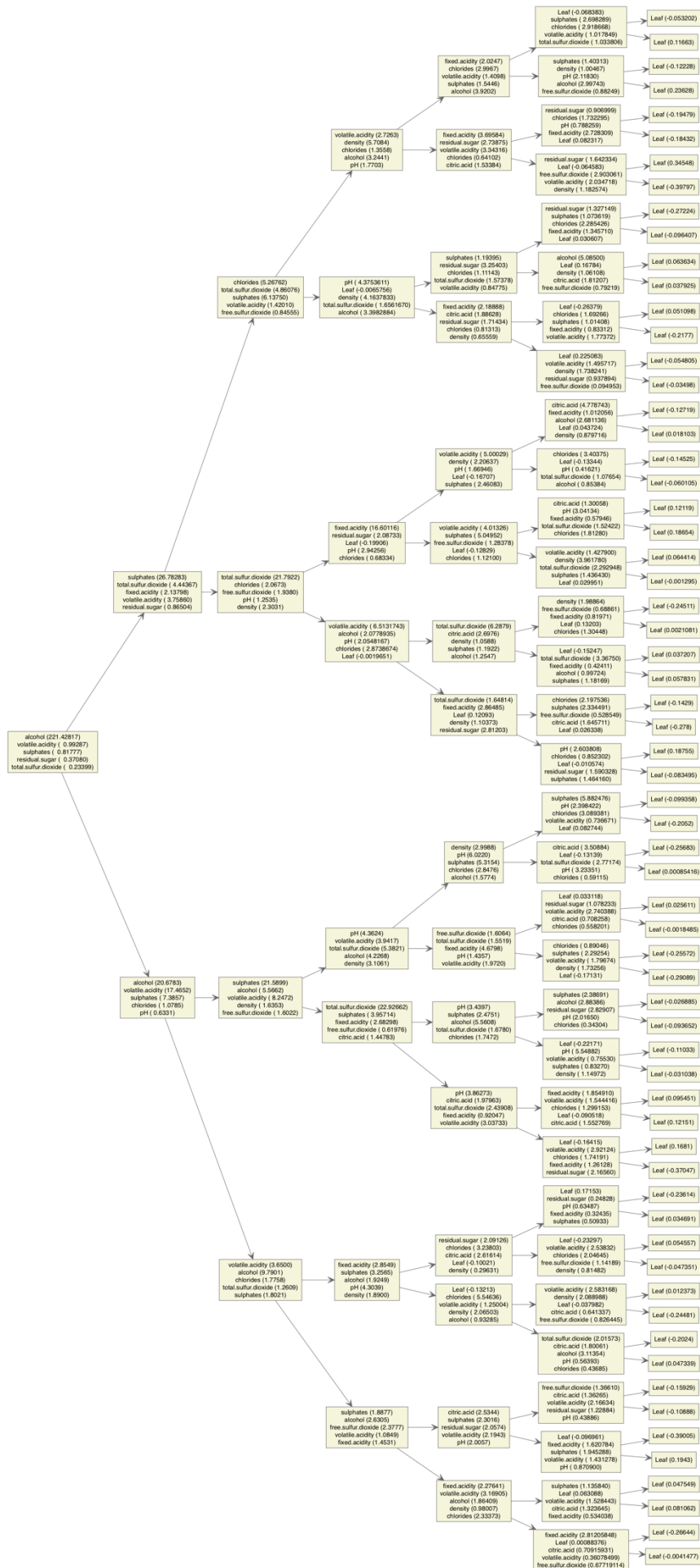
*Figure 5 Neural Network Model*

*Figure 6 XGBoosting Model*

# [APPENDIX II: Formula Presentation]

$$quality_{red} = \beta_0 + \beta_1 fixed.acidity + \beta_2 volatile.acidity + \beta_3 free.sulfur.dioxide + \beta_4 total.sulfur.dioxide$$
$$+ \beta_5 alcohol + \beta_6 residual.sugar^2 + \beta_7 chlorides^2 + \beta_8 pH^2 + \beta_9 alcohol \cdot sulphates + \varepsilon$$

$$quality_{white} = \beta_0 + \beta_1 volatile.acidity + \beta_2 residual.sugar + \beta_3 chlorides + \beta_4 free.sulfur.dioxide + \beta_5 pH$$
$$+ \beta_6 sulphates + \beta_7 alcohol + \beta_8 fixed.acidity^2 + \beta_9 total.sulfur.dioxide^2 + \beta_{10} pH^2 + \varepsilon$$

*Fomula 1, Polynomial Classification Logistic Regression Model*
*quality here resembles its sigmoid form*
*(Red Wine Upper, White Wine Lower)*

$$quality = \beta_0 + \beta_1 fixed.acidity + \beta_2 volatile.acidity + \beta_3 citric.acid + \beta_4 residual.sugar$$
$$+ \beta_5 chlorides + \beta_6 free.sulfur.dioxide + \beta_7 total.sulfur.dioxide + \beta_8 density$$
$$+ \beta_9 pH + \beta_{10} sulphates + \beta_{11} alcohol + \varepsilon$$

*Fomula 2, Empirical Research into Both Wine*
*quality here resembles its sigmoid form*

# [APPENDIX III: Reference]

Chen T, Guestrin C. XGBoost: A Scalable Tree Boosting System[C]. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 2016.

Cortez P, António Cerdeira, Almeida F, et al. Modeling wine preferences by data mining from physicochemical properties[J]. Decision Support Systems, 2009, 47(4):547-553.

Biau, Gérard. Analysis of a Random Forests Model[J]. Journal of Machine Learning Research, 2010, 13(2):1063-1095.

Cutler K, Breiman L. Random forests[J]. Machine Learning, 2004, 45(1):157-176.

Schapire R E. A Brief Introduction to Boosting[C]. Sixteenth International Joint Conference on Artificial Intelligence. Morgan Kaufmann Publishers Inc. 1999.

Zhu R L, Cao Y C, Pu Q M. Research on Information Technology with Wine Quality Evaluation Based on Neural Network[J]. Advanced Materials Research, 2014, 886:532-536.

James G, Witten D, Hastie T, et al. An Introduction to Statistical Learning[M]. Springer New York, 2013.

Keating K A, Cherry S. Use and Interpretation of Logistic Regression in Habitat-Selection[J]. Journal of Wildlife Management, 2011, 68(4):774-789.

Jaccard J, Turrisi R, Wan C K. Interaction effects in multiple regression[J]. Newbury Park California Sage Publications, 2003, 40(4):461.

# [APPENDIX IV: Reproducible codes]