

DATA 621 Homework 5

Critical Thinking Group 1

November 29, 2021

Contents

Introduction	3
Problem	3
Data Exploration	4
Basic Statistics	4
Histogram of Variables	5
Relationship of Predictors to Target	5
Boxplots	6
Data Preparation	7
Identify Missing Values	7
Impute Missing Values	9
Identifying Multicollinearity	13
Build Models	14
Model Selection	20
Appendix	22

Prepared for:
Prof. Dr. Nasrin Khansari
City University of New York, School of Professional Studies - Data 621

DATA 621 – Business Analytics and Data Mining

Home Work 5

Prepared by:
Critical Thinking Group 1

Vic Chan
Gehad Gad
Evan McLaughlin
Bruno de Melo
Anjal Hussan
Zhouxin Shi
Sie Siong Wong

Introduction

Problem

Our goal is to explore, analyze and model a dataset containing information on approximately 12,000 commercially available wines. The variables are mostly related to the chemical properties of the wine being sold. The response variable is the number of sample cases of wine that were purchased by wine distribution companies after sampling a wine. These cases would be used to provide tasting samples to restaurants and wine stores around the United States. The more sample cases purchased, the more likely is a wine to be sold at a high end restaurant.

A large wine manufacturer is studying the data in order to predict the number of wine cases ordered based upon the wine characteristics. If the wine manufacturer can predict the number of cases, then that manufacturer will be able to adjust their wine offering to maximize sales.

The objective is to build a count regression model to predict the number of cases of wine that will be sold given certain properties of the wine.

Data Exploration

Below we'll display a few basic EDA techniques to gain insight into our wine dataset.

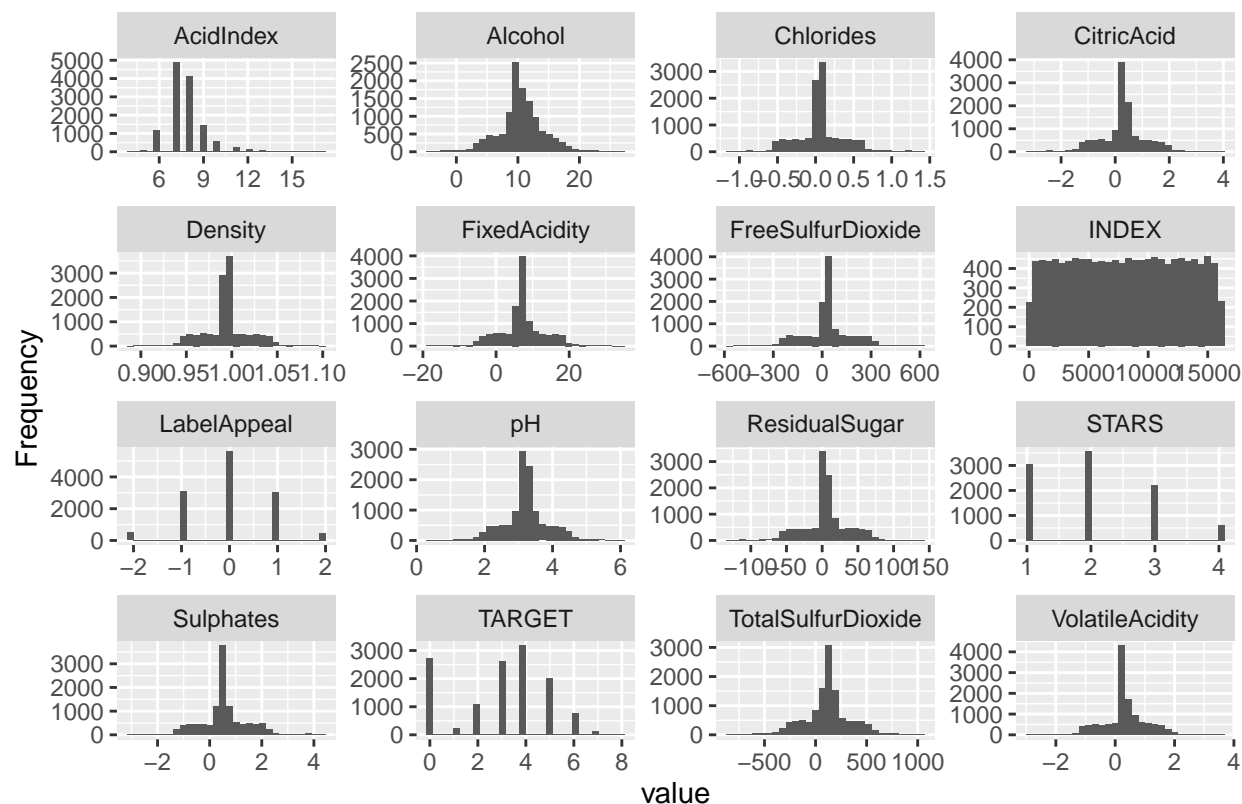
Basic Statistics

The data is 1.3 Mb in size. There are 12,795 rows and 15 columns (features). Of all 15 columns, 0 are discrete, 15 are continuous, and 0 are all missing. There are 8,200 missing values out of 191,925 data points.

```
##              n    mean    sd  median    min      max  skew
## INDEX      12795 8069.98 4656.91 8110.00    1.00 16129.00  0.00
## TARGET      12795   3.03   1.93   3.00    0.00   8.00 -0.33
## FixedAcidity 12795   7.08   6.32   6.90  -18.10  34.40 -0.02
## VolatileAcidity 12795   0.32   0.78   0.28   -2.79   3.68  0.02
## CitricAcid   12795   0.31   0.86   0.31   -3.24   3.86 -0.05
## ResidualSugar 12179   5.42  33.75   3.90 -127.80 141.15 -0.05
## Chlorides    12157   0.05   0.32   0.05   -1.17   1.35  0.03
## FreeSulfurDioxide 12148  30.85 148.71  30.00 -555.00 623.00  0.01
## TotalSulfurDioxide 12113 120.71 231.91 123.00 -823.00 1057.00 -0.01
## Density      12795   0.99   0.03   0.99   0.89   1.10 -0.02
## pH           12400   3.21   0.68   3.20   0.48   6.13  0.04
## Sulphates    11585   0.53   0.93   0.50   -3.13   4.24  0.01
## Alcohol      12142  10.49   3.73  10.40   -4.70  26.50 -0.03
## LabelAppeal  12795  -0.01   0.89   0.00   -2.00   2.00  0.01
## AcidIndex    12795   7.77   1.32   8.00   4.00  17.00  1.65
##              kurtosis
## INDEX          -1.20
## TARGET          -0.88
## FixedAcidity     1.67
## VolatileAcidity   1.83
## CitricAcid        1.84
## ResidualSugar     1.88
## Chlorides         1.79
## FreeSulfurDioxide 1.84
## TotalSulfurDioxide 1.67
## Density           1.90
## pH                1.65
## Sulphates         1.75
## Alcohol           1.54
## LabelAppeal       -0.26
## AcidIndex         5.19
```

It's useful to note a couple of things right off the bat with regard to our dataset: - There are several variables that have negative values. - ResidualSugar, Chlorides, FreeSulfurDioxide, and TotalSulfurDioxide all have quite a few missing values that we are going to need to deal with in order to assess the variables. - The Index column is useless and can be ignored.

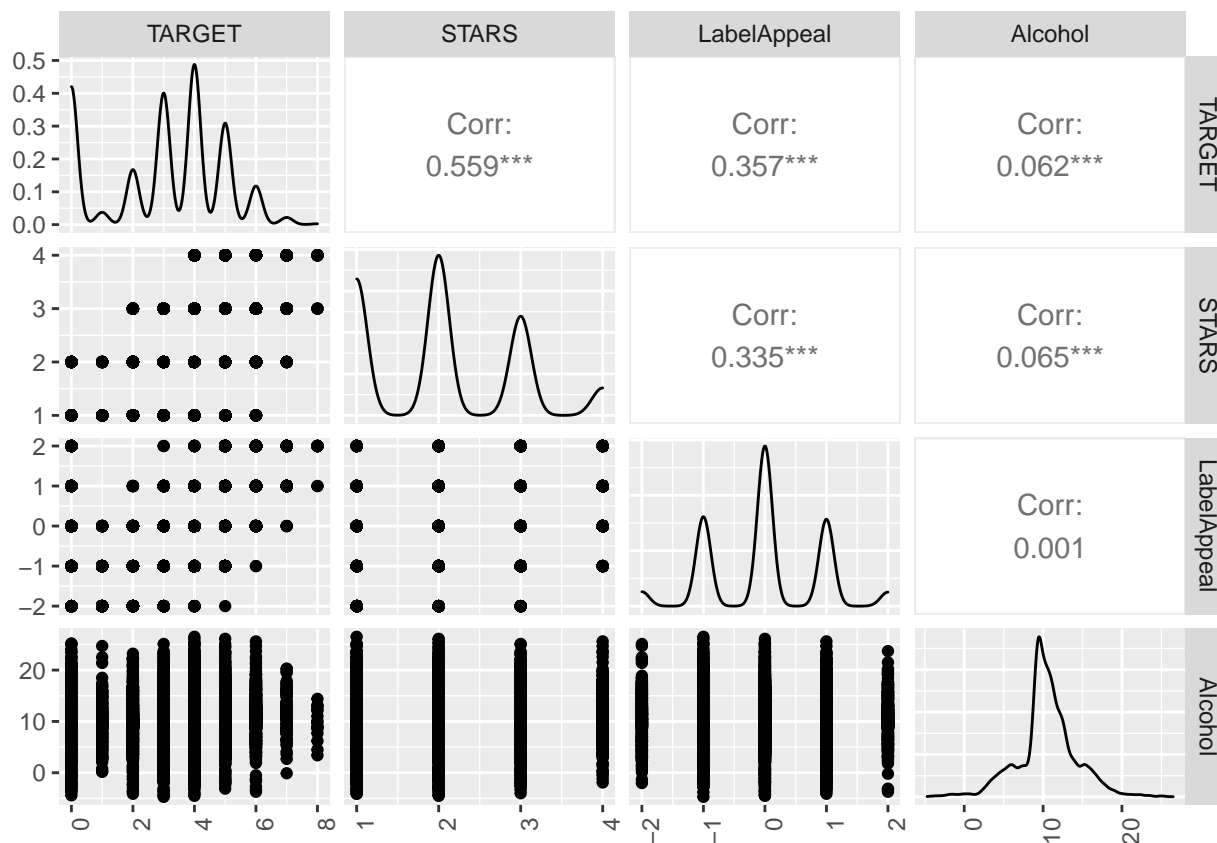
Histogram of Variables



Based on the histograms we can see that a lot of the variables distributions looks to be a normal distribution. We can see that **AcidIndex**, **STARS**, and **TARGET** are a bit skewed. One thing to note is that the **TARGET** variable has a lot of 0 cases sold. These 0 **TARGET** variables will need to be cleaned during the data prep phase as they can skew the results of the model.

Relationship of Predictors to Target

It is useful to assess the plots of each variable against the target variable. Using the GGPairs function from GGally we can plot some of the variables of interest to see if any of the variables correlates with the response variable **TARGET**. We will be making sure to include the variables **STARS** and **LabelAppeal** as it is believed that these two variables affect sales numbers



We can see that **STARS** and **Alcohol** has a bit of correlation with **TARGET**.

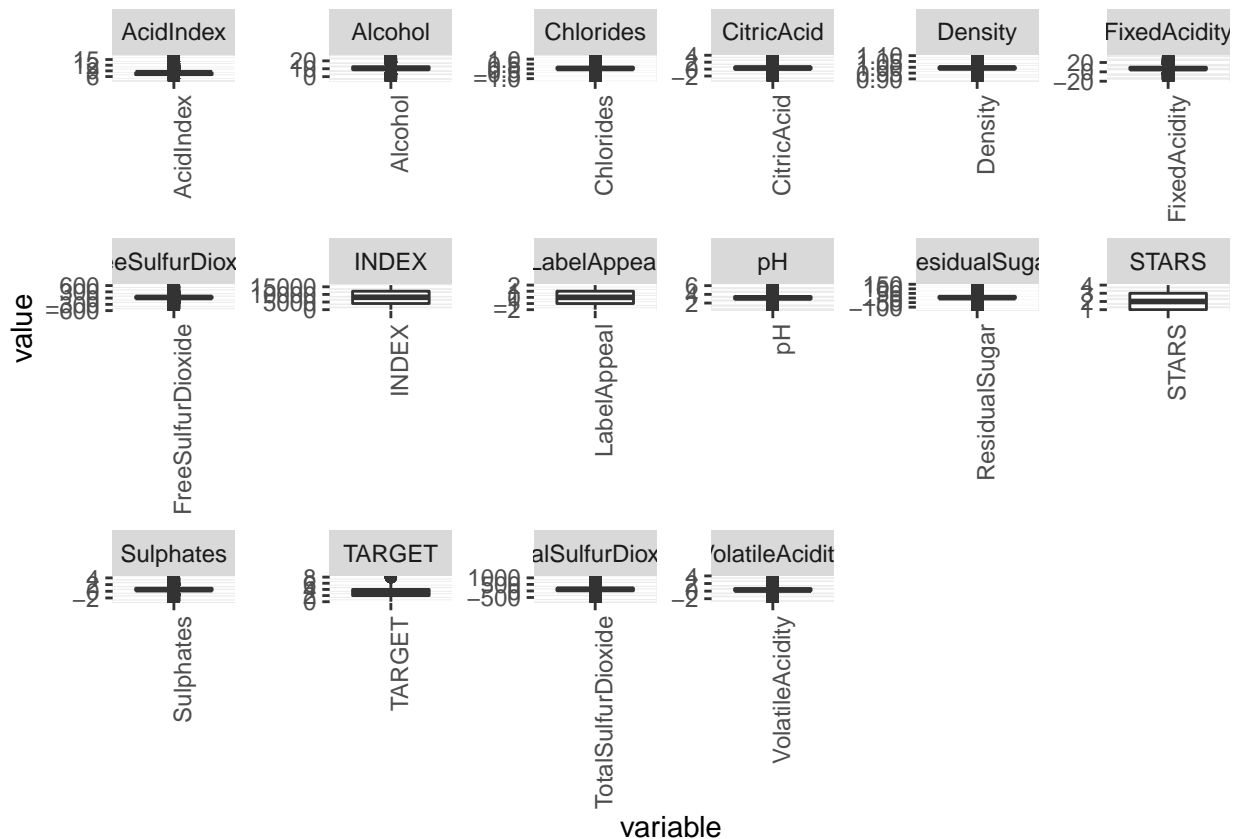
Boxplots

After observing our distributions, we can next assess the variables' relationship with our target variable (TARGET). It would be useful to establish the

```
bp_train <- train_data %>%
  gather(key = 'variable', value = 'value')

ggplot(bp_train, aes(variable, value)) +
  geom_boxplot() +
  facet_wrap(~variable, scales='free', ncol=6)
```

```
## Warning: Removed 8200 rows containing non-finite values (stat_boxplot).
```



When looking at the boxplot for **TARGET** we can see a very different picture compared to looking at the histogram. In the histogram it shows all the 0 **TARGETS** which can skew the modeling results while in the boxplot one can not easily point that out. This is the reason why we need to look at the data in multiple different ways.

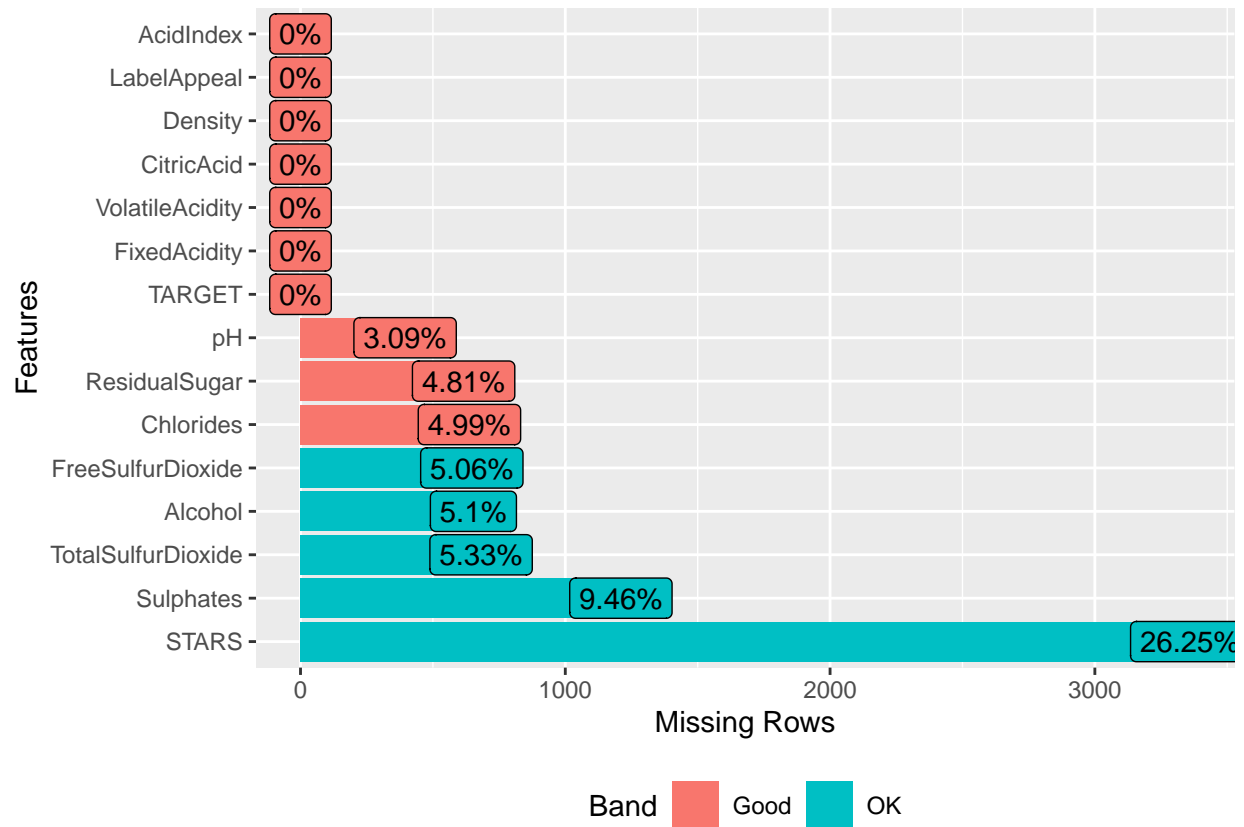
Data Preparation

Identify Missing Values

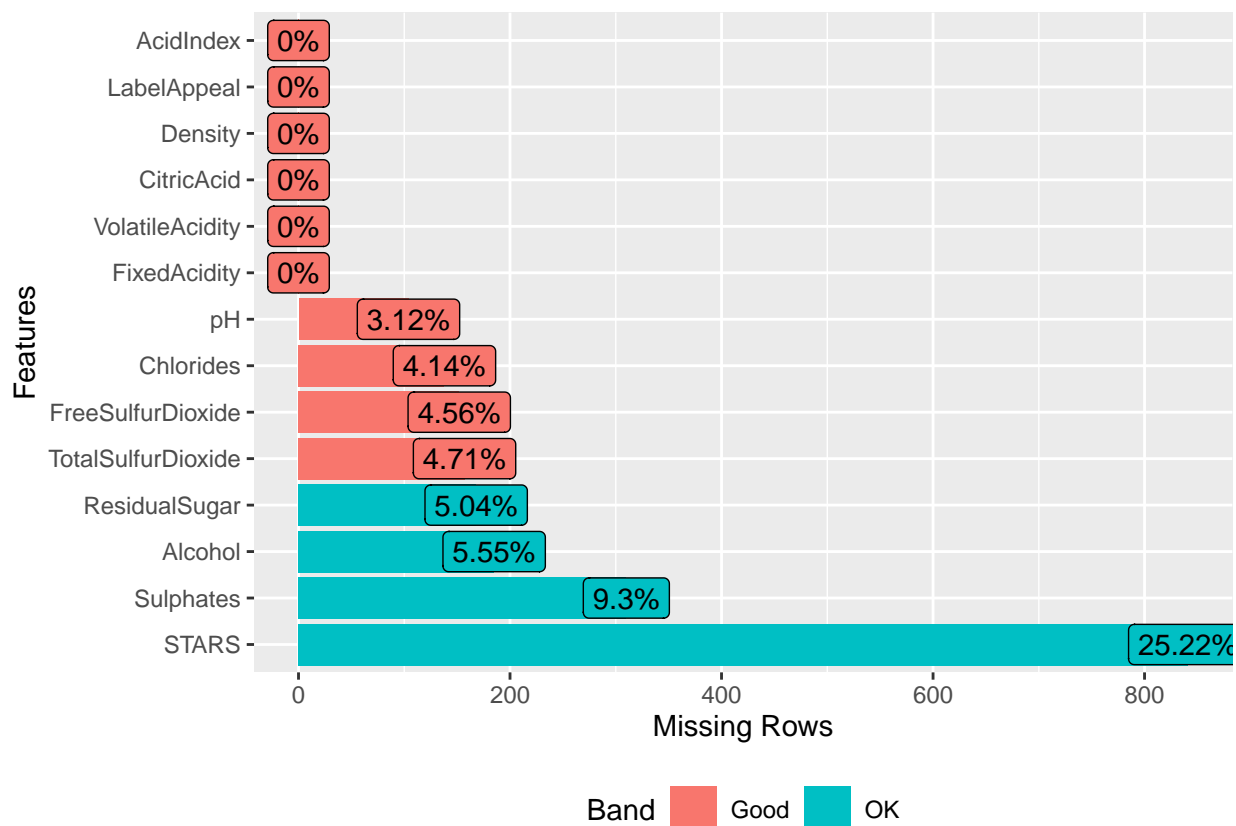
We can see that the same variables for both train data set and evaluation data set contains missing values. The variable that contains the most missing values is the **STAR** followed by **Sulphates**, **Alcohol**, **ResidualSugar**, **TotalSulfurDioxide**, **FreeSulfurDioxide**, **Chlorides**, and **pH**.

```
## Train Data Set Variables Missing Count
## 1 TARGET 0
## 2 FixedAcidity 0
## 3 VolatileAcidity 0
## 4 CitricAcid 0
## 5 ResidualSugar 616
## 6 Chlorides 638
## 7 FreeSulfurDioxide 647
## 8 TotalSulfurDioxide 682
## 9 Density 0
## 10 pH 395
```

## 11	Sulphates	1210
## 12	Alcohol	653
## 13	LabelAppeal	0
## 14	AcidIndex	0
## 15	STARS	3359



##	Train Data Set Variables	Missing Count
## 1	TARGET	3335
## 2	FixedAcidity	0
## 3	VolatileAcidity	0
## 4	CitricAcid	0
## 5	ResidualSugar	168
## 6	Chlorides	138
## 7	FreeSulfurDioxide	152
## 8	TotalSulfurDioxide	157
## 9	Density	0
## 10	pH	104
## 11	Sulphates	310
## 12	Alcohol	185
## 13	LabelAppeal	0
## 14	AcidIndex	0
## 15	STARS	841

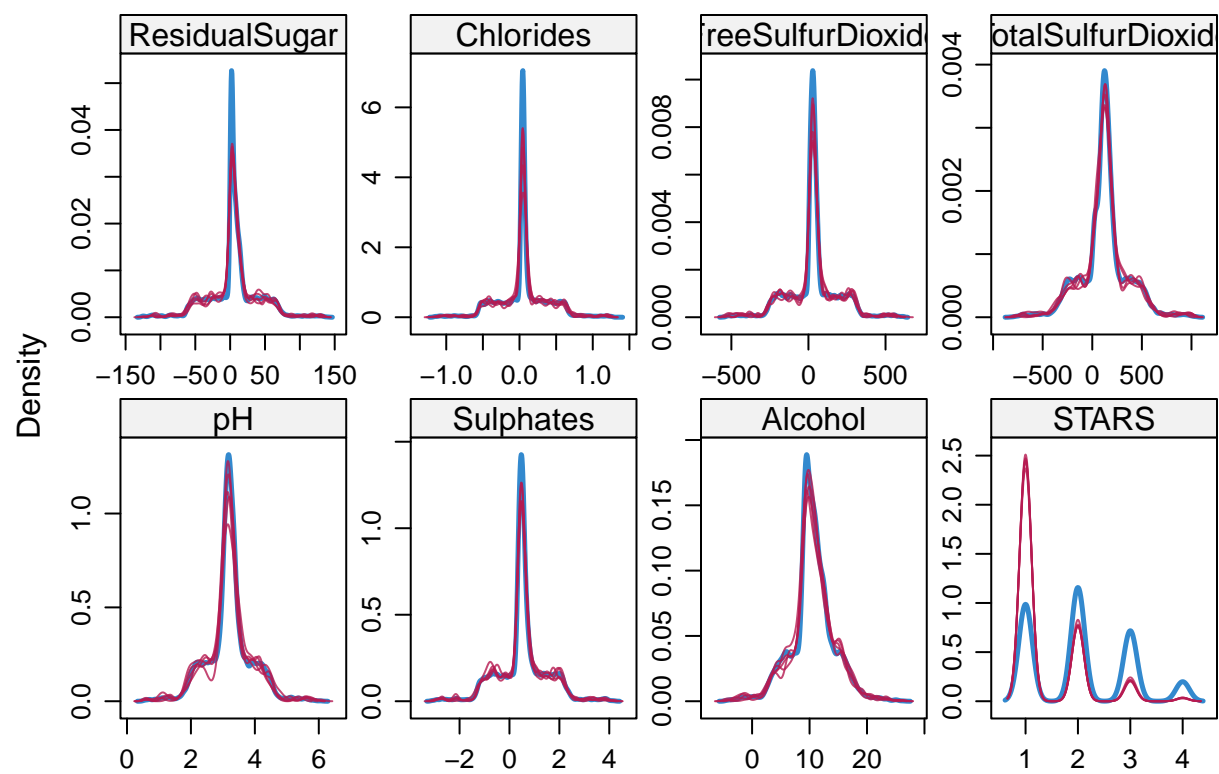


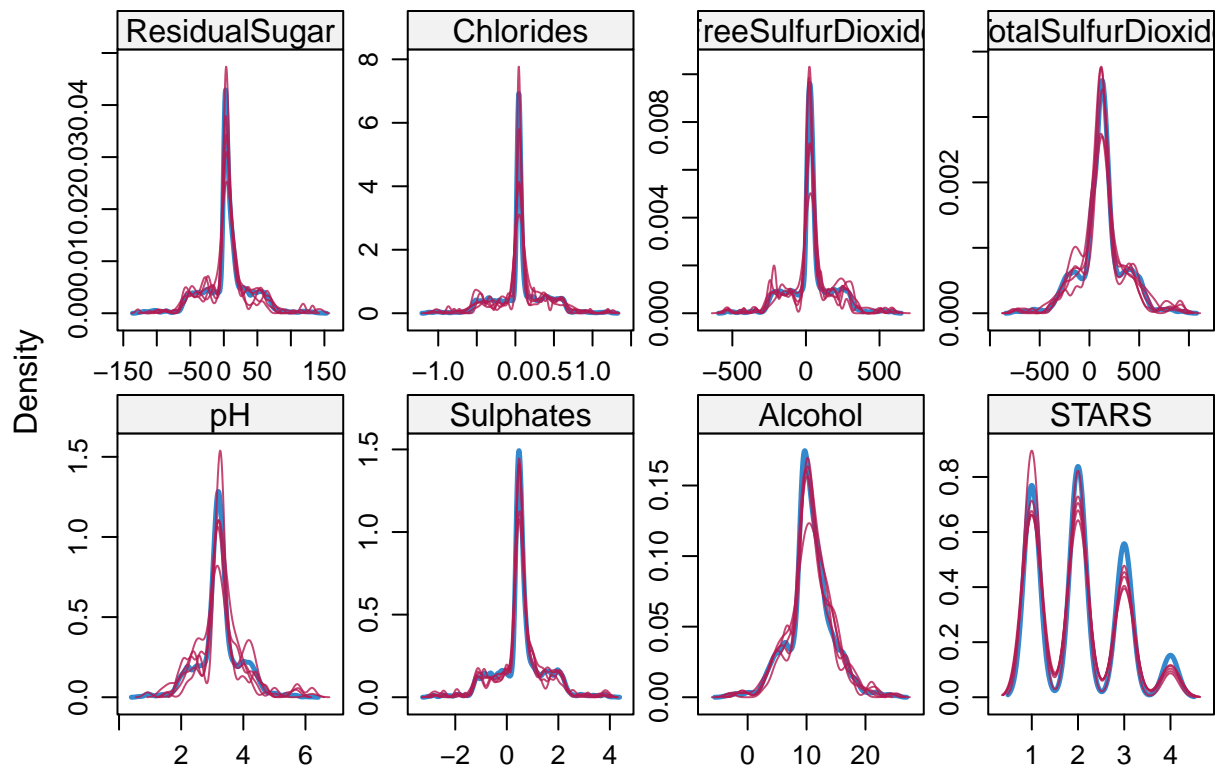
Our chart above does a good job of highlighting the missing values that will no doubt impact our analysis if we don't deal with them. In particular:

- STARS is missing 25% of its records. We could guess that the wines haven't been assessed and rated by experts.
- Sulphates is missing 9% of its records.
- Alcohol is missing 5% of its records. It is unlikely that a 0 here would indicate no alcohol in the wine, given that it's wine, so we can assume these values are missing.
- ResidualSugar is missing 5% of its records. There are some records that have 0 for ResidualSugar so these flagged records most certainly have missing values.
- A few more variables experiences missing values as well such as TotalSulfurDioxide (<5%), FreeSulfurDioxide (<5%), Chlorides (<5%), and pH (<4%).

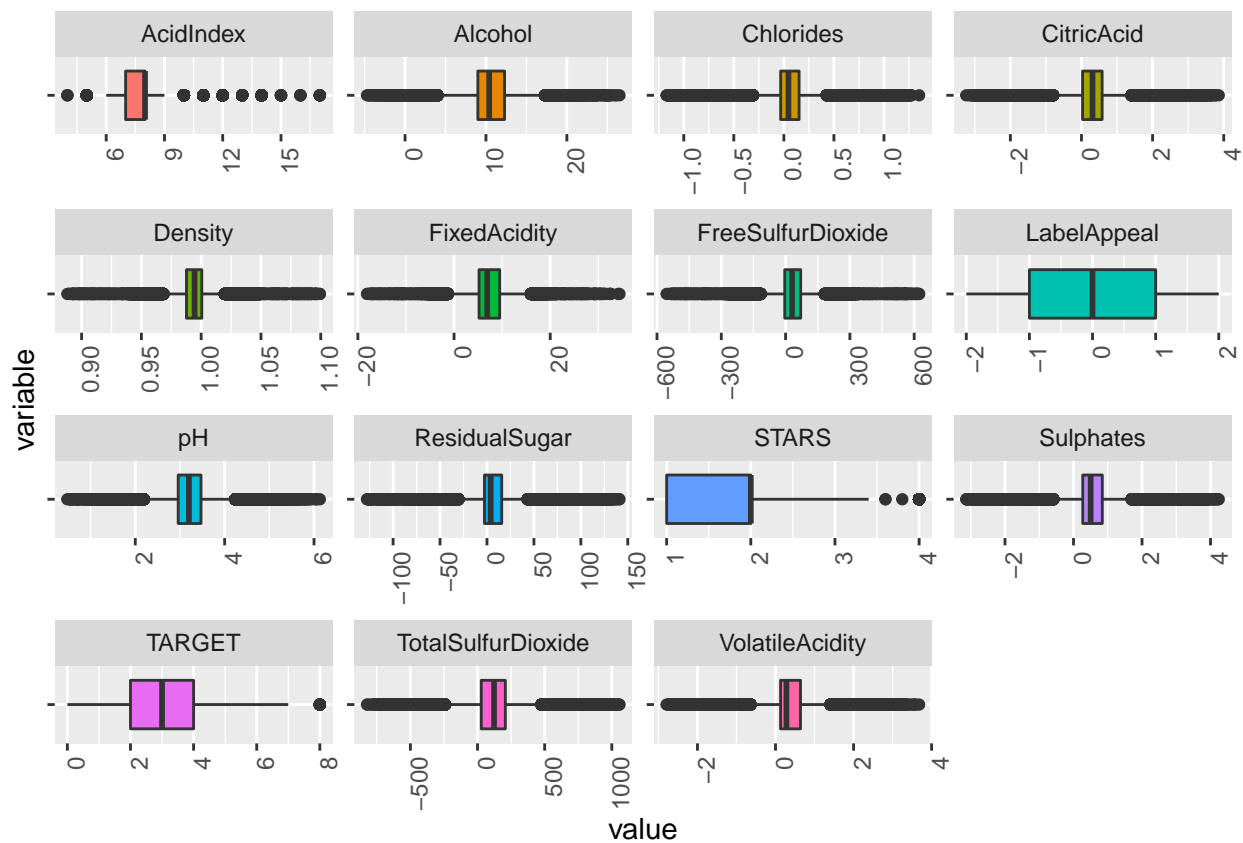
Impute Missing Values

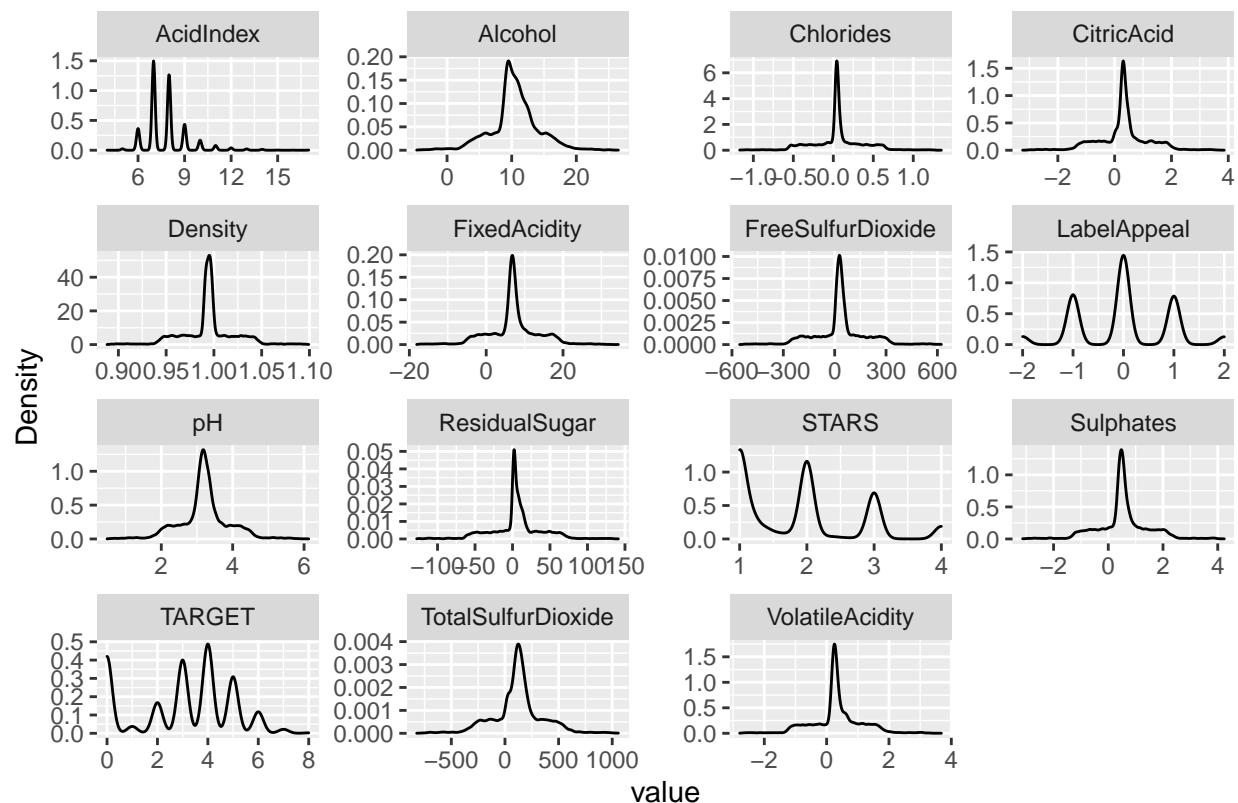
We assume the missing data are Missing at Random and choose to impute. The reason we want to impute the missing data rather than replacing with mean or median because of large number of missing values. If we're replacing with mean or median on the large number of missing values, can result in loss of variation in data. We're imputing the missing data using the MICE package. The method of predictive mean matching (PMM) is selected for continuous variables.





Next, we take average values of the 5 imputed data set as a final train data set used for building models.





Because we'll use the Poisson or Negative Binomial regression to build count regression models in GLM approach, transformation for each variable to make them look normal is not required. Diagnostics of actual outliers or influential points can be identified in the build models section through plots such as residuals vs fitted, standardized residuals vs fitted, etc.

Identifying Multicollinearity

```
stack(sort(cor(complete_train_data[,1], complete_train_data[,2:ncol(complete_train_data)]), decreasing = TRUE))
```

```
##          values          ind
## 1  0.646881500          STARS
## 2  0.356500469      LabelAppeal
## 3  0.063601452          Alcohol
## 4  0.051714048 TotalSulfurDioxide
## 5  0.044495907 FreeSulfurDioxide
## 6  0.015992585      ResidualSugar
## 7  0.008684633      CitricAcid
## 8 -0.009220723           pH
## 9 -0.035517502          Density
## 10 -0.039429538        Chlorides
## 11 -0.041250213        Sulphates
## 12 -0.049010939      FixedAcidity
## 13 -0.088793212 VolatileAcidity
## 14 -0.246049449      AcidIndex
```

```
correlation = cor(complete_train_data, use = 'pairwise.complete.obs')
#corrplot(correlation, method='ellipse', type = 'lower', order = 'hclust', col=brewer.pal(n=8, name="RdY
```

After our EDA and Data Prep, we can render some judgments on the dataset, including that there are some variables with a weak enough relationship with our TARGET that we can plan to drop them. We also need to plan on dealing with the many outliers that could skew our models and pay close attention to multicollinearity, especially as it relates to the relationship between STARS and LabelAppeal. These two features also happen to have the strongest correlation with the TARGET.

Build Models

Poisson Model 1

```
Model1 <- glm(TARGET ~ FixedAcidity + VolatileAcidity + CitricAcid + ResidualSugar +
Chlorides + FreeSulfurDioxide + TotalSulfurDioxide + Density +
pH + Sulphates + Alcohol +
               as.factor(LabelAppeal) +
               as.factor(AcidIndex) +
               as.factor(STARS),
               data=complete_train_data,
               family=poisson)
summary(Model1)
```

```
##
## Call:
## glm(formula = TARGET ~ FixedAcidity + VolatileAcidity + CitricAcid +
##      ResidualSugar + Chlorides + FreeSulfurDioxide + TotalSulfurDioxide +
##      Density + pH + Sulphates + Alcohol + as.factor(LabelAppeal) +
##      as.factor(AcidIndex) + as.factor(STARS), family = poisson,
##      data = complete_train_data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.1912  -0.6238  -0.0108   0.4660   3.3938
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      5.461e-01  3.738e-01   1.461 0.144054
## FixedAcidity      2.398e-04  8.205e-04   0.292 0.770142
## VolatileAcidity  -3.085e-02  6.567e-03 -4.698 2.62e-06 ***
## CitricAcid        5.314e-03  5.911e-03   0.899 0.368640
## ResidualSugar    -2.938e-05  1.543e-04 -0.190 0.848990
## Chlorides        -5.119e-02  1.643e-02 -3.116 0.001836 **
## FreeSulfurDioxide  8.369e-05  3.483e-05   2.403 0.016276 *
## TotalSulfurDioxide  5.640e-05  2.259e-05   2.497 0.012527 *
## Density          -2.480e-01  1.925e-01 -1.289 0.197464
## pH               -8.492e-03  7.641e-03 -1.111 0.266418
## Sulphates        -8.037e-03  5.690e-03 -1.412 0.157819
## Alcohol           4.109e-03  1.406e-03   2.923 0.003463 **
## as.factor(LabelAppeal)-1 1.665e-01  3.805e-02  4.374 1.22e-05 ***
## as.factor(LabelAppeal)0  3.269e-01  3.714e-02  8.802 < 2e-16 ***
```

```

## as.factor(LabelAppeal)1    4.487e-01  3.778e-02  11.875 < 2e-16 ***
## as.factor(LabelAppeal)2    5.753e-01  4.260e-02  13.506 < 2e-16 ***
## as.factor(AcidIndex)5      1.050e-01  3.241e-01   0.324 0.746061
## as.factor(AcidIndex)6      1.327e-01  3.187e-01   0.416 0.677218
## as.factor(AcidIndex)7      1.095e-01  3.184e-01   0.344 0.730993
## as.factor(AcidIndex)8      7.783e-02  3.185e-01   0.244 0.806921
## as.factor(AcidIndex)9     -9.913e-04  3.188e-01  -0.003 0.997519
## as.factor(AcidIndex)10    -2.060e-01  3.199e-01  -0.644 0.519615
## as.factor(AcidIndex)11    -4.961e-01  3.235e-01  -1.533 0.125182
## as.factor(AcidIndex)12    -5.717e-01  3.292e-01  -1.737 0.082449 .
## as.factor(AcidIndex)13    -3.635e-01  3.320e-01  -1.095 0.273666
## as.factor(AcidIndex)14    -4.180e-01  3.446e-01  -1.213 0.225164
## as.factor(AcidIndex)15    -2.626e-03  4.050e-01  -0.006 0.994827
## as.factor(AcidIndex)16    -5.236e-01  5.494e-01  -0.953 0.340597
## as.factor(AcidIndex)17    -1.080e+00  5.494e-01  -1.965 0.049378 *
## as.factor(STARS)1.2      -1.716e+00  6.003e-02 -28.590 < 2e-16 ***
## as.factor(STARS)1.4      -4.462e-01  4.677e-02  -9.539 < 2e-16 ***
## as.factor(STARS)1.6       2.183e-01  4.633e-02   4.712 2.45e-06 ***
## as.factor(STARS)1.8       4.515e-01  4.211e-02  10.722 < 2e-16 ***
## as.factor(STARS)2         6.209e-01  1.420e-02  43.733 < 2e-16 ***
## as.factor(STARS)2.2       6.353e-01  5.056e-02  12.565 < 2e-16 ***
## as.factor(STARS)2.4       6.916e-01  5.322e-02  12.994 < 2e-16 ***
## as.factor(STARS)2.6       7.541e-01  7.030e-02  10.727 < 2e-16 ***
## as.factor(STARS)2.8       7.628e-01  8.548e-02   8.924 < 2e-16 ***
## as.factor(STARS)3         7.493e-01  1.552e-02  48.269 < 2e-16 ***
## as.factor(STARS)3.2       9.766e-01  1.207e-01   8.094 5.78e-16 ***
## as.factor(STARS)3.4       9.576e-01  2.300e-01   4.163 3.14e-05 ***
## as.factor(STARS)3.6       9.860e-01  2.686e-01   3.672 0.000241 ***
## as.factor(STARS)3.8       1.226e+00  3.541e-01   3.462 0.000535 ***
## as.factor(STARS)4         8.739e-01  2.160e-02  40.454 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##    Null deviance: 22861  on 12794  degrees of freedom
## Residual deviance: 12514  on 12751  degrees of freedom
## AIC: 44544
##
## Number of Fisher Scoring iterations: 6

```

Poisson Model 2

```

Model2 <- glm(TARGET ~ VolatileAcidity + TotalSulfurDioxide + Alcohol +
              as.factor(LabelAppeal) +
              as.factor(AcidIndex) +
              as.factor(STARS),
              data=complete_train_data,
              family=poisson)
summary(Model2)

```

```

##
## Call:

```

```

## glm(formula = TARGET ~ VolatileAcidity + TotalSulfurDioxide +
##     Alcohol + as.factor(LabelAppeal) + as.factor(AcidIndex) +
##     as.factor(STARS), family = poisson, data = complete_train_data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -3.1880  -0.6187  -0.0066   0.4671   3.4316
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      2.467e-01  3.207e-01   0.769 0.441674
## VolatileAcidity  -3.124e-02  6.565e-03 -4.759 1.95e-06 ***
## TotalSulfurDioxide  5.658e-05  2.257e-05   2.507 0.012185 *
## Alcohol          4.221e-03  1.404e-03   3.006 0.002650 **
## as.factor(LabelAppeal)-1  1.674e-01  3.805e-02   4.399 1.09e-05 ***
## as.factor(LabelAppeal)0  3.275e-01  3.714e-02   8.820 < 2e-16 ***
## as.factor(LabelAppeal)1  4.492e-01  3.778e-02  11.891 < 2e-16 ***
## as.factor(LabelAppeal)2  5.750e-01  4.259e-02  13.500 < 2e-16 ***
## as.factor(AcidIndex)5    1.280e-01  3.238e-01   0.395 0.692566
## as.factor(AcidIndex)6    1.548e-01  3.184e-01   0.486 0.626799
## as.factor(AcidIndex)7    1.323e-01  3.181e-01   0.416 0.677543
## as.factor(AcidIndex)8    1.006e-01  3.181e-01   0.316 0.751891
## as.factor(AcidIndex)9    2.187e-02  3.184e-01   0.069 0.945233
## as.factor(AcidIndex)10  -1.861e-01  3.195e-01  -0.582 0.560342
## as.factor(AcidIndex)11  -4.772e-01  3.231e-01  -1.477 0.139758
## as.factor(AcidIndex)12  -5.529e-01  3.288e-01  -1.682 0.092638 .
## as.factor(AcidIndex)13  -3.421e-01  3.316e-01  -1.031 0.302335
## as.factor(AcidIndex)14  -3.909e-01  3.441e-01  -1.136 0.255956
## as.factor(AcidIndex)15    4.118e-02  4.047e-01   0.102 0.918940
## as.factor(AcidIndex)16  -4.899e-01  5.489e-01  -0.892 0.372135
## as.factor(AcidIndex)17  -1.071e+00  5.490e-01  -1.950 0.051154 .
## as.factor(STARS)1.2     -1.719e+00  6.002e-02 -28.638 < 2e-16 ***
## as.factor(STARS)1.4     -4.476e-01  4.676e-02  -9.572 < 2e-16 ***
## as.factor(STARS)1.6      2.173e-01  4.632e-02   4.691 2.72e-06 ***
## as.factor(STARS)1.8      4.515e-01  4.210e-02  10.723 < 2e-16 ***
## as.factor(STARS)2        6.218e-01  1.419e-02  43.817 < 2e-16 ***
## as.factor(STARS)2.2      6.332e-01  5.054e-02  12.528 < 2e-16 ***
## as.factor(STARS)2.4      6.891e-01  5.319e-02  12.956 < 2e-16 ***
## as.factor(STARS)2.6      7.553e-01  7.028e-02  10.747 < 2e-16 ***
## as.factor(STARS)2.8      7.691e-01  8.546e-02   9.000 < 2e-16 ***
## as.factor(STARS)3        7.509e-01  1.551e-02  48.406 < 2e-16 ***
## as.factor(STARS)3.2      9.800e-01  1.206e-01   8.124 4.50e-16 ***
## as.factor(STARS)3.4      9.582e-01  2.300e-01   4.166 3.09e-05 ***
## as.factor(STARS)3.6      9.889e-01  2.685e-01   3.683 0.000230 ***
## as.factor(STARS)3.8      1.229e+00  3.539e-01   3.474 0.000513 ***
## as.factor(STARS)4        8.745e-01  2.159e-02  40.501 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 22861  on 12794  degrees of freedom
## Residual deviance: 12536  on 12759  degrees of freedom
## AIC: 44550

```



```
##
## Number of Fisher Scoring iterations: 6
```

Baseline model

```
Model3 <- glm(TARGET ~ ., complete_train_data, family = poisson(link = "log"))
summary(Model3)
```

```
##
## Call:
## glm(formula = TARGET ~ ., family = poisson(link = "log"), data = complete_train_data)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.8571  -0.6913   0.1207   0.6226   2.6738
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.477e+00  1.960e-01   7.534 4.91e-14 ***
## FixedAcidity   -5.313e-04  8.197e-04  -0.648 0.516928
## VolatileAcidity -3.910e-02  6.518e-03  -5.998 1.99e-09 ***
## CitricAcid      1.033e-02  5.896e-03   1.752 0.079813 .
## ResidualSugar    5.835e-05  1.539e-04   0.379 0.704539
## Chlorides       -5.080e-02  1.637e-02  -3.104 0.001910 **
## FreeSulfurDioxide 1.412e-04  3.493e-05   4.043 5.28e-05 ***
## TotalSulfurDioxide 8.434e-05  2.258e-05   3.735 0.000187 ***
## Density         -3.413e-01  1.923e-01  -1.774 0.075998 .
## pH              -1.812e-02  7.611e-03  -2.381 0.017285 *
## Sulphates       -1.532e-02  5.686e-03  -2.694 0.007060 **
## Alcohol         2.676e-03  1.404e-03   1.906 0.056683 .
## LabelAppeal     1.365e-01  6.111e-03  22.333 < 2e-16 ***
## AcidIndex       -9.587e-02  4.525e-03 -21.187 < 2e-16 ***
## STARS           3.582e-01  5.754e-03  62.244 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for poisson family taken to be 1)
##
##      Null deviance: 22861  on 12794  degrees of freedom
## Residual deviance: 15774  on 12780  degrees of freedom
## AIC: 47746
##
## Number of Fisher Scoring iterations: 5
```

quasi-Poisson Model

```
Model4 <- glm(TARGET ~ ., family=quasipoisson, complete_train_data)
summary(Model4)
```

```
##
## Call:
## glm(formula = TARGET ~ ., family = quasipoisson, data = complete_train_data)
```

```
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.8571  -0.6913   0.1207   0.6226   2.6738
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.477e+00  1.833e-01   8.057 8.54e-16 ***
## FixedAcidity   -5.313e-04  7.666e-04  -0.693 0.488320
## VolatileAcidity -3.910e-02  6.096e-03  -6.414 1.47e-10 ***
## CitricAcid      1.033e-02  5.514e-03   1.873 0.061067 .
## ResidualSugar    5.835e-05  1.439e-04   0.405 0.685132
## Chlorides      -5.080e-02  1.531e-02  -3.319 0.000906 ***
## FreeSulfurDioxide 1.412e-04  3.267e-05   4.323 1.55e-05 ***
## TotalSulfurDioxide 8.434e-05  2.112e-05   3.994 6.53e-05 ***
## Density        -3.413e-01  1.799e-01  -1.897 0.057802 .
## pH             -1.812e-02  7.118e-03  -2.546 0.010922 *
## Sulphates      -1.532e-02  5.318e-03  -2.881 0.003975 **
## Alcohol         2.676e-03  1.313e-03   2.038 0.041589 *
## LabelAppeal     1.365e-01  5.715e-03  23.880 < 2e-16 ***
## AcidIndex       -9.587e-02  4.231e-03 -22.656 < 2e-16 ***
## STARS           3.582e-01  5.381e-03  66.558 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for quasipoisson family taken to be 0.874578)
##
##      Null deviance: 22861  on 12794  degrees of freedom
## Residual deviance: 15774  on 12780  degrees of freedom
## AIC: NA
##
## Number of Fisher Scoring iterations: 5
```

Negative binomial regression

```
Model15 <- glm.nb(TARGET ~ ., data = complete_train_data)
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
```

```
## Warning in theta.ml(Y, mu, sum(w), w, limit = control$maxit, trace =
## control$trace > : iteration limit reached
```

```
summary(Model15)
```

```
##
## Call:
## glm.nb(formula = TARGET ~ ., data = complete_train_data, init.theta = 49193.33129,
##      link = log)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -2.8571 -0.6913 0.1207 0.6226 2.6738
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.477e+00  1.960e-01   7.534 4.91e-14 ***
## FixedAcidity   -5.313e-04  8.197e-04  -0.648 0.516935
## VolatileAcidity -3.910e-02  6.518e-03  -5.998 1.99e-09 ***
## CitricAcid      1.033e-02  5.896e-03   1.752 0.079821 .
## ResidualSugar    5.835e-05  1.539e-04   0.379 0.704534
## Chlorides       -5.080e-02  1.637e-02  -3.104 0.001911 **
## FreeSulfurDioxide 1.412e-04  3.493e-05   4.043 5.28e-05 ***
## TotalSulfurDioxide 8.434e-05  2.258e-05   3.735 0.000188 ***
## Density        -3.413e-01  1.923e-01  -1.774 0.076002 .
## pH              -1.812e-02  7.611e-03  -2.381 0.017286 *
## Sulphates       -1.532e-02  5.687e-03  -2.694 0.007061 **
## Alcohol         2.676e-03  1.404e-03   1.906 0.056694 .
## LabelAppeal     1.365e-01  6.111e-03  22.332 < 2e-16 ***
## AcidIndex       -9.587e-02  4.525e-03 -21.187 < 2e-16 ***
## STARS           3.582e-01  5.755e-03  62.242 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for Negative Binomial(49193.33) family taken to be 1)
##
##      Null deviance: 22860  on 12794  degrees of freedom
## Residual deviance: 15774  on 12780  degrees of freedom
## AIC: 47749
##
## Number of Fisher Scoring iterations: 1
##
##              Theta: 49193
##          Std. Err.: 55792
## Warning while fitting theta: iteration limit reached
##
## 2 x log-likelihood: -47716.51
```

Zero Inflated Count Models

```
Model6 <- zeroinfl(TARGET ~ ., data = complete_train_data)
summary(Model6)
```

```
##
## Call:
## zeroinfl(formula = TARGET ~ ., data = complete_train_data)
##
## Pearson residuals:
##      Min      1Q   Median      3Q      Max
## -1.983021 -0.453229  0.001259  0.403750  6.459031
##
## Count model coefficients (poisson with log link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.370e+00  2.025e-01   6.767 1.31e-11 ***
```

```
## FixedAcidity      2.006e-04  8.416e-04   0.238 0.811564
## VolatileAcidity  -1.158e-02  6.736e-03  -1.719 0.085681 .
## CitricAcid       7.397e-04  6.038e-03   0.122 0.902506
## ResidualSugar    -8.595e-05  1.582e-04  -0.543 0.587026
## Chlorides        -2.231e-02  1.687e-02  -1.322 0.186025
## FreeSulfurDioxide 3.138e-05  3.528e-05   0.890 0.373635
## TotalSulfurDioxide -2.307e-05  2.254e-05  -1.024 0.305910
## Density          -3.042e-01  1.986e-01  -1.531 0.125692
## pH               5.824e-03  7.850e-03   0.742 0.458106
## Sulphates        -2.918e-04  5.872e-03  -0.050 0.960364
## Alcohol          6.982e-03  1.434e-03   4.869 1.12e-06 ***
## LabelAppeal      2.305e-01  6.370e-03  36.189 < 2e-16 ***
## AcidIndex        -1.671e-02  4.887e-03  -3.420 0.000626 ***
## STARS            1.282e-01  6.413e-03  19.995 < 2e-16 ***
##
## Zero-inflation model coefficients (binomial with logit link):
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.9047850  1.2126749  -1.571 0.116245
## FixedAcidity    0.0039855  0.0050301   0.792 0.428176
## VolatileAcidity  0.2149929  0.0395880   5.431 5.61e-08 ***
## CitricAcid     -0.0570646  0.0366172  -1.558 0.119136
## ResidualSugar  -0.0006423  0.0009427  -0.681 0.495631
## Chlorides       0.2528988  0.1009827   2.504 0.012267 *
## FreeSulfurDioxide -0.0008383  0.0002185  -3.837 0.000125 ***
## TotalSulfurDioxide -0.0008014  0.0001397  -5.736 9.70e-09 ***
## Density         0.8825322  1.1936681   0.739 0.459698
## pH              0.1910630  0.0466618   4.095 4.23e-05 ***
## Sulphates       0.1525424  0.0355848   4.287 1.81e-05 ***
## Alcohol         0.0203647  0.0085429   2.384 0.017135 *
## LabelAppeal     0.6812879  0.0390907  17.428 < 2e-16 ***
## AcidIndex       0.4359032  0.0235896  18.479 < 2e-16 ***
## STARS           -3.4137368  0.0969201 -35.222 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Number of iterations in BFGS optimization: 36
## Log-likelihood: -2.081e+04 on 30 Df
```

Model Selection

To aid in model selection, let's test each of our models against the holdout validation set.

```
#train_data <- train_data[-split,]

getMAE <- function(x) {
  mean(abs(train_data$TARGET - x))
}
getRMSE <- function(x) {
  sqrt(mean((train_data$TARGET - x)^2))
}
results <- data.frame(Model = c("Poisson1",
                                "Poisson2",
```

```

      "Baseline-Model",
      "quasi-Poisson-Model",
      "Negative-binomial-regression",
      "Zero-Inflated-Count-Models"),
MAE = c(getMAE(predict.lm(Model1, train_data)),
        getMAE(predict.lm(Model2, train_data)),
        getMAE(predict.lm(Model3, train_data)),
        getMAE(predict.lm(Model4, train_data)),
        getMAE(predict(Model5, train_data,
                        type="response")),
        getMAE(predict(Model6, train_data,
                        type="response"))
      ),
RMSE = c(getRMSE(predict.lm(Model1, train_data)),
        getRMSE(predict.lm(Model2, train_data)),
        getRMSE(predict.lm(Model3, train_data)),
        getRMSE(predict.lm(Model4, train_data)),
        getRMSE(predict(Model5, train_data,
                        type="response")),
        getRMSE(predict(Model6, train_data,
                        type="response"))
      )
)
knitr::kable(results)

```

Model	MAE	RMSE
Poisson1	NA	NA
Poisson2	NA	NA
Baseline-Model	NA	NA
quasi-Poisson-Model	NA	NA
Negative-binomial-regression	NA	NA
Zero-Inflated-Count-Models	NA	NA

Here we see a preference for our full zero-counts Model. We'll make our predictions using our the entire model.

```
summary(Model6)
```

```

##
## Call:
## zeroinfl(formula = TARGET ~ ., data = complete_train_data)
##
## Pearson residuals:
##      Min      1Q   Median      3Q      Max
## -1.983021 -0.453229  0.001259  0.403750  6.459031
##
## Count model coefficients (poisson with log link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    1.370e+00  2.025e-01   6.767 1.31e-11 ***
## FixedAcidity    2.006e-04  8.416e-04   0.238 0.811564
## VolatileAcidity -1.158e-02  6.736e-03  -1.719 0.085681 .
## CitricAcid      7.397e-04  6.038e-03   0.122 0.902506
## ResidualSugar  -8.595e-05  1.582e-04  -0.543 0.587026

```

```

## Chlorides          -2.231e-02  1.687e-02  -1.322  0.186025
## FreeSulfurDioxide  3.138e-05  3.528e-05   0.890  0.373635
## TotalSulfurDioxide -2.307e-05  2.254e-05  -1.024  0.305910
## Density            -3.042e-01  1.986e-01  -1.531  0.125692
## pH                 5.824e-03  7.850e-03   0.742  0.458106
## Sulphates          -2.918e-04  5.872e-03  -0.050  0.960364
## Alcohol             6.982e-03  1.434e-03   4.869  1.12e-06 ***
## LabelAppeal        2.305e-01  6.370e-03  36.189  < 2e-16 ***
## AcidIndex          -1.671e-02  4.887e-03  -3.420  0.000626 ***
## STARS              1.282e-01  6.413e-03  19.995  < 2e-16 ***
##
## Zero-inflation model coefficients (binomial with logit link):
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.9047850  1.2126749  -1.571  0.116245
## FixedAcidity    0.0039855  0.0050301   0.792  0.428176
## VolatileAcidity  0.2149929  0.0395880   5.431  5.61e-08 ***
## CitricAcid     -0.0570646  0.0366172  -1.558  0.119136
## ResidualSugar  -0.0006423  0.0009427  -0.681  0.495631
## Chlorides       0.2528988  0.1009827   2.504  0.012267 *
## FreeSulfurDioxide -0.0008383  0.0002185  -3.837  0.000125 ***
## TotalSulfurDioxide -0.0008014  0.0001397  -5.736  9.70e-09 ***
## Density         0.8825322  1.1936681   0.739  0.459698
## pH              0.1910630  0.0466618   4.095  4.23e-05 ***
## Sulphates       0.1525424  0.0355848   4.287  1.81e-05 ***
## Alcohol         0.0203647  0.0085429   2.384  0.017135 *
## LabelAppeal     0.6812879  0.0390907  17.428  < 2e-16 ***
## AcidIndex       0.4359032  0.0235896  18.479  < 2e-16 ***
## STARS           -3.4137368  0.0969201 -35.222  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Number of iterations in BFGS optimization: 36
## Log-likelihood: -2.081e+04 on 30 Df

```

Appendix