

Predicting Wine Quality

Sandra & Chris & Evan

12/8/2022

Contents

Background	4
Analysis	4
Visualization	4
Correlations	4
Input Variable versus Output Variable	5
Base Model	5
Transformations	5
Model Selection	7
CV and Test Set Assessment	7
Prediction Values	7
Prediction Results	7
Conclusion	8

Executive Summary:

We aimed to predict the quality of 4547 red and white vinho verde wine samples, from the north of Portugal. Using a data set of 1950 wines, we analysed the data to determine the best subset of predictors and fit them to several different models. Using RMSE as a measure of model performance, we chose a Lasso model to make the final predictions.

Background

Background and Goals:

We received data on 1464 red and white vinho verde wine samples, from the north of Portugal. We have been given 13 different variables to help predict quality of wine (0-10 scale). Note that although the scale of wine is 0-10, all observations have a scale from 3-9.

Our goal is to build a model that will accurately predict wine quality from an additional 4547 red and white vinho verde wine samples, from the north of Portugal.

Analysis

We began by splitting the data given to us into a training and test set, this resampling of data was done randomly. The training set has 75% of our data and the test set has the remaining 25%. We will be trying to create a model based on the 75% of the data to predict the final 25%. Once we have a plausible model that accurately predicts the remaining 25% of our data, we will fit the model on the entire data set and predict wine quality of the remaining 4547 wine samples that we have not touched.

Visualization

To explore our data, we started with visualization. This gives us some understanding of the relationships existing within all the variables. We will observe characteristics that will provide valuable information, as well as checking if there exists some non-linear relationships. Given that they exists, the appropriate transformations will be applied to the respective variable.

Correlations

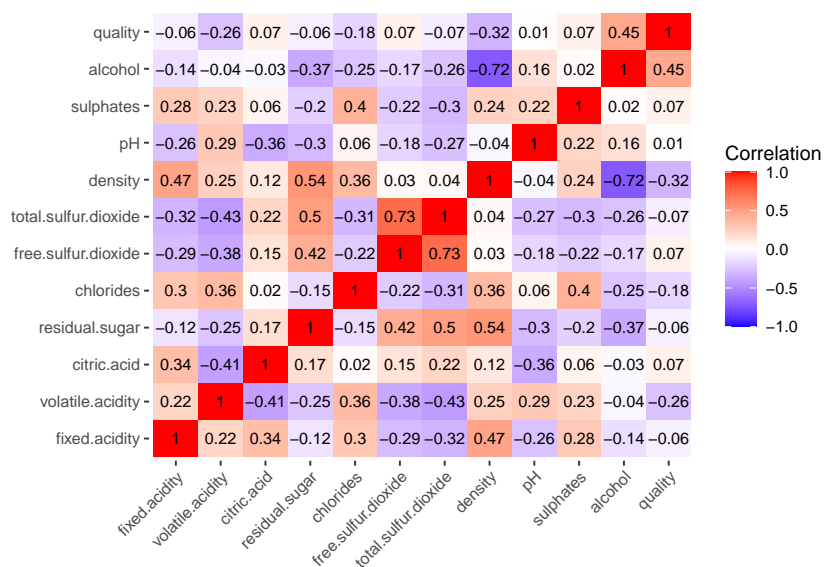


Figure 1: Correlation heatmap.

Figure 1 displays the correlations between each pairwise combination of numeric variables. We observe variables having multiple notable relationships with other variables; **residual sugar**, **free sulfur dioxide**,

total sulfur dioxide, density, and alcohol. For our response variable, `quality` is most related to `alcohol`. These relationships are kept in mind as we proceed forward in our data analysis.

Input Variable versus Output Variable

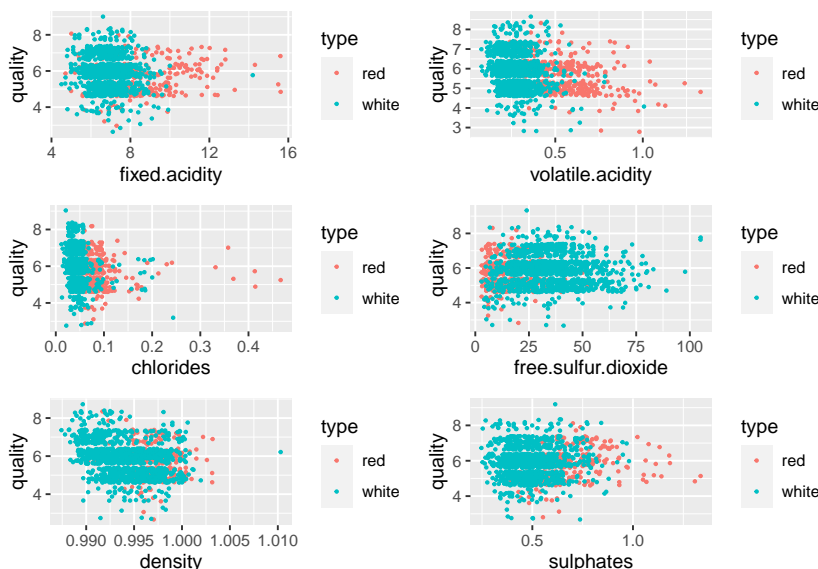


Figure 2: Variables with non-linear relationships.

Scatter plots for all the input variables against `quality` were visualized. Figure 2 includes plots for variables we suspect might have a possible non-linear relationship. Although we suspect a possible fault in our observation, we deemed a quadratic relationship to be fairly or most present within these selected variables.

Base Model

A linear model was trained using all input variables, with no transformation, to seek a CV RMSE value that will be referred to as the error we are trying to improve. The error we will attempt to improve is **0.7196**.

Viewing the plots with variables we suspected had non-linear relationships, we observe some pattern. Figure 3 includes the plots we refer to in this section. Although, we still caution that our observations may not be accurate. Now, we will proceed implementing the transformations from our initial assumptions.

Transformations

Following the base model, we included the addition of quadratic variables that were previously mentioned. Using this model with transformed variables, we plotted the residuals to observe any leftover patterns as well as interactions. The plots for each variables were separated by the observation's wine `type` to show any distinction between them that can provide an indication of interaction.

The scatter plots in Figure 4 only show a few of the plots we observed have some evidence for interaction. However, we had the same observation for all the other variables. Thus, we detected an interaction for all input variables with wine `type`.

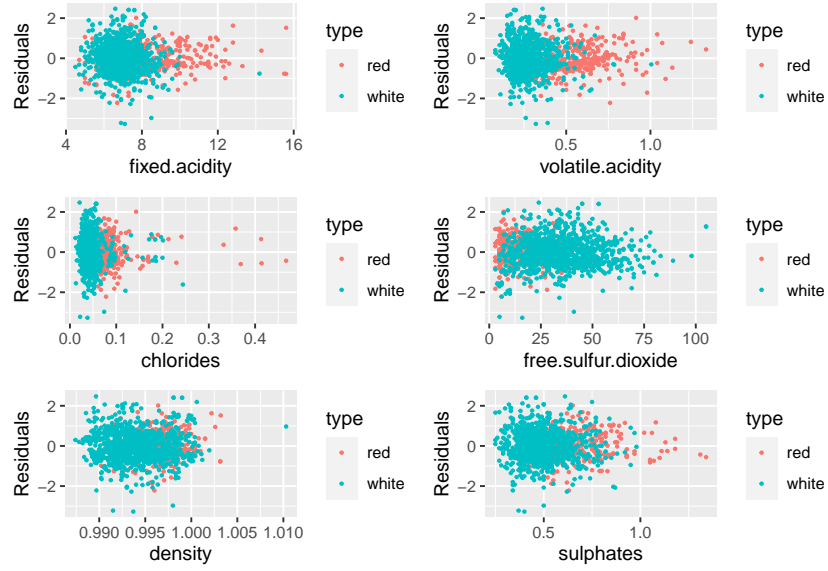


Figure 3: Select residual plots from base model.

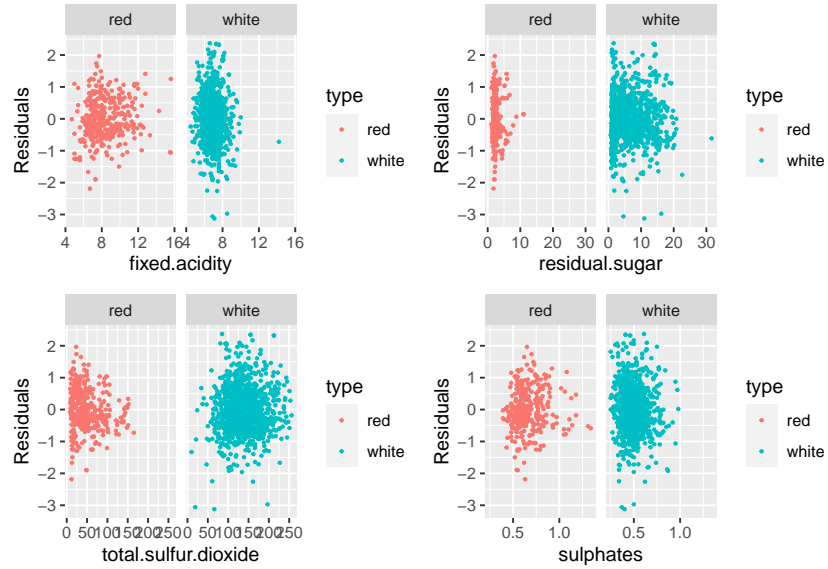


Figure 4: Select residual plots from model with transformed variables, colored by wine type.

Model Selection

We trained five different models using 3-fold cross validation.

- Model 1: LASSO- all variables.
- Model 2: LASSO- all variables and all possible interactions.
- Model 3: LASSO- all variables interacting wine `type`.
- Model 4: LASSO- all variables and transformed variables interacting wine `type`.
- Model 5: Linear Model- select variables; `density` and `alcohol` interacting wine `type`.

CV and Test Set Assessment

Table 1: RMSE Values

Model	CV	Test
1	0.7151	0.7189
2	0.6972	0.7148
3	0.7015	0.7182
4	0.7040	0.7153
5	0.7427	0.7421

Prediction Values

Once we finished testing out different models, we decided to go with the Lasso model which includes all variables and transformed variables interacting wine `type` to make predictions. The reason we decided to go with this model is because it is comparable to all other models, but it took out a lot of variables. Therefore it would be slightly easier to interpret. Once we decided to go with this model, we trained our entire data, the training and test set combined (with ID), using this model. We then used this model to predict the additional 4547 wines on quality. The table below shows the first 6 wines, with their predicted quality.

Our next step is to use the model on the entire data set, and then predict the quality of wine on the 4547 wines.

Prediction Results

```
## Warning in nominalTrainWorkflow(x = x, y = y, wts = weights, info = trainInfo, :  
## There were missing values in resampled performance measures.
```

```
##      [,1] [,2]  
## [1,] 1951  5.1  
## [2,] 1952  5.8  
## [3,] 1953  5.1  
## [4,] 1954  5.1  
## [5,] 1955  5.0  
## [6,] 1956  5.2
```

Conclusion

What we discovered in this project, is that predicting wine is difficult to do. Quality of wine is subject to change. Some people may like a specific wine, and others may not. It was interesting to see how using a range of 3 variables up to 30 variables would give us the roughly the same RMSE. It was clear however, that type had an enormous impact. The interactions with type and most variables made that clear. The average wine quality on the date we were given was 5.4 and the average wine we predicted was 5.7. We feel happy with our model of choice, knowing that our mean was close, and every model was giving us a similar RMSE.