

# Final Report for Red Wine Analysis

*Ruiqiang Chen, Michael DeWitt, David Williams, Alex Vannoy*

*7/28/2017*

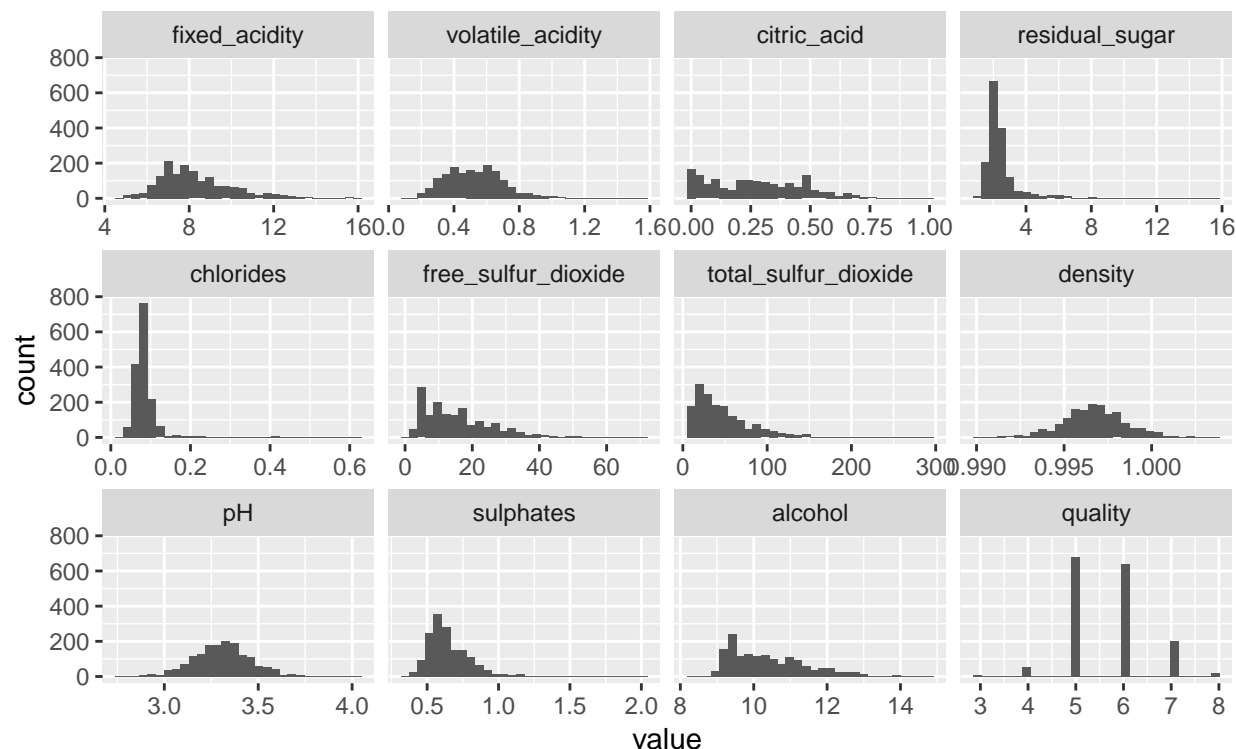
## 0.1 Introduction

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.

## 0.2 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:

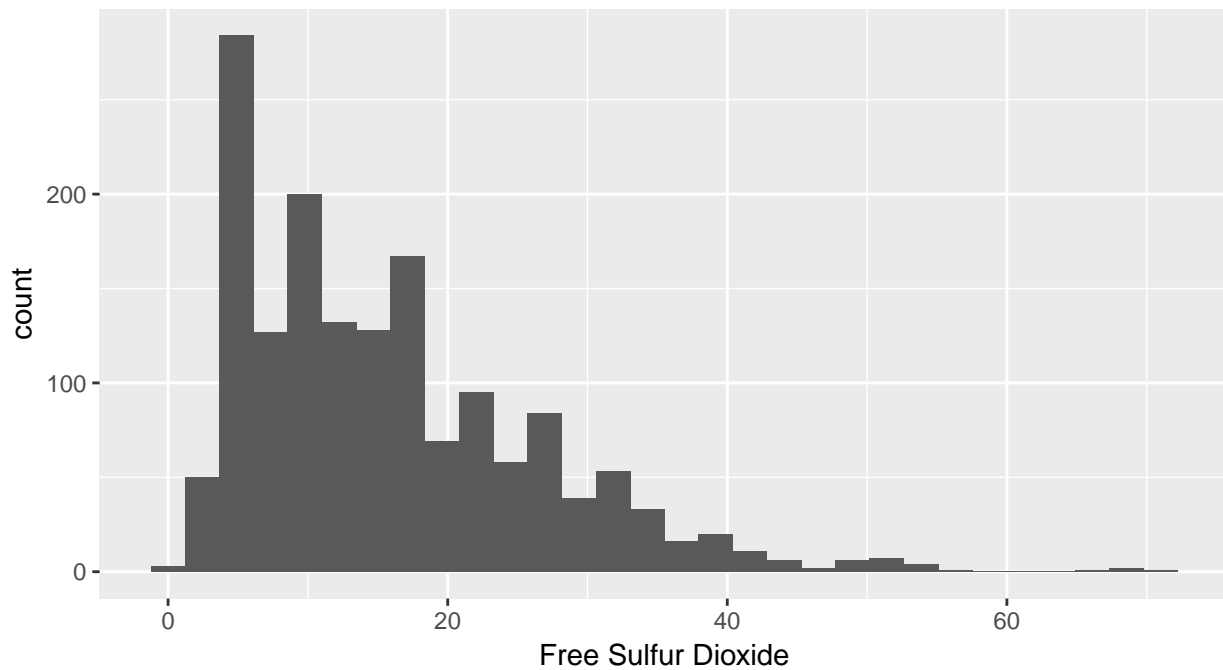
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



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
```

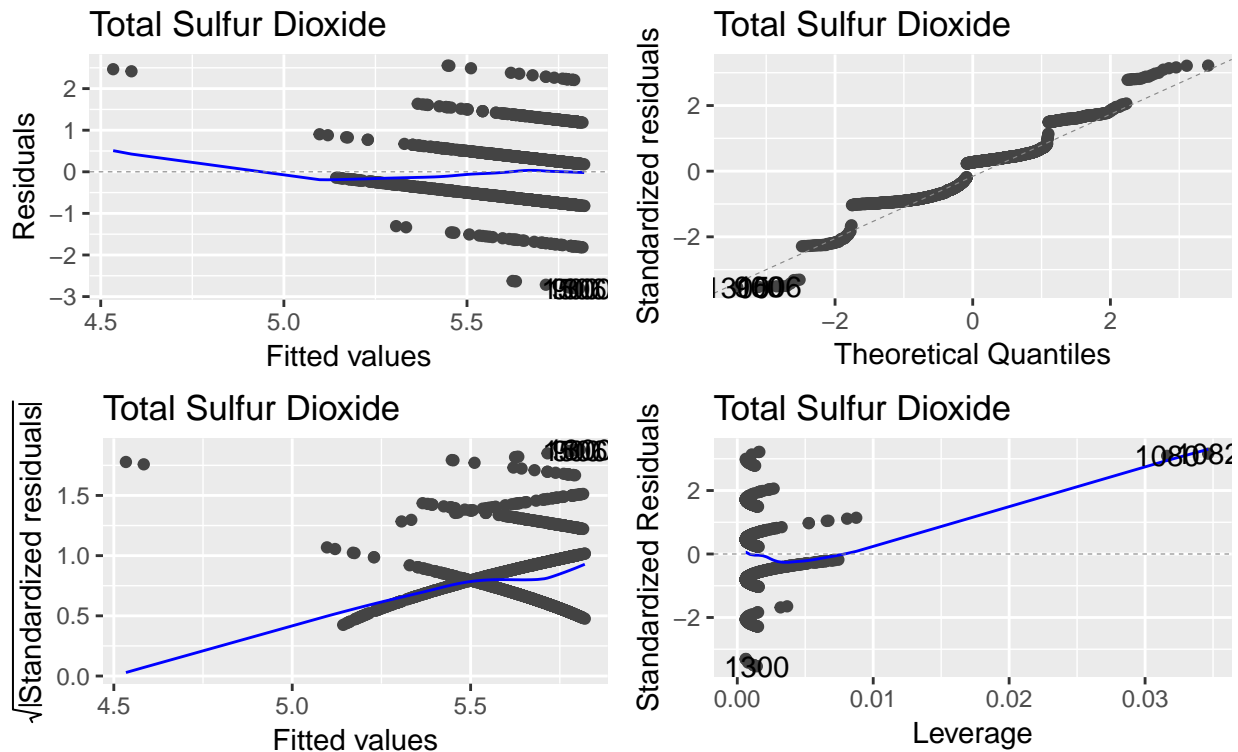
# Histogram of Free Sulfur Dioxide



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.

```
autoplot(fit2)+
  labs(title = "Total Sulfur Dioxide")
```



## 0.3 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.

### 0.3.1 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.

#### 0.3.1.1 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 3 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:

##### 0.3.1.1.1 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

#### 0.3.1.2 Ridge Regression

Ridge regression was performed on the dataset. cross validation was performed on the training data set to determine the optimum lambda for the ridge regression. This lambda was then using in a ridge regression model with the testing dataset.

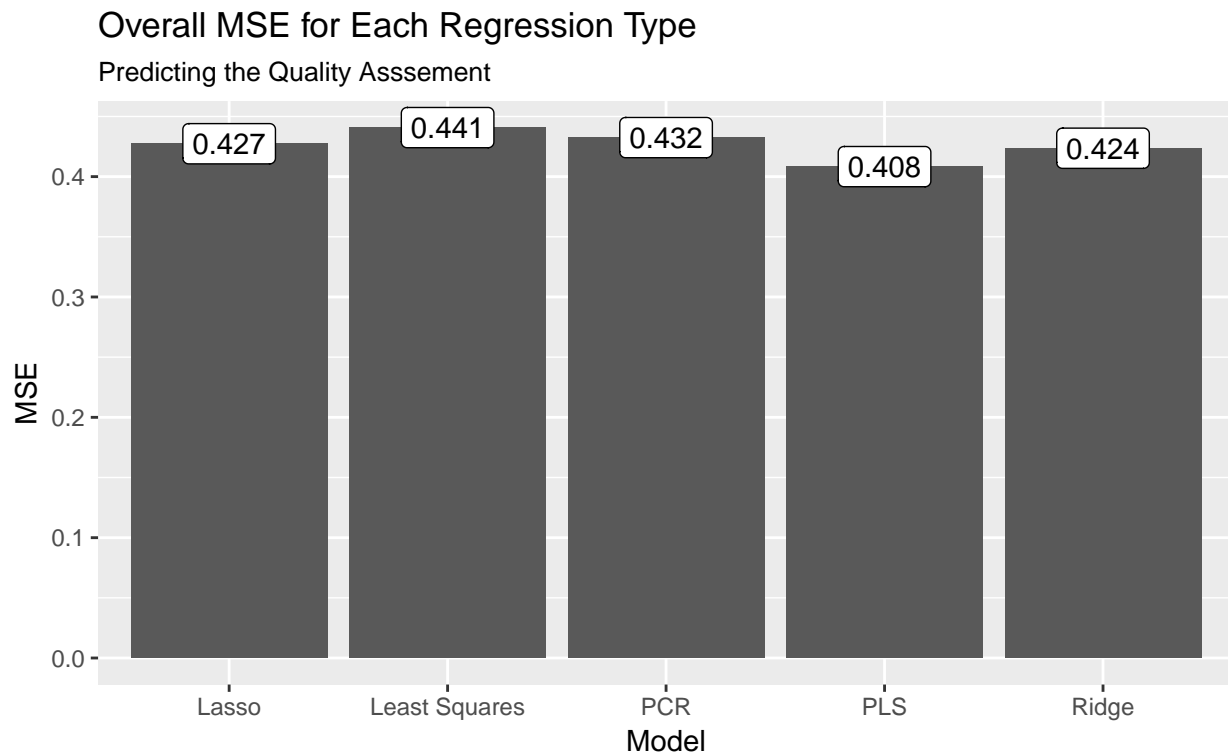
#### 0.3.1.3 Lasso Regression

Which variables were eliminated

#### 0.3.1.4 Principle Components Regression

Why did we choose the # we did PC #####Partial Least Squares Regression

#### 0.3.1.5 Model Selection



#### 0.3.1.6 Residual Analysis

### 0.3.2 Classification

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

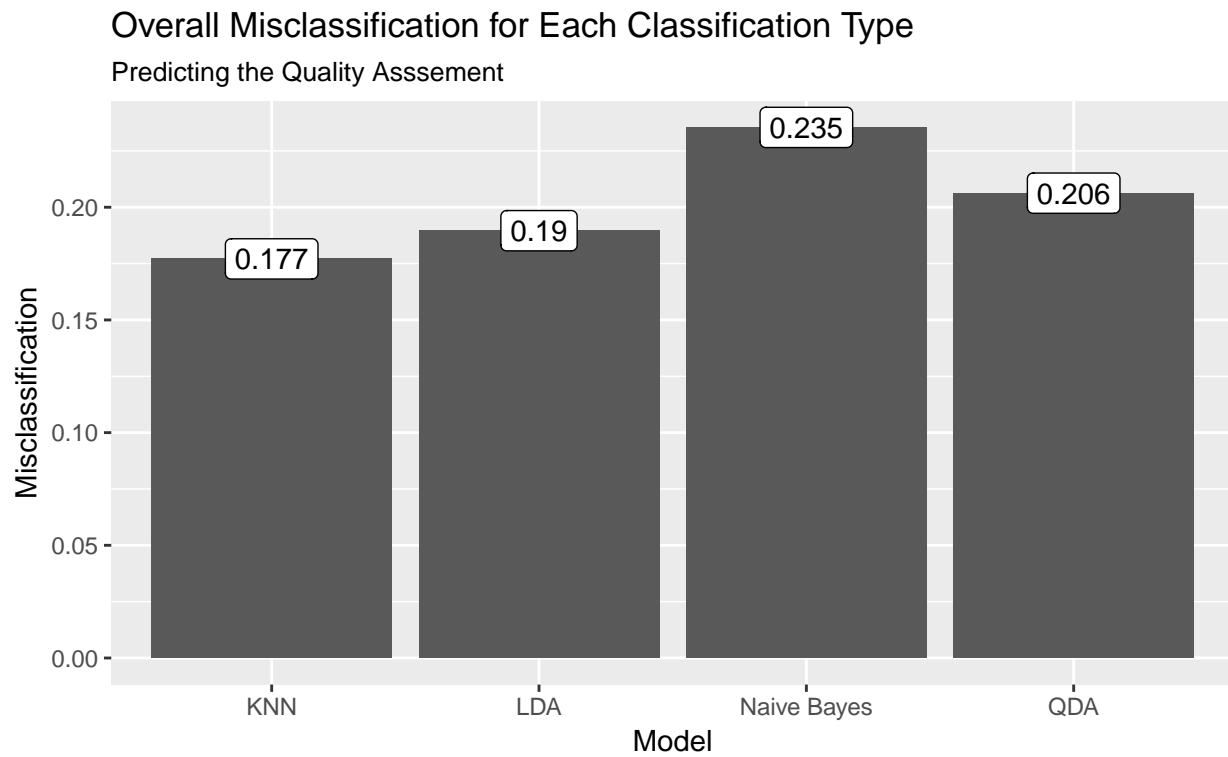
#### 0.3.2.1 Model Selection

#### 0.3.2.2 Residual Analysis

### 0.3.3 Comparison of Models

```
print(classification_plot)
```

```
## Warning: Width not defined. Set with `position_dodge(width = ?)`
```



**0.4 Discussion**

**0.5 Conclusion**

**0.6 Issues**