

Final Report for Red Wine Analysis

Ruiqiang Chen, Michael DeWitt, David Williams, Alex Vannoy

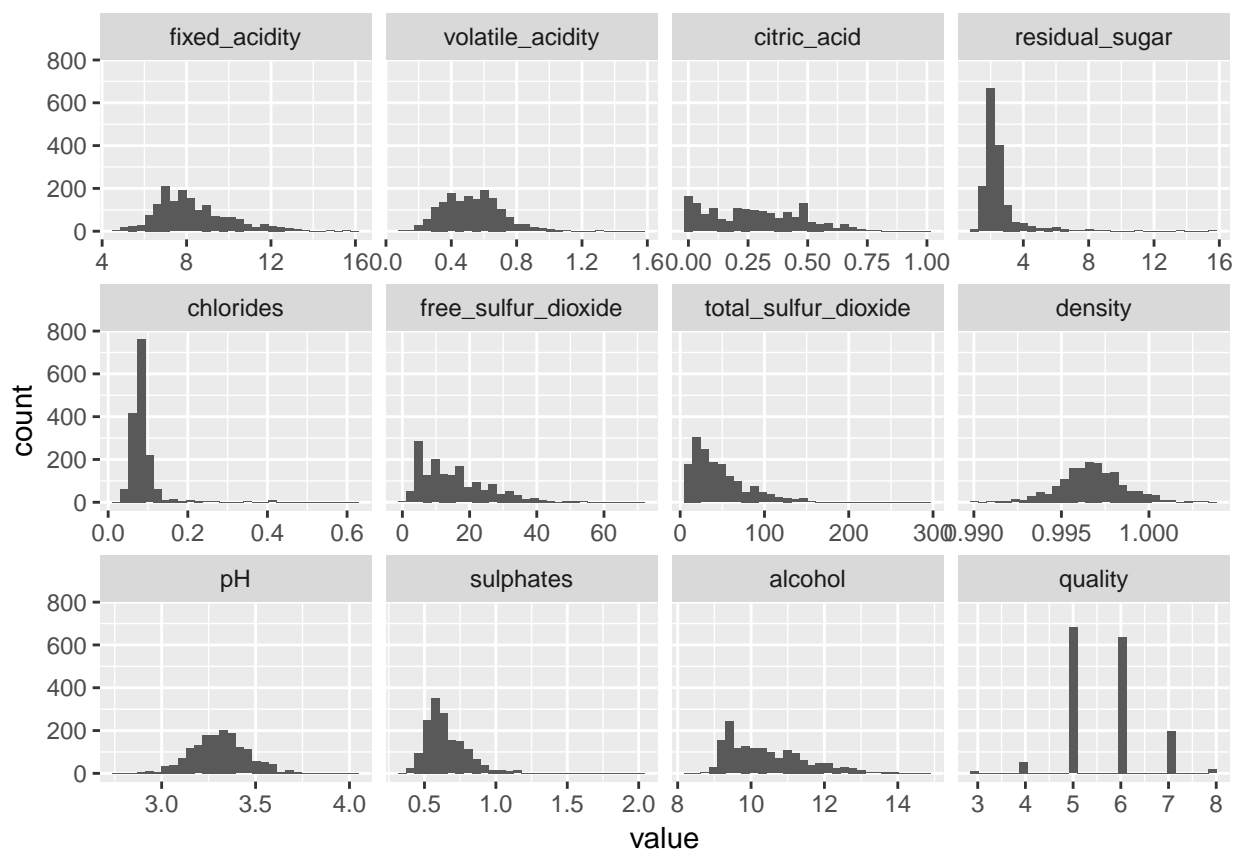
7/28/2017

Introduction

The purpose of this document is to report the proposed statistical models for classification of red wine bases on 11 predictors. The purpose of this analysis is to provide a model to the vinters in order for them to better predict the quality rating for their product. Analysis will be performing using both regression techniques and classification techniques.

Description of Data

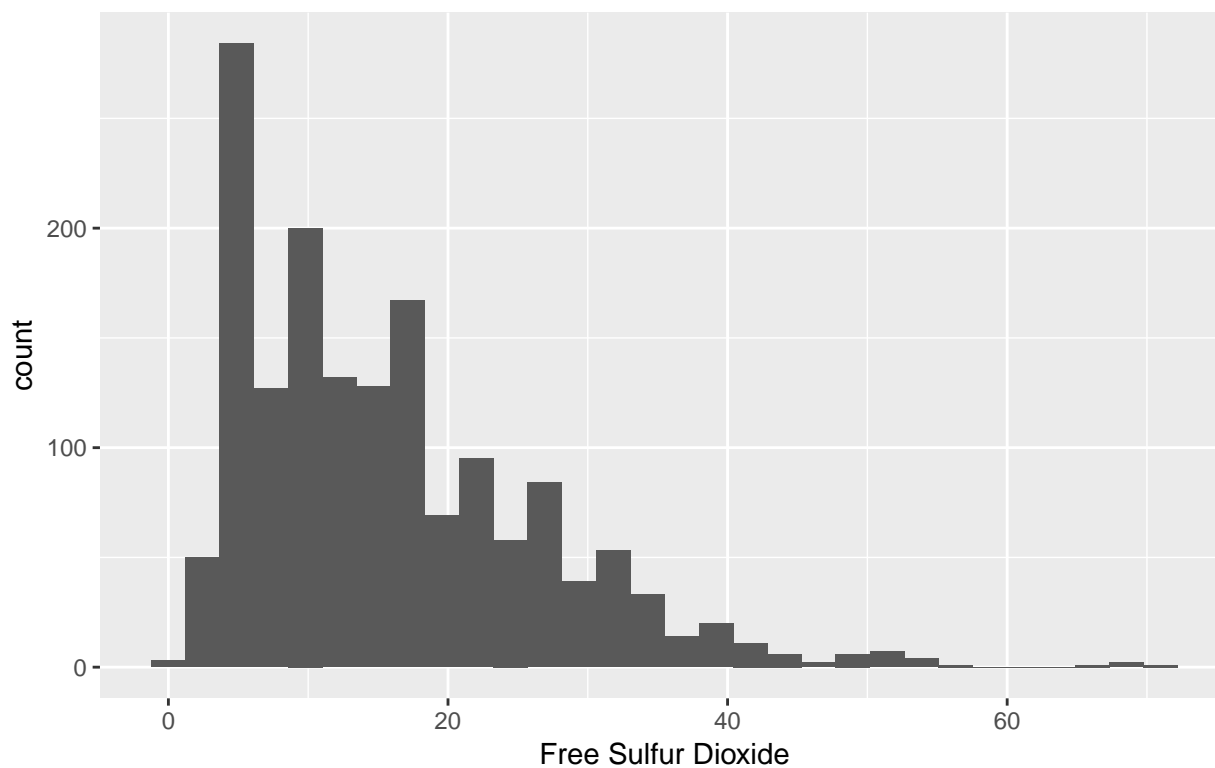
The data set provided is the Wine dataset from UC Irvine. It consists of 1599 with a total of 12 predictors. These predictors include the following fixed_acidity, volatile_acidity, citric_acid, residual_sugar, chlorides, free_sulfur_dioxide, total_sulfur_dioxide, density, pH, sulphates, alcohol, quality with the quality feature being associated with the judgement of the individual wine's quality. Quality is the feature of interest for the dataset as the vinter is interested in judging the wine's quality through objective means rather than todays subjective method of averaging the 1-10 point judgment of tastetesters. The distribution of these different criteria can be seen below:



The following are slightly right skewed: Fixed Acidity, Volatile Acidity, Citric Acid, , Free Sulphur Dioxide, Total Sulphur Dioxide, Sulphates, and Alcohol. Residual Sugar and Chlorides are heavily right skewed with density and pH appearing more normally distributed. Reviewing the individual components there appears to be a slight irregularity with total free sulfur dioxide. This can be seen in the histogram of this variable.

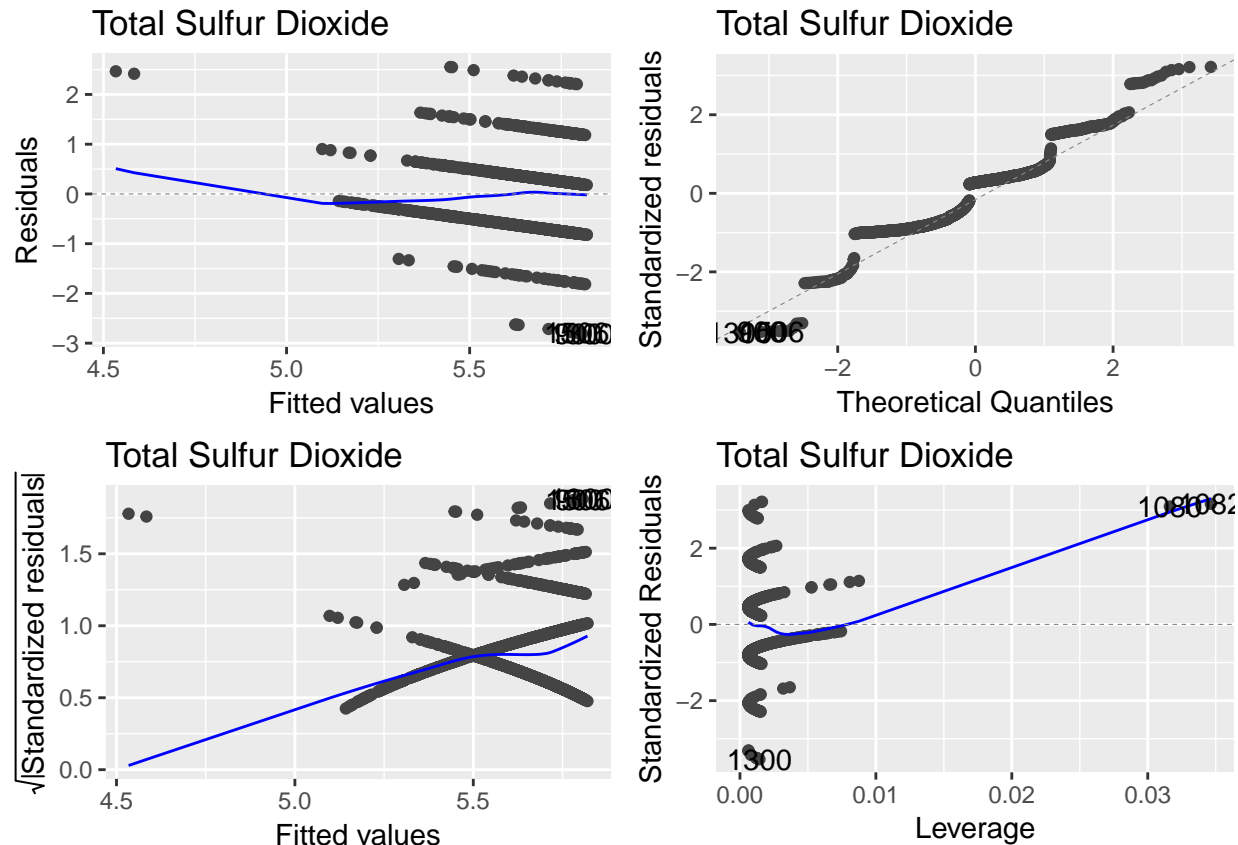
```
## Using classification as id variables
```

Histogram of Free Sulfur Dioxide



From UCI Wine Data Set

As well as the fit of thithat display high studentized residuals and leverage and thus should be considered for removal in the modeling process. These wines are 1080 and 1082.



These two wines have been removed from the clean dataset in order to be better predictors. The presence of these two wines may result in incorrect or inaccurate predictions. As we did not gather this dataset, we do not know if this information was incorrectly captured or if these values are real.

Method

In order to understand the testing error of any of the modeling used the data was divided in testing and training data sets with which to train then models and then test and estimate the testing error. Seventy percent of the raw data was randomly selected and placed in the training set with the remaining 30% used in the testing data set.

Regression

In order to select the best fit regression model several different modeling methods were tested. These include Least Squares Regression, Ridge Regression, Lasso Regression, Principle Components Regression and Partial Least Squares Regression. For each of these methods the quality integer was the value that the model was attempting to predict. The data was divided into two sets, a training set to train the model and a testing set for model validation. We will now go deeper in the model generation process for each of these different modeling types and methods.

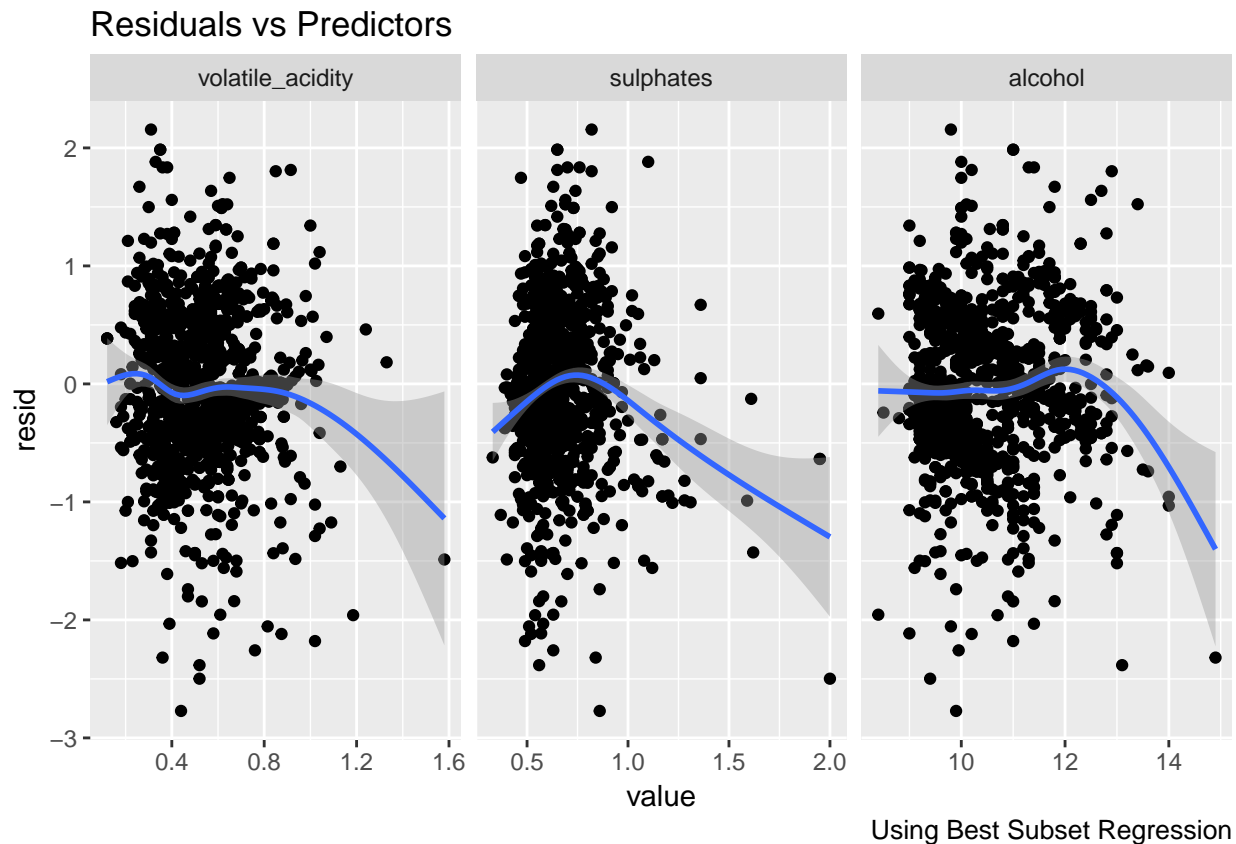
Least Squares

The least squares regression method that was tested was the best subset selection. The methodology used to determine the best subset model was to first run cross validation on the training set in order to determine the

number of predictors to include in the model. Once this analysis indicated that any added predictor after four variables were selected did not increase the accuracy of the model greatly using this cross validated method. The training data was then used to determine the best subset of the linear model with three predictors. The best subset included:

Residual Analysis

Here we need to make some plots against of the fit vs predictors and fit vs prediction to cross off that we considered our residuals



The residuals appear to have no distinct pattern which is a positive sign that there are not lurking relationships that have not been treated by the modeling.

Ridge Regression

Ridge regression was performed on the dataset as well. Cross validation was performed on the training data set to determine the optimum value for lambda for the ridge regression. This lambda, 0.079 was then using in a ridge regression model with the testing dataset.

Lasso Regression

Lasso regression was used with cross validation on the data set. Cross validation was used to determine the best lambda which was 0.005. As a function of the lasso regression only pH was shrunk to zero with total sulphur dioxide and free sulphur dioxide being spring to near zero.

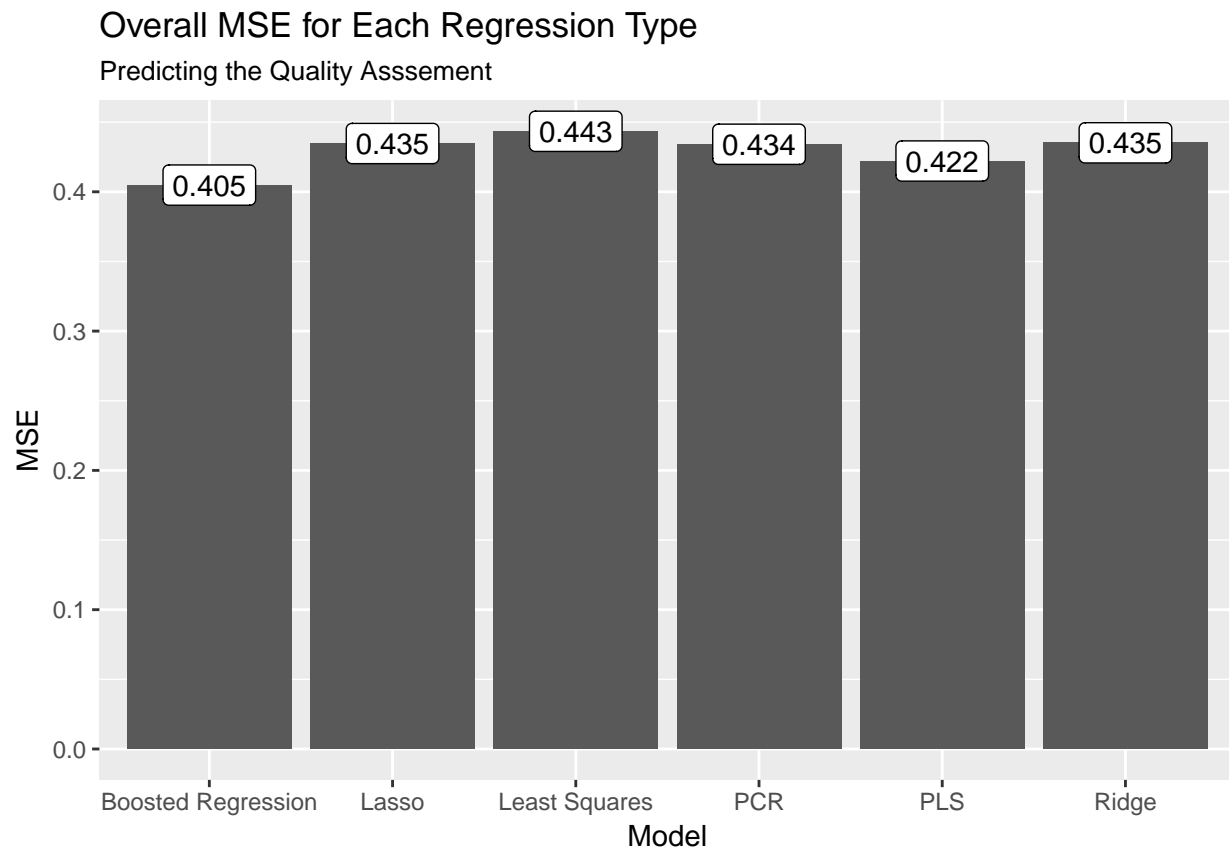
Principal Components Regression

Principal components regression was used. Based on the analysis of the principal components, the first nine principal components were used to be trained on the training set. This was done because 90% of the variation could be explained by these first nine components.

Partial Least Squares Regression

Boosted Regression

Model Selection



Residual Analysis

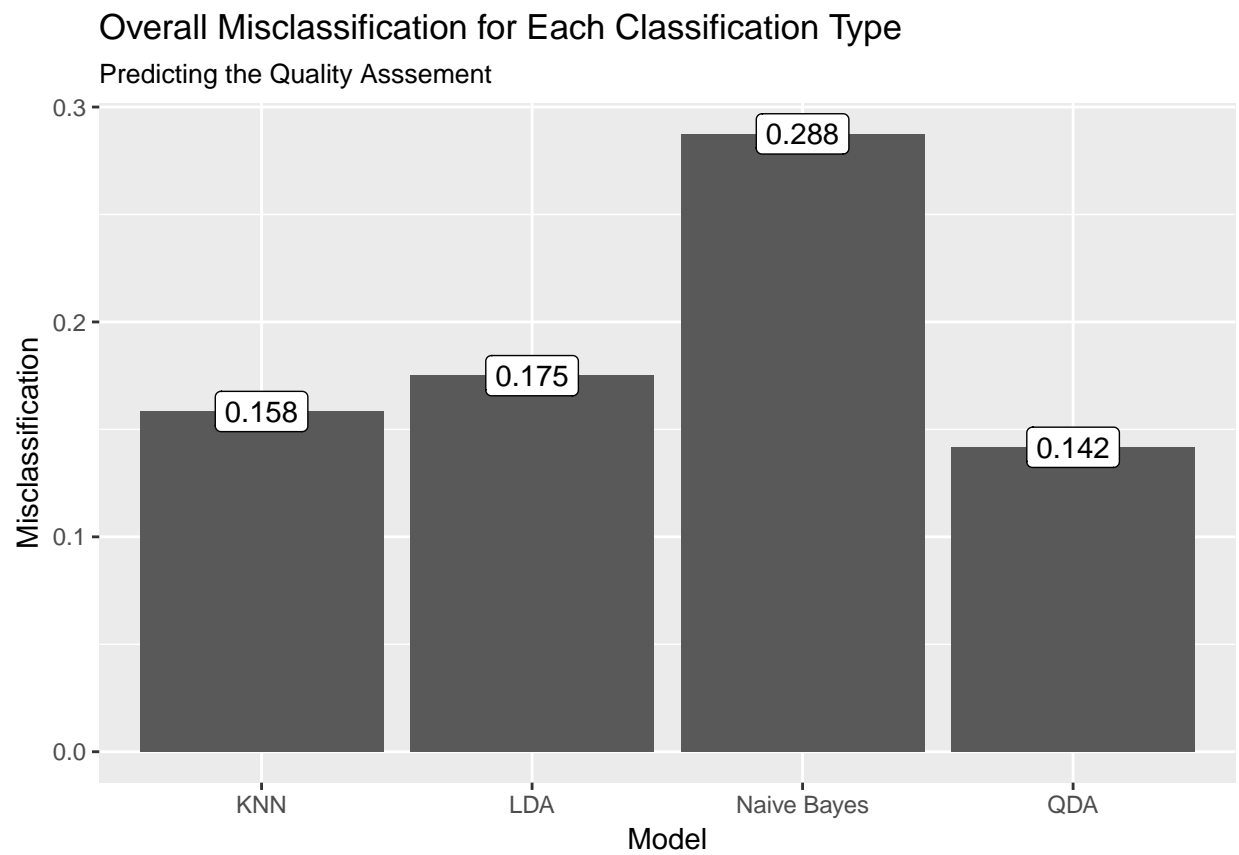
Classification

For classification purposes the wines were segregated into three different classes. These classes include “good” (quality >7), “medium” (quality between 4 and 7) and “poor”(quality < 4).

Model Selection

Residual Analysis

Comparison of Models



Discussion

Conclusion

Issues