

Machine Learning Methods to Predict Wine Quality

Joseph Despres

April 12, 2022

Contents

Introduction	1
Problem Description	1
Methodology	1
Results	3
Conclusions	3
References	5
Appendix	5

Introduction

Anyone looking to purchase a bottle of wine will be facing a number of options so large that only possible for a computer to search through. Wine varies significantly from region to region. Therefore, this is an ideal problem for machine learning to handle. In particular, this project aims to specify a function that maps chemical properties of various wines to the judges quality scores. This report begins with a discussion of the problem followed up by a description of the data. To model I use the sklearn package (Pedregosa et al., 2011) in the Python Programming Language (Van Rossum & Drake, 2009). After that, I devote several paragraphs to comparing different feature scaling strategies against logistic regression because there is a significant difference in runtime and performance. After that, I discuss the use of an elastic net regression, XGBoost, and Support Vector Machines in an ensemble to reach a prediction accuracy of . From there, I will go back into the dataset and attempt to improve the model's the accuracy through some more feature engineering.

Problem Description

This project aims to discover the relationship between wine quality and chemical properties. The motivation for this project is that one purchasing wine seriously must consider different regions, grapes, growing climate, seasons, storage time, and substantial variations in quality and price. This is a vast and complex search space with the potential to find something of significant value. This is a problem well suited to Machine Learning because of the search. The benefits of an accurate predictive wine model are that one could select excellent wine at a low price, lower the risk of purchasing low-quality wine, or even determine undervalued wine at auctions. I am not the first person to think of this, hence the published dataset.

Methodology

To determine the quality score based on chemical properties of a wine we will use a dataset hosted by the University of California Irvines Machine Learning repository (Lichman, 2013). These data are from a study conducted by the portugese government studying the ability to predict human wine taste (Cortez, Cerdeira, Almeida, Matos, & Reis, 2009). Using support vector machines, they achieve an out of sample accuracy of 0.33. Given this is in 2009, the purpose of this study is to employ a mixture of newer methods and attempt to outperform the published by Cortez. First, I do not have the domain knowledge to assess the validity of the data. I have no idea what is an unreasonable level of acidity, for instance. However, given that this is a published dataset, I assuming all the observations are actual observations.

This is a structured, curated, tabular dataset conveniently in First Normal Form. The outcome we wish to predict is the quality as rated by a professional wine judge. One thing that was not obvious from the UCI repository is the fact that there are duplicated rows, which means multiple judges rated the same wine the same quality. However, judges were anonymous in this study, therefore there is no feature indicating this. I will come back to this point after discussing model specifications. There are 11 additional features, which are listed in Table 1.

the target is wine quality as scored by a professional judge. Wine quality is an ordinal categorical variable. The class assignments are quite unbalanced. However, the nature of this problem is that very few wines will be of very high quality and many will be average. See Figure 2 for a plot of the class distributions. At this point, I will split the data and stash 30% of the data to reserve for testing and I will not touch those data until it is time to report a final result.

After some difficulties with fitting initial models and performance plateaus associated with the min max scaling method. I decided to test and document the results of using several different scaling strategies. I decided to run a 2,500 different partitions of training and testing data with a 50% split to investigate

Table 1: Descriptive Statistics

Feature	Mean	Std
Fixed Acidity	7.215	1.296
Volatile Acidity	0.340	0.165
Citric Acid	0.319	0.145
Residual Sugar	5.443	4.758
Chlorides	0.056	0.035
Free Sulfur Dioxide	30.525	17.749
Total Sulfur Dioxide	115.745	56.522
Density	0.995	0.003
Ph	3.219	0.161
Sulphates	0.531	0.149
Alcohol	10.492	1.193
Quality	5.818	0.873
Is Red	0.246	0.431

Note: 6497 observations

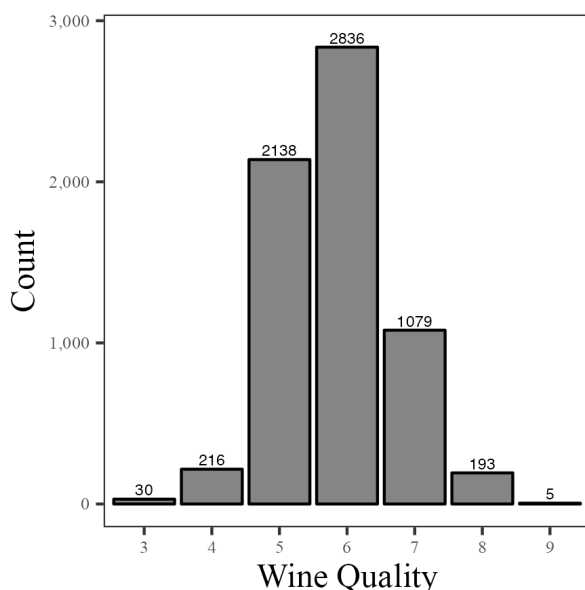


Figure 1: Comparing Resampled Prediction Accuracy by Transformation

the effects different scaling schemes. There was a significant difference in results as well as runtime.

This is due to the design matrices having difficulties inverting. The runtimes will differ based on packages, hardware and datasets so I don't list them.

This paper will consider two performance metrics. First accuracy measured as the percentage of correct predictions. Note, this is an unbalanced assignment problem, accuracy should not be compared to the number of classes rather the frequency of the most frequent class. The most frequent class is a quality score of 6. Which occurs in about 43% of the observations. This metric is intuitive, but does not capture our true intentions. I will also use a balance between precision and recall known as AUC. AUC is simply the area under the ROC curve.

Table 2: Demo table

Scaling Method	Mean	Std	Median
No Transformation	0.590	0.024	0.592
Dropped Outliers	0.645	0.026	0.641
Min Max Scale	0.690	0.018	0.691
Rescaled to Uniform	0.707	0.015	0.707
Standard Scale	0.726	0.017	0.725
Yao Johnson Transform	0.733	0.015	0.732

The Yao-Johnson outperforms the Standard scaling methods however, not by much and is a far more complicated transformation process. It is a Box-Cox transformation method modified to take values less than 0. The Box-Cox Method is for linearizing vectors by some exponent that will convert the variable to a linear scale. I do not see any indication that wine chemical properties are anything other than linear in practice. My reasoning for this is that one does not add an exponentially increasing or decreasing amount of sugar for instance. Although it has slightly better accuracy and lower standard deviation, both are desirable properties, The increase did not justify the additional complexity and a deviation from what I felt is reasonable. Therefore, I will use the standard scaler.

After studying the features and exploratory plotting, we are now ready to specify the models. The plan is to use an ensemble of models to predict

wine quality. The idea is that a collection of voting classifiers will be more accurate than a single one. The reasoning is that one model could be wrong in many ways however it would be correct in only one. Add a collection together and they will be collectively right more than wrong.

The first model is the multinomial logistic regression used during the feature scaling trials. This model is based on estimating the cumulative log odds of belonging to a given category based on the features. I add some regularization parameters L^1 and L^2 norms, commonly referred to as elastic net regression. A gridsearch over a five fold cross validation shows that an optimal $L^1 = 0.0315$ and $L^2 = 0.013$. This achieves a training set accuracy of 54.53% and an AUC of 0.775. These are not impressive results, so we will fit a Support Vector Machines.

A support vector machine model aims to find the maximum margin hyperplane between features that separate points. I use the radial basis function kernel and tune only the regularization parameter squared L^2 norm. After crossvalidating I find the best regularization parameter to be 6.58. This reaches an in-sample accuracy of 60% and an AUC of 0.801. A slight improvement.

After that I use XGBoost (Chen & Guestrin, 2016) an efficient implementation of gradient boosted trees. There are many parameters to search through, however I will only tune max tree depth, column sample by tree, and learning rate. After tuning, I found column sample to be 50% to be the best a learning rate of about 0.034 and a max depth of 13. These yield an astonishing in sample accuracy of 99.5% and an AUC of 0.998. This is a substantial improvement over the past two models. From here I will wrap all three in an ensemble. Test the accuracy. Given these matrices, it may be better to not even ensemble and just use XGBoost.

Then I save these models and wrap them into one ensemble model. Then fit that to the training data and compare different approaches. I want to reengineer some features and report some results.

Results

The XGBoost model wildly outperforms the linear I found it has an accuracy of the final testing set of. The ensemble. The model predictions are in Figure . This shows that the model could be used to

Conclusions

The main take-away of this study

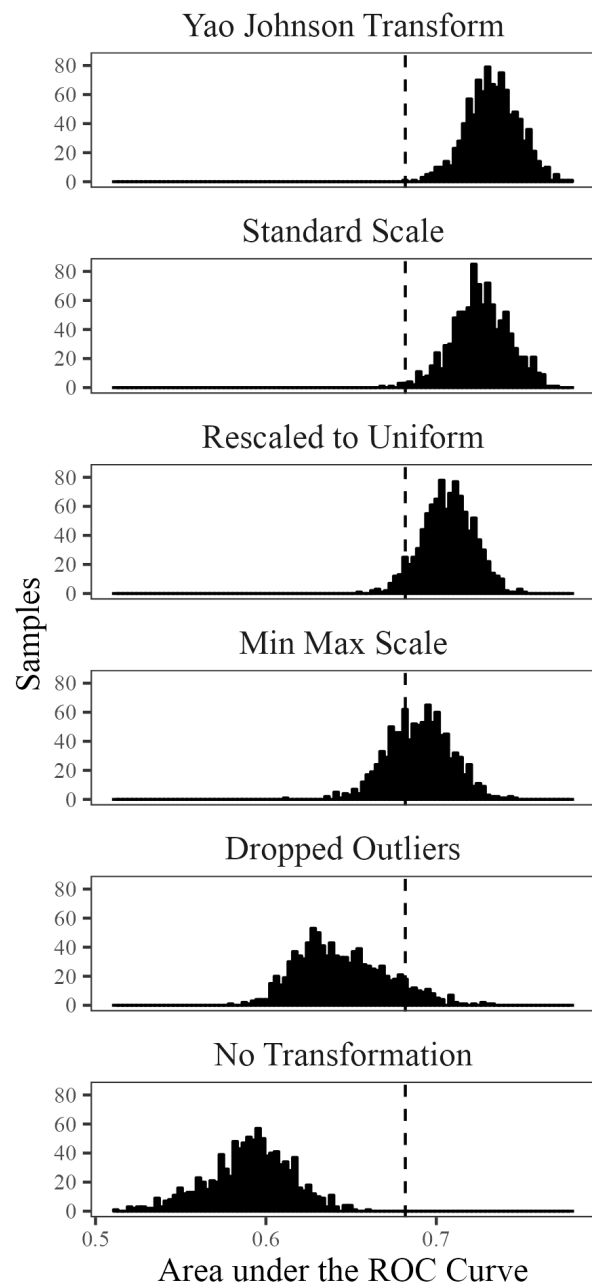


Figure 2: Comparing Resampled Prediction Accuracy by Transformation

References

- Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining* (pp. 785–794). New York, NY, USA: ACM. Retrieved from <http://doi.acm.org/10.1145/2939672.2939785> doi: 10.1145/2939672.2939785
- Cortez, P., Cerdeira, A., Almeida, F., Matos, T., & Reis, J. (2009). Modeling wine preferences by data mining from physicochemical properties. *Decision Support Systems*, 47(4), 547-553. Retrieved from <https://www.sciencedirect.com/science/article/pii/S0167923609001377> (Smart Business Networks: Concepts and Empirical Evidence) doi: <https://doi.org/10.1016/j.dss.2009.05.016>
- Lichman, M. (2013). *UCI machine learning repository*. Retrieved from <http://archive.ics.uci.edu/ml>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... others (2011). Scikit-learn: Machine learning in python. *Journal of machine learning research*, 12(Oct), 2825–2830.
- Van Rossum, G., & Drake, F. L. (2009). *Python 3 reference manual*. Scotts Valley, CA: CreateSpace.

Appendix

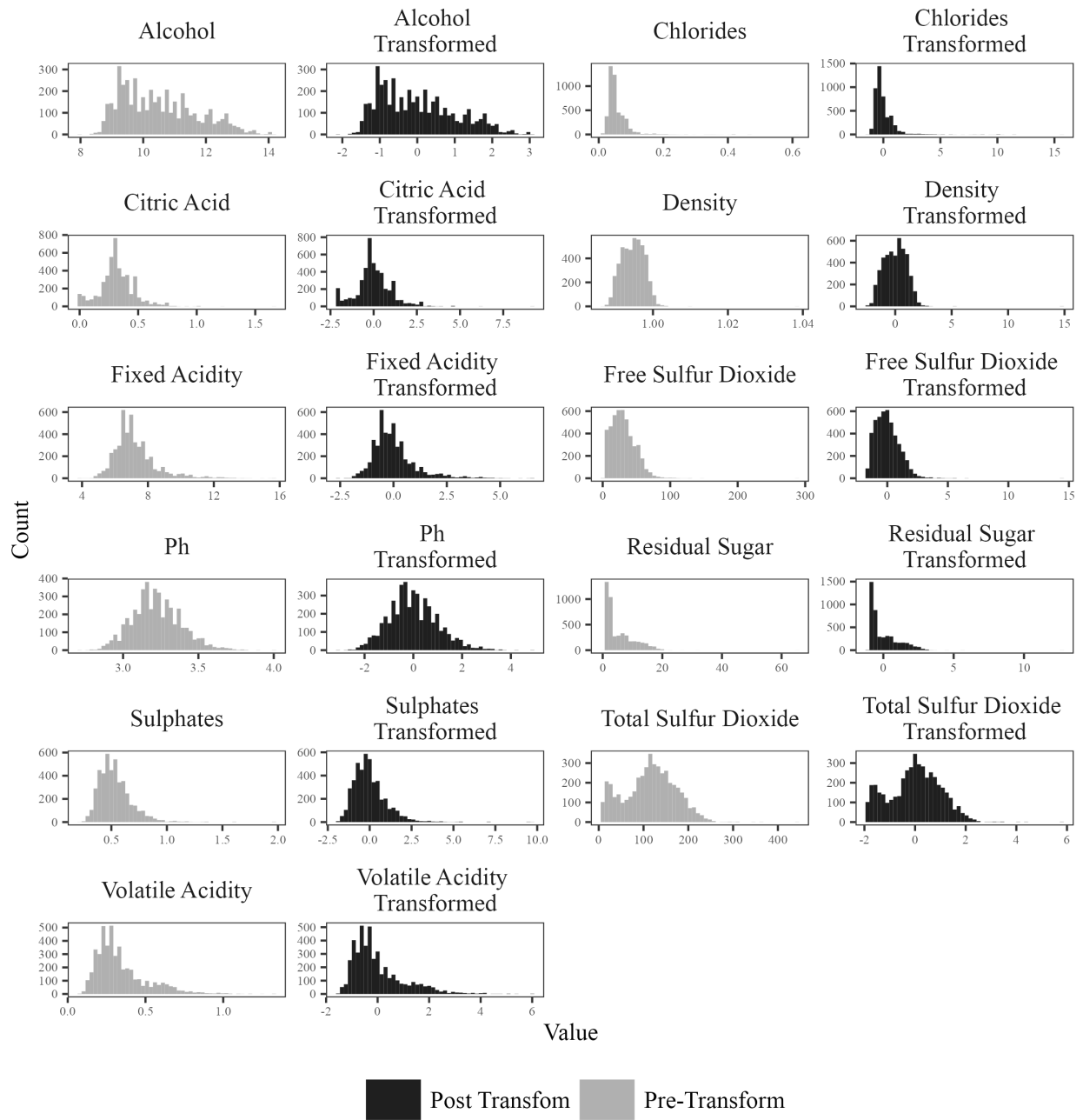


Figure 3: Comparing Resampled Prediction Accuracy by Transformation