

Udacity DAND Project: ExploreAndSummarizeData

Jubin Soni

August 27, 2017

Loading all the packages

```
library(ggplot2)
library(grid)
library(gridExtra)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following object is masked from 'package:gridExtra':
##
##      combine
```

```
## The following objects are masked from 'package:stats':
##
##      filter, lag
```

```
## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union
```

```
library(tidyr)
```

Loading the Dataset

```
wine <- read.csv("wineQualityWhites.csv", row.names = NULL)
```

Initial Data Exploration

```
names(wine)
```

```
## [1] "X"                "fixed.acidity"      "volatile.acidity"
## [4] "citric.acid"       "residual.sugar"     "chlorides"
## [7] "free.sulfur.dioxide" "total.sulfur.dioxide" "density"
## [10] "pH"               "sulphates"          "alcohol"
## [13] "quality"
```

```
str(wine)
```

```
## 'data.frame': 4898 obs. of 13 variables:
## $ X : int 1 2 3 4 5 6 7 8 9 10 ...
## $ fixed.acidity : num 7 6.3 8.1 7.2 7.2 8.1 6.2 7 6.3 8.1 ...
## $ volatile.acidity : num 0.27 0.3 0.28 0.23 0.23 0.28 0.32 0.27 0.3 0.22 ...
## $ citric.acid : num 0.36 0.34 0.4 0.32 0.32 0.4 0.16 0.36 0.34 0.43 ...
## $ residual.sugar : num 20.7 1.6 6.9 8.5 8.5 6.9 7 20.7 1.6 1.5 ...
## $ chlorides : num 0.045 0.049 0.05 0.058 0.058 0.05 0.045 0.045 0.049 0.044 ...
## $ free.sulfur.dioxide : num 45 14 30 47 47 30 30 45 14 28 ...
## $ total.sulfur.dioxide: num 170 132 97 186 186 97 136 170 132 129 ...
## $ density : num 1.001 0.994 0.995 0.996 0.996 ...
## $ pH : num 3 3.3 3.26 3.19 3.19 3.26 3.18 3 3.3 3.22 ...
## $ sulphates : num 0.45 0.49 0.44 0.4 0.4 0.44 0.47 0.45 0.49 0.45 ...
## $ alcohol : num 8.8 9.5 10.1 9.9 9.9 10.1 9.6 8.8 9.5 11 ...
## $ quality : int 6 6 6 6 6 6 6 6 6 6 ...
```

```
summary(wine)
```

```
##           X           fixed.acidity  volatile.acidity  citric.acid
## Min.      : 1   Min.      : 3.800   Min.      :0.0800   Min.      :0.0000
## 1st Qu.:1225   1st Qu.: 6.300   1st Qu.:0.2100   1st Qu.:0.2700
## Median :2450   Median : 6.800   Median :0.2600   Median :0.3200
## Mean      :2450   Mean      : 6.855   Mean      :0.2782   Mean      :0.3342
## 3rd Qu.:3674   3rd Qu.: 7.300   3rd Qu.:0.3200   3rd Qu.:0.3900
## Max.      :4898   Max.      :14.200   Max.      :1.1000   Max.      :1.6600
## residual.sugar  chlorides      free.sulfur.dioxide
## Min.      : 0.600   Min.      :0.00900   Min.      : 2.00
## 1st Qu.: 1.700   1st Qu.:0.03600   1st Qu.: 23.00
## Median : 5.200   Median :0.04300   Median : 34.00
## Mean      : 6.391   Mean      :0.04577   Mean      : 35.31
## 3rd Qu.: 9.900   3rd Qu.:0.05000   3rd Qu.: 46.00
## Max.      :65.800   Max.      :0.34600   Max.      :289.00
## total.sulfur.dioxide  density      pH      sulphates
## Min.      : 9.0      Min.      :0.9871   Min.      :2.720   Min.      :0.2200
## 1st Qu.:108.0      1st Qu.:0.9917   1st Qu.:3.090   1st Qu.:0.4100
## Median :134.0      Median :0.9937   Median :3.180   Median :0.4700
## Mean      :138.4      Mean      :0.9940   Mean      :3.188   Mean      :0.4898
## 3rd Qu.:167.0      3rd Qu.:0.9961   3rd Qu.:3.280   3rd Qu.:0.5500
## Max.      :440.0      Max.      :1.0390   Max.      :3.820   Max.      :1.0800
## alcohol      quality
## Min.      : 8.00   Min.      :3.000
## 1st Qu.: 9.50   1st Qu.:5.000
## Median :10.40   Median :6.000
## Mean      :10.51   Mean      :5.878
## 3rd Qu.:11.40   3rd Qu.:6.000
## Max.      :14.20   Max.      :9.000
```

Structure of Dataset

The dataset on white wines has 4898 observations of 13 variables. The main feature of interest is quality and alcohol content.

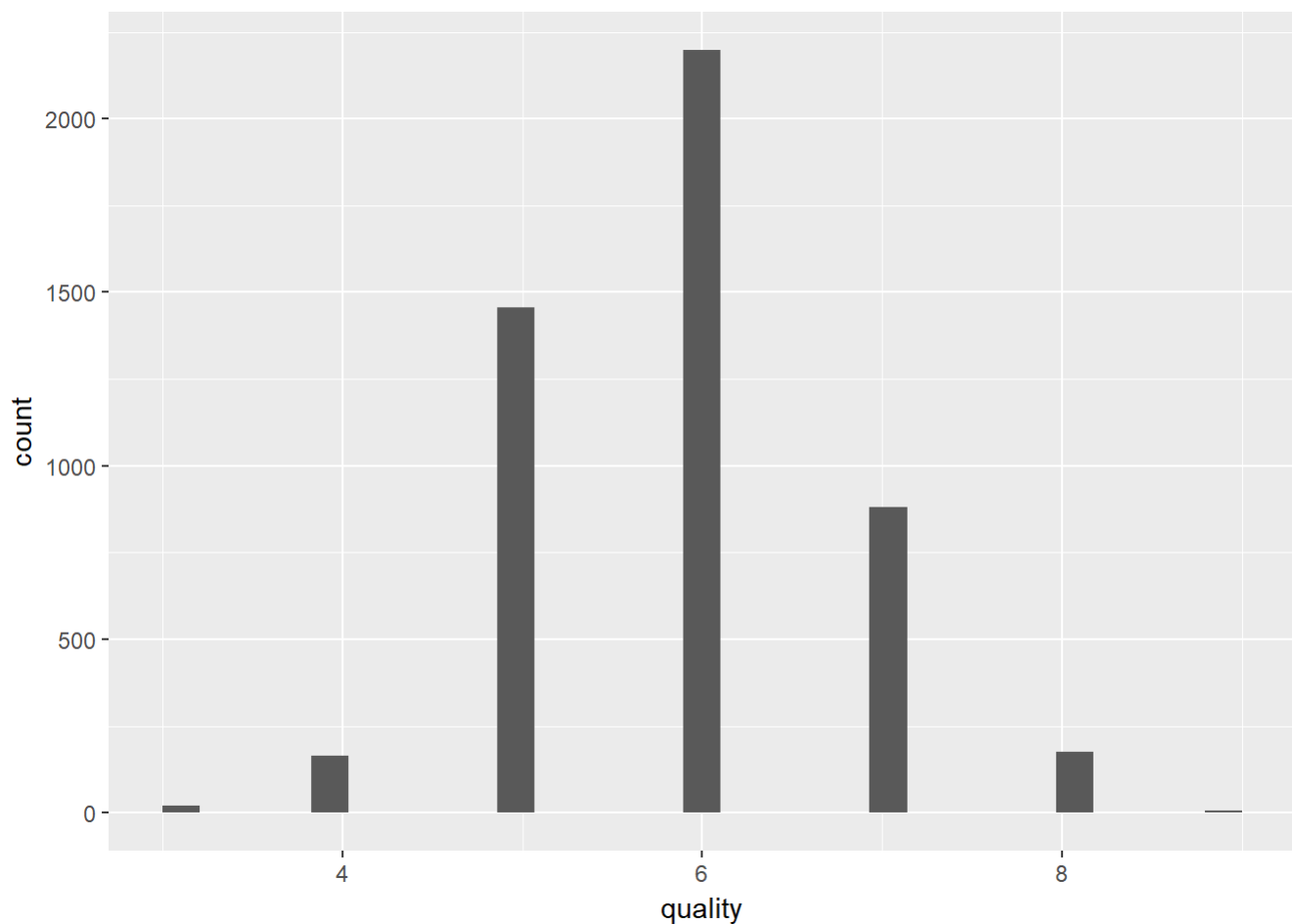
Univariate Analysis

Plotting quality of wine

We see that the quality is close to a normal distribution, the graph is which is not fully symmetric because the quality 5 is higher than its corresponding value. However looking at the quality of 6 being highest we can say that the quality is above average.

```
library(ggplot2)
ggplot(aes(x = quality), data = wine) +
  geom_histogram()
```

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

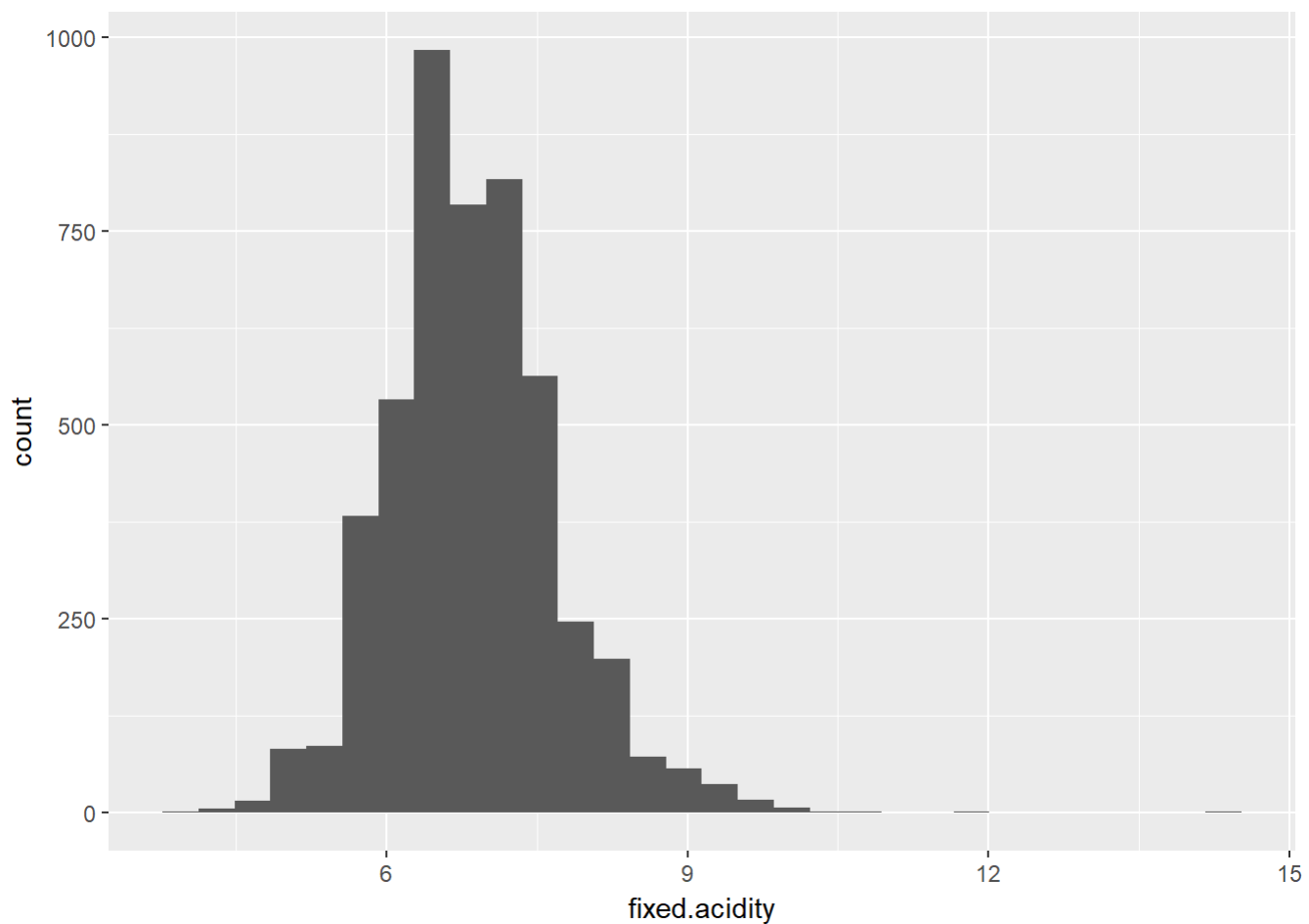


Histogram of fixed.acidity value of wine

The histogram of acidity of potential shows normal distribution.

```
ggplot(aes(x = fixed.acidity), data = wine) +  
  geom_histogram()
```

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

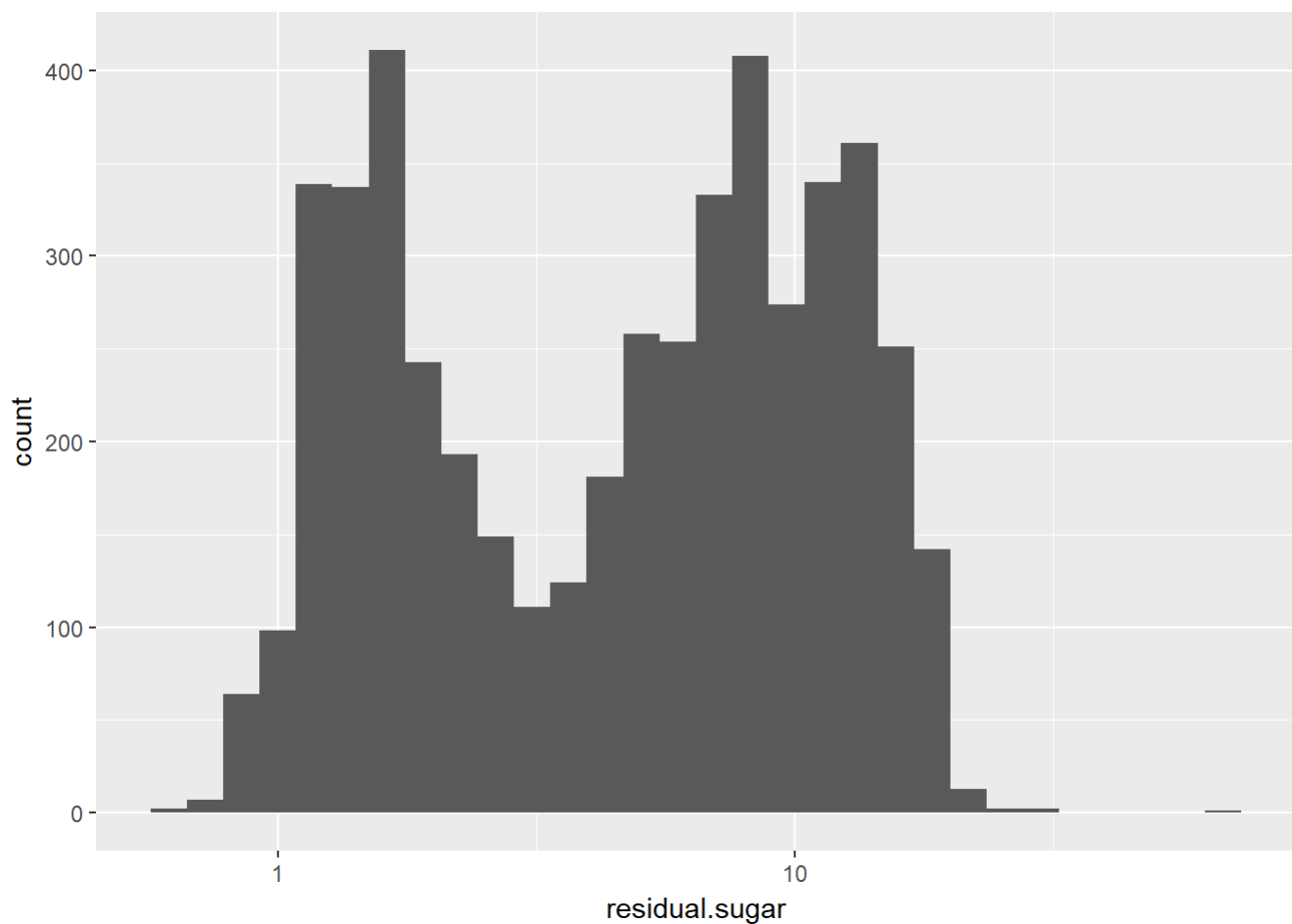


Histogram of residual sugar

The histogram of residual sugar shows that as the count of sugar in whitewine does not follow any standard distribution. It has low peaks at left and right and low peak in the middle of distribution.

```
ggplot(aes(x = residual.sugar), data = wine) +  
  geom_histogram() +  
  scale_x_log10()
```

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

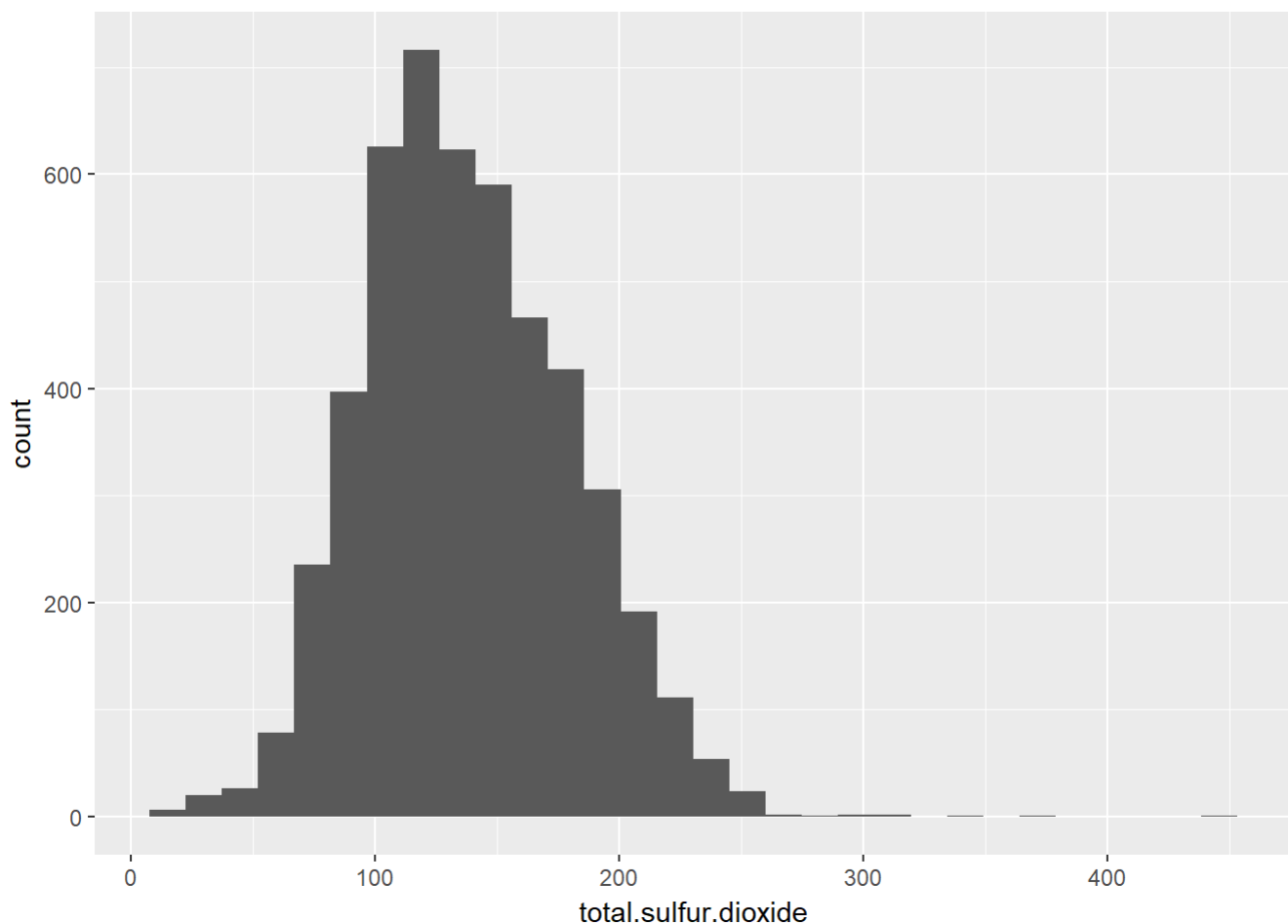


Histogram of Sulphur dioxide

The histogram of Sulfur Dioxide shows normal distribution in the White wine.

```
ggplot(aes(x = total.sulfur.dioxide), data = wine) +  
  geom_histogram()
```

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

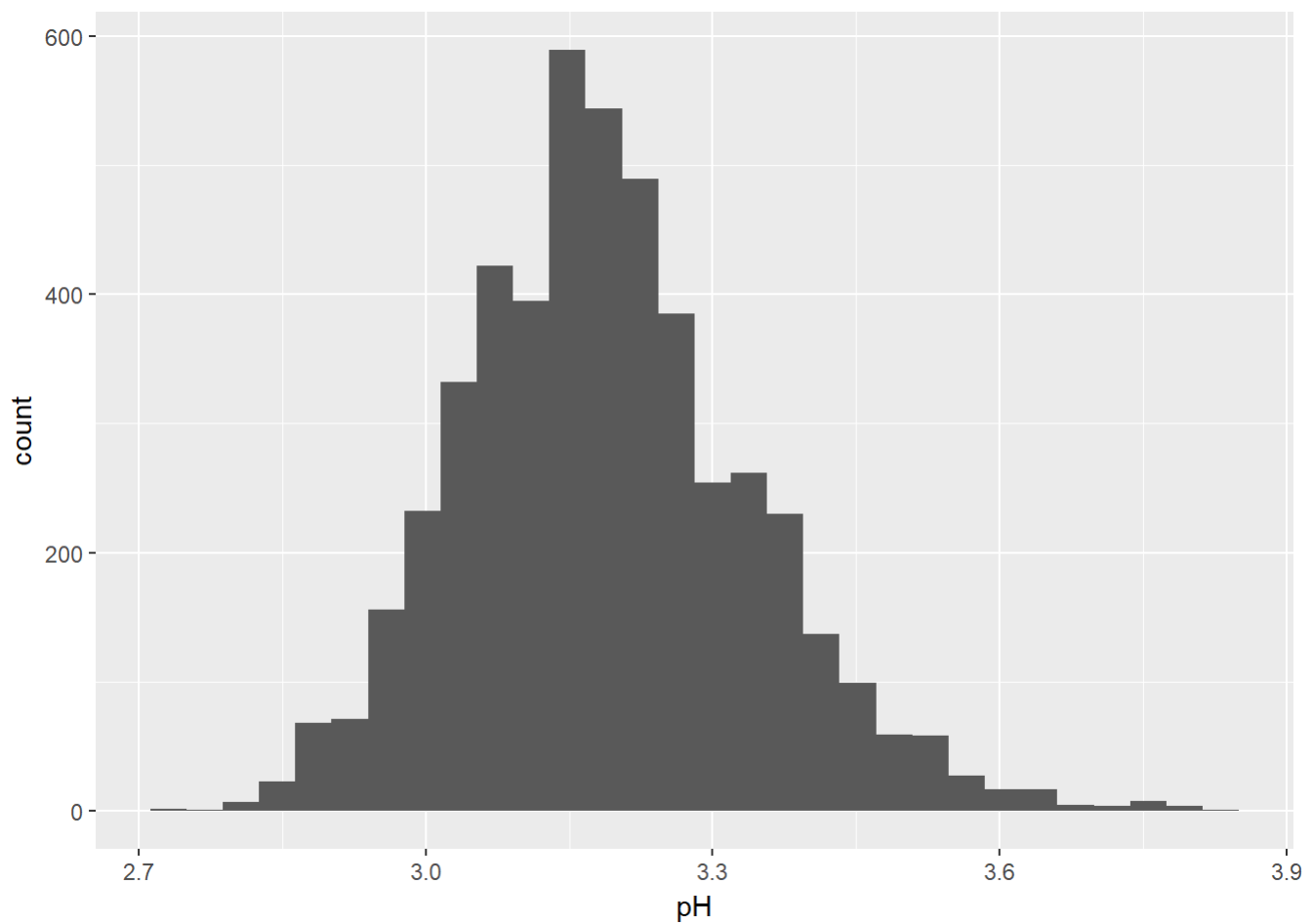


Plotting pH of wine

The pH is very narrow for White Wine between 2.7 to 3.8, so it seems like these values are within the acidity range and no alkaline (which is above 7)

```
ggplot(aes(x = pH), data = wine) +  
  geom_histogram()
```

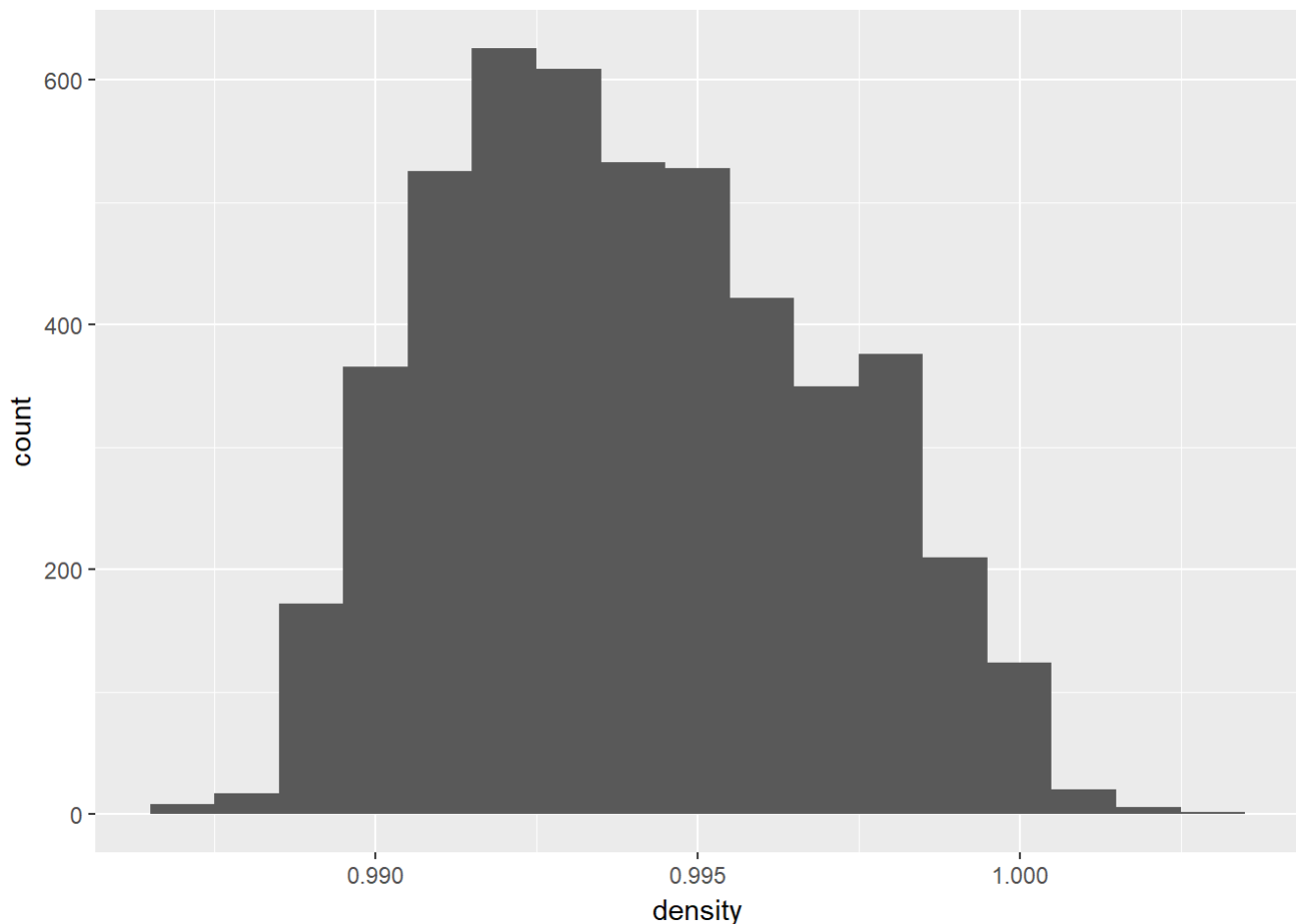
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



Checking Wine Density

The distribution below shows that density is within 0.9 to 1.0 range.

```
ggplot(aes(x = density), data = subset(wine, density<1.01)) +  
  geom_histogram(binwidth = 0.001)
```

Univariate Analysis Questions and Findings

What other features in the dataset do you think will help support your investigation into your feature(s) of interest?

I think further investigation whether pH could be classified into categorical (acid, base, neutral) could potentially be of interest. Also the relationship between fixed.acidity, volatile.acidity, and citric.acid (in particular this one due to the large number of 0s) could be interesting to further investigate in order to understand if the values are properly reported. This would require to understand the theoretical relationship between these variables.

Did you create any new variables from existing variables in the dataset?

I only created the variable (ordered factor) rating to classify wine as good, bad or average. The information that the other variables store there does not seem to be any other variable that would fit well a sub-classification into categorical variables, perhaps with the exception of residual pH (neutral, base, acid).

Of the features you investigated, were there any unusual distributions?

I created Boxplots in Bivariate analysis to gain a better understanding wrt. outliers (except for X for being an index, and quality for being categorical). I've also plotted the data using a log10 scale in Univariate analysis (I plotted for all variables, although only those with long tails are interesting for the log10 scale – plots below the box plots). These plots have shown that fixed acidity and to some extent pH, chlorides, density, sulphates, volatile acidity to follow a normal distribution. As for the acidity variables, this is aligned with the fact that pH seems to be normally distributed, apart from the citric acid. The reason for the latter might be the number of 0 (potentially non-responses) discussed earlier. pH is normally distribution which suggests that the data is good, since by definition it is a measure of acidity and is on a logarithmic scale.

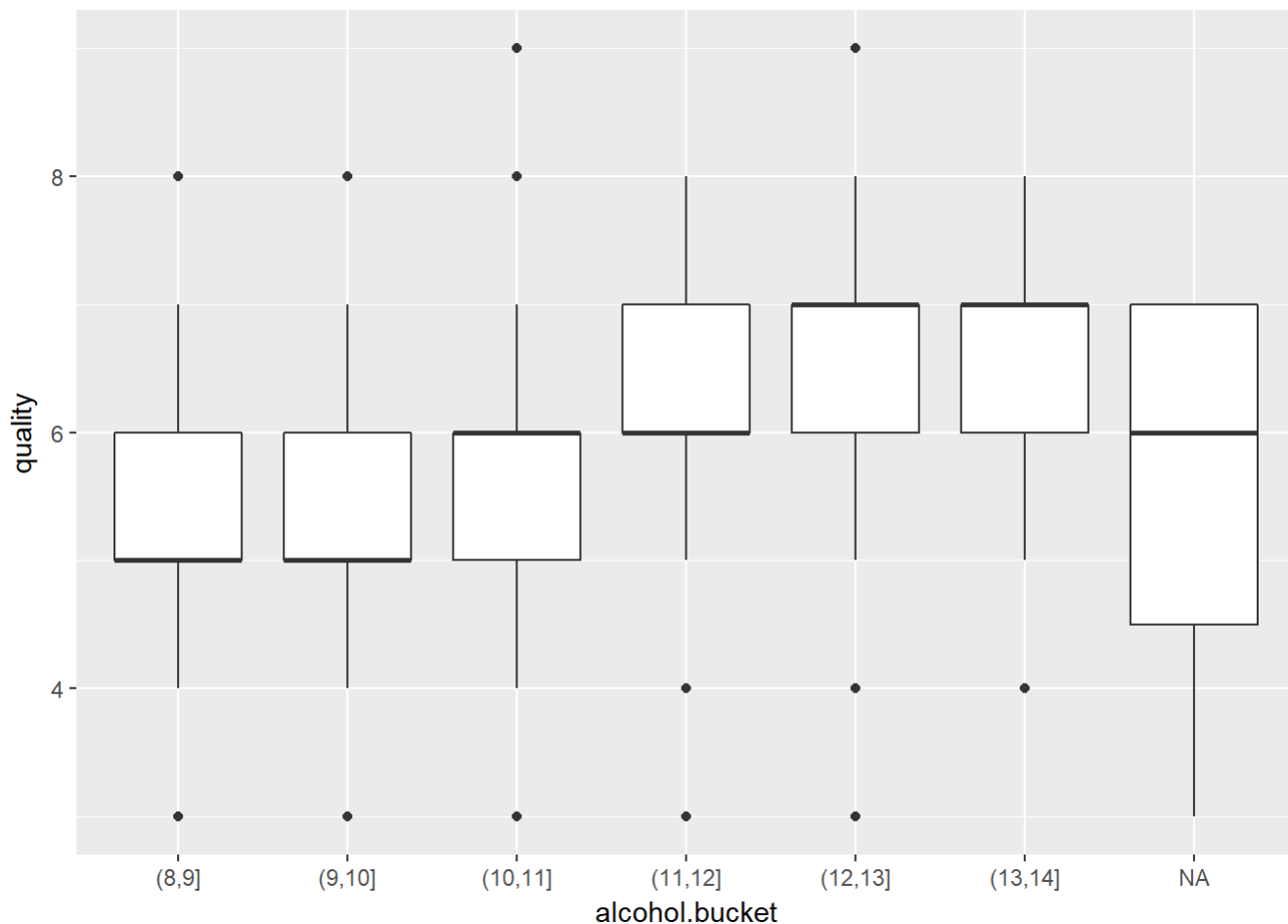
Did you perform any operations on the data to tidy, adjust, or change the form of the data? If so, why did you do this?

I've not found the need to tidy, adjust, or change the form of the data.

Bivariate Analysis

Let us check how is the distribution of quality and alcohol content in wine. To do this creating a new variable `alcohol.bucket`.

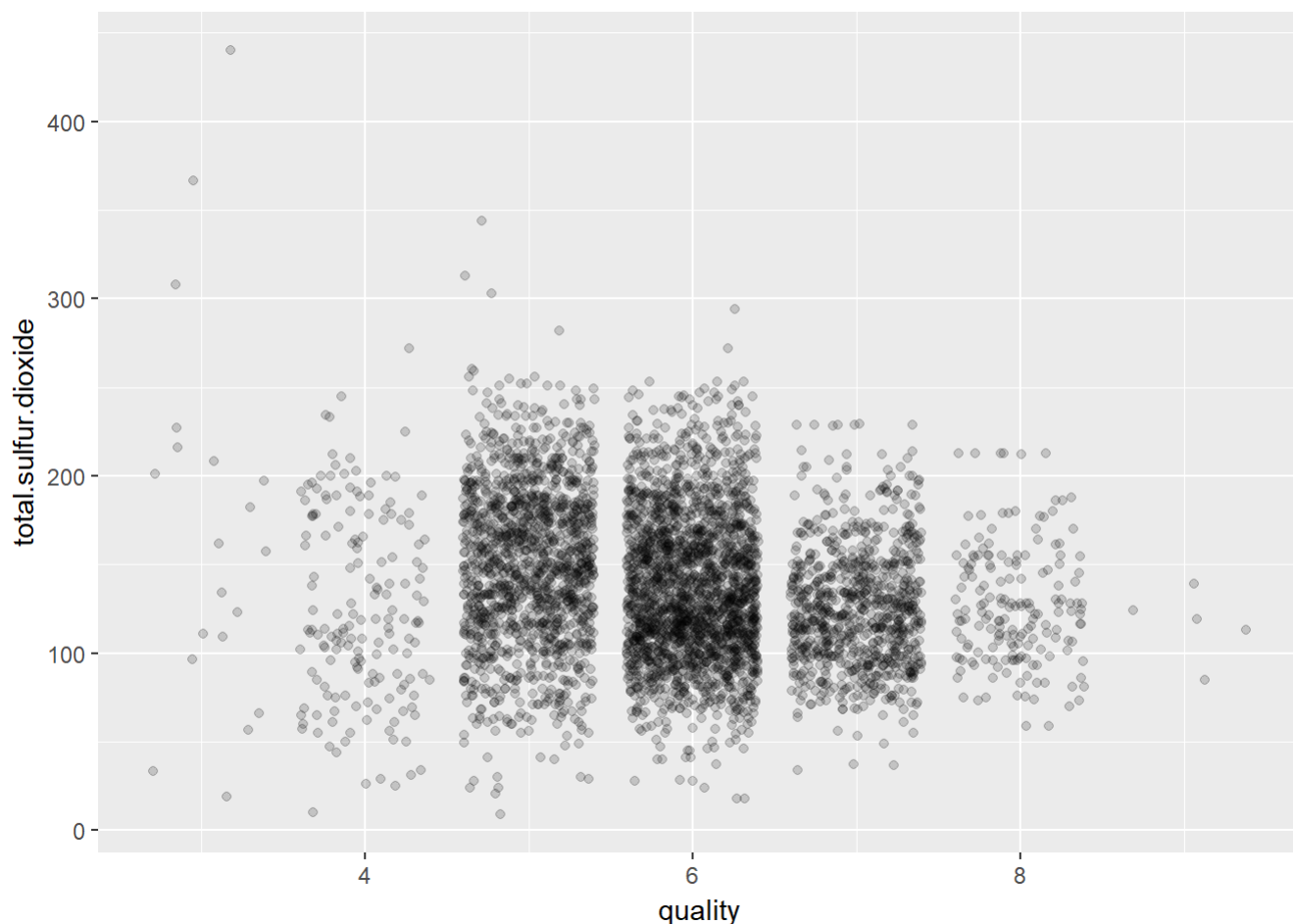
```
wine$alcohol.bucket <- cut(wine$alcohol, c(8,9,10,11,12,13,14))  
qplot(x = alcohol.bucket, y= quality, data = wine, geom = 'boxplot')
```



The graphs shows us that alcohol and wine quality have clear correlation.

Scatterplot of quality and Sulfur.Dioxide

```
ggplot(aes(x=quality, y=total.sulfur.dioxide), data=wine)+  
  geom_jitter(alpha=1/6)
```



This plot shows that in white wine quality of 5,6,7 has majority of Sulfur.Dioxide in it.

Talk about some of the relationships you observed in this part of the investigation. How did the feature(s) of interest vary with other features in the dataset?

I noticed that the following characteristics yield 'good' wines:

- Acidity: higher fixed acidity and citric acid; lower volative acidity
- Lower pH (~3.5)
- Higher Sulphates
- Higher Alcohol
- Lower density
- Regular chlorides and sugar do not seem to impact much on the quality.

Did you observe any interesting relationships between other features and not the main features of interest?

Below I have calculated the correlation between variables that correlate the most with quality: citric acid, sulphates, alcohol and volatile acidity.

Also the Pearson test revealed that volatile acidity and citric acid have strong negative correlation.

```
with(wine, cor.test(citric.acid, alcohol))
```

```
##
## Pearson's product-moment correlation
##
## data: citric.acid and alcohol
## t = -5.3141, df = 4896, p-value = 1.119e-07
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.1035156 -0.0478237
## sample estimates:
##          cor
## -0.07572873
```

```
with(wine, cor.test(sulphates, alcohol))
```

```
##
## Pearson's product-moment correlation
##
## data: sulphates and alcohol
## t = -1.22, df = 4896, p-value = 0.2225
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.04541705 0.01057885
## sample estimates:
##          cor
## -0.01743277
```

```
with(wine, cor.test(volatile.acidity, alcohol))
```

```
##
## Pearson's product-moment correlation
##
## data: volatile.acidity and alcohol
## t = 4.7492, df = 4896, p-value = 2.1e-06
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.03978694 0.09554320
## sample estimates:
##          cor
## 0.06771794
```

```
with(wine, cor.test(sulphates, citric.acid))
```

```
##
## Pearson's product-moment correlation
##
## data: sulphates and citric.acid
## t = 4.3699, df = 4896, p-value = 1.269e-05
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## 0.03438451 0.09017997
## sample estimates:
## cor
## 0.06233094
```

```
with(wine, cor.test(sulphates, volatile.acidity))
```

```
##
## Pearson's product-moment correlation
##
## data: sulphates and volatile.acidity
## t = -2.5015, df = 4896, p-value = 0.0124
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.063670893 -0.007729425
## sample estimates:
## cor
## -0.03572815
```

```
with(wine, cor.test(volatile.acidity, citric.acid))
```

```
##
## Pearson's product-moment correlation
##
## data: volatile.acidity and citric.acid
## t = -10.578, df = 4896, p-value < 2.2e-16
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
## -0.1767384 -0.1219760
## sample estimates:
## cor
## -0.1494718
```

What was the strongest relationship you found?

The quality shows the strongest correlation with alcohol (0.4761663). I also observed that volatile.acidity and alcohol, and sulphates and alcohol show weakest bi-variate relationship.

```
with(wine, cor.test(quality, alcohol))
```

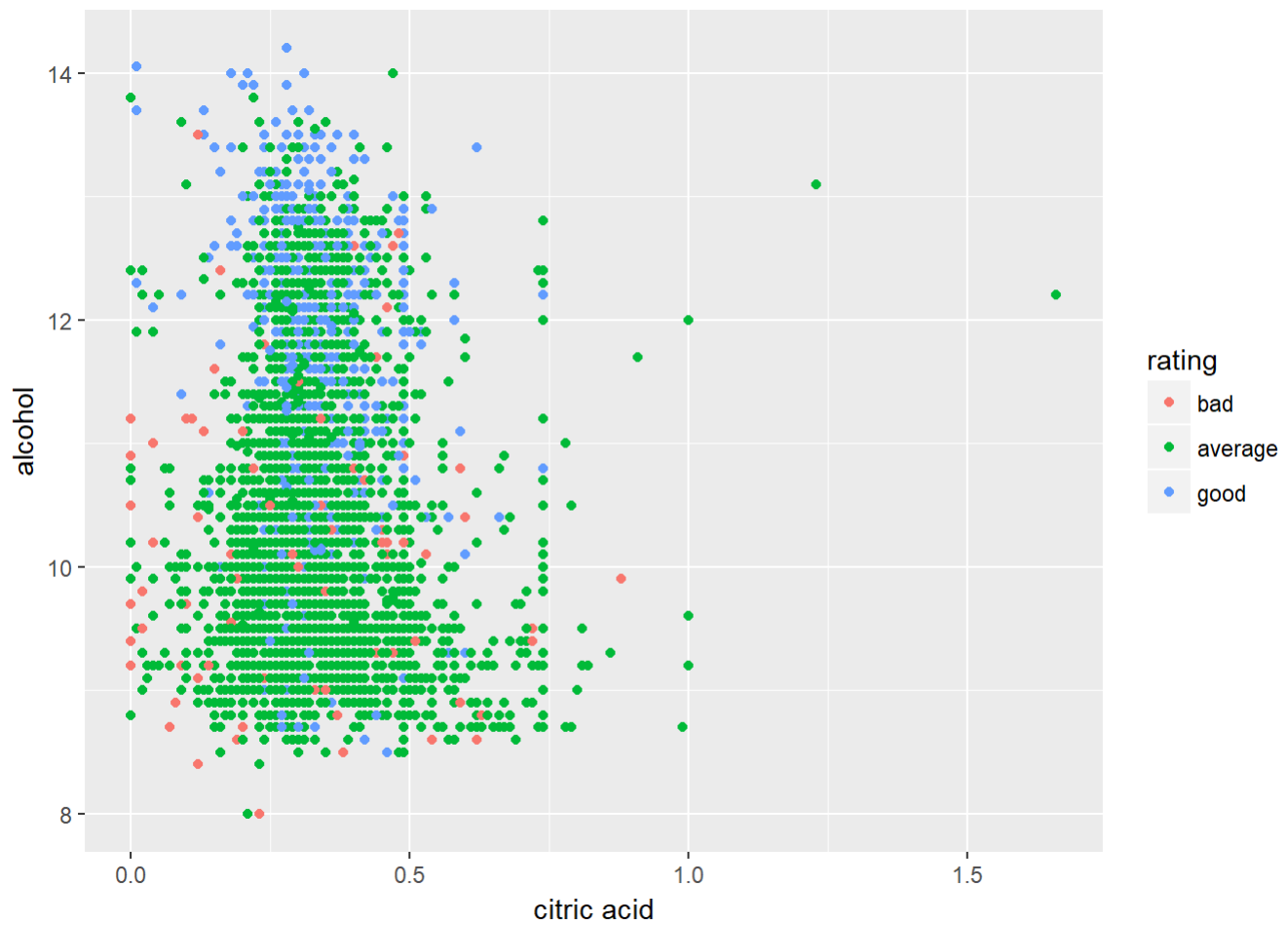
```
##  
## Pearson's product-moment correlation  
##  
## data: quality and alcohol  
## t = 33.858, df = 4896, p-value < 2.2e-16  
## alternative hypothesis: true correlation is not equal to 0  
## 95 percent confidence interval:  
## 0.4126015 0.4579941  
## sample estimates:  
## cor  
## 0.4355747
```

Multivariate Analysis

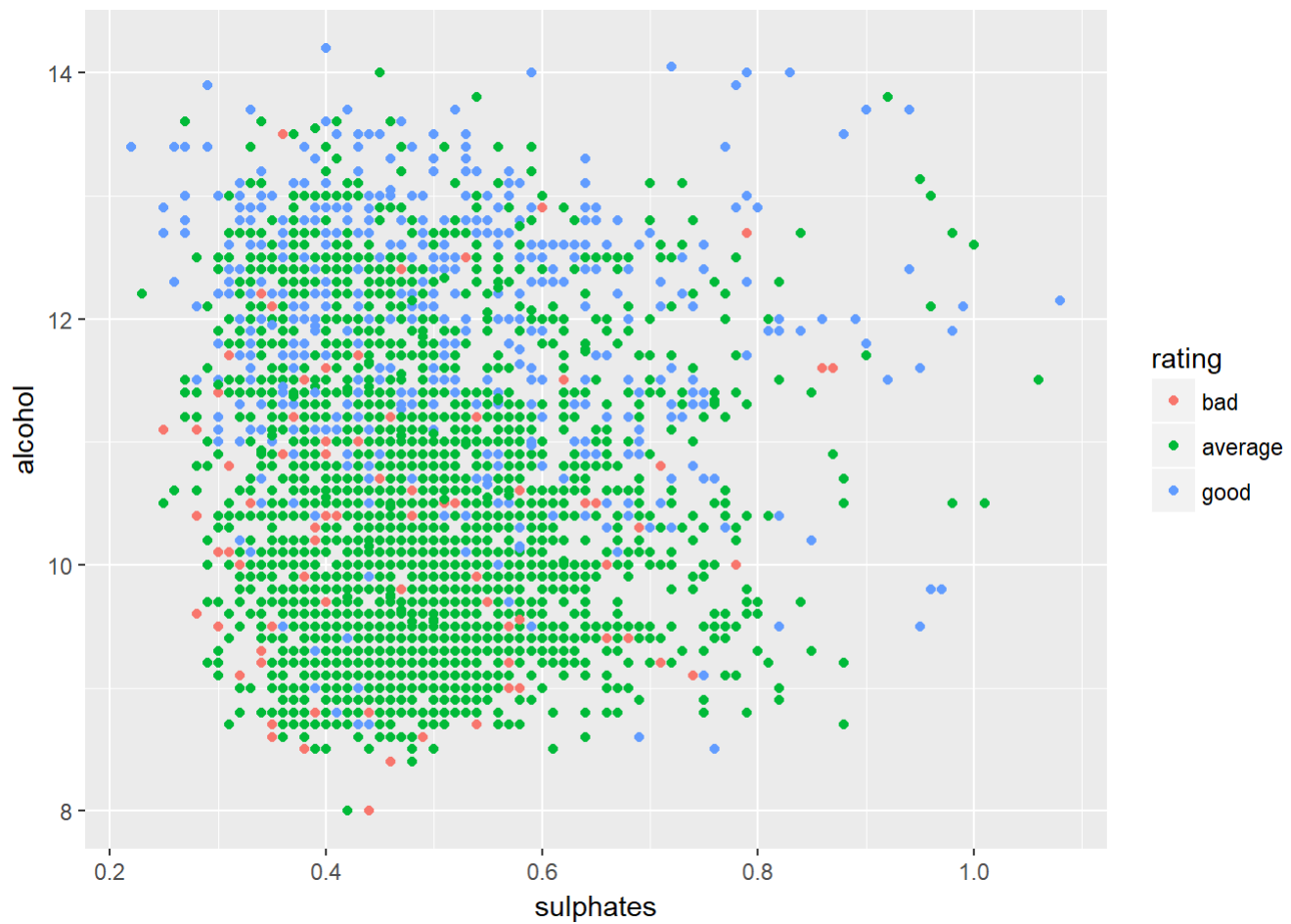
Talk about some of the relationships you observed in this investigation. Were there features that strengthened each other in terms of looking at your feature(s) of interest?

I focused on 4 features that showed strong correlation with quality and pH, I found that the key factors that help in achieving good wine quality are higher sulphates, higher alcohol, lower volatile acid and higher citric acid.

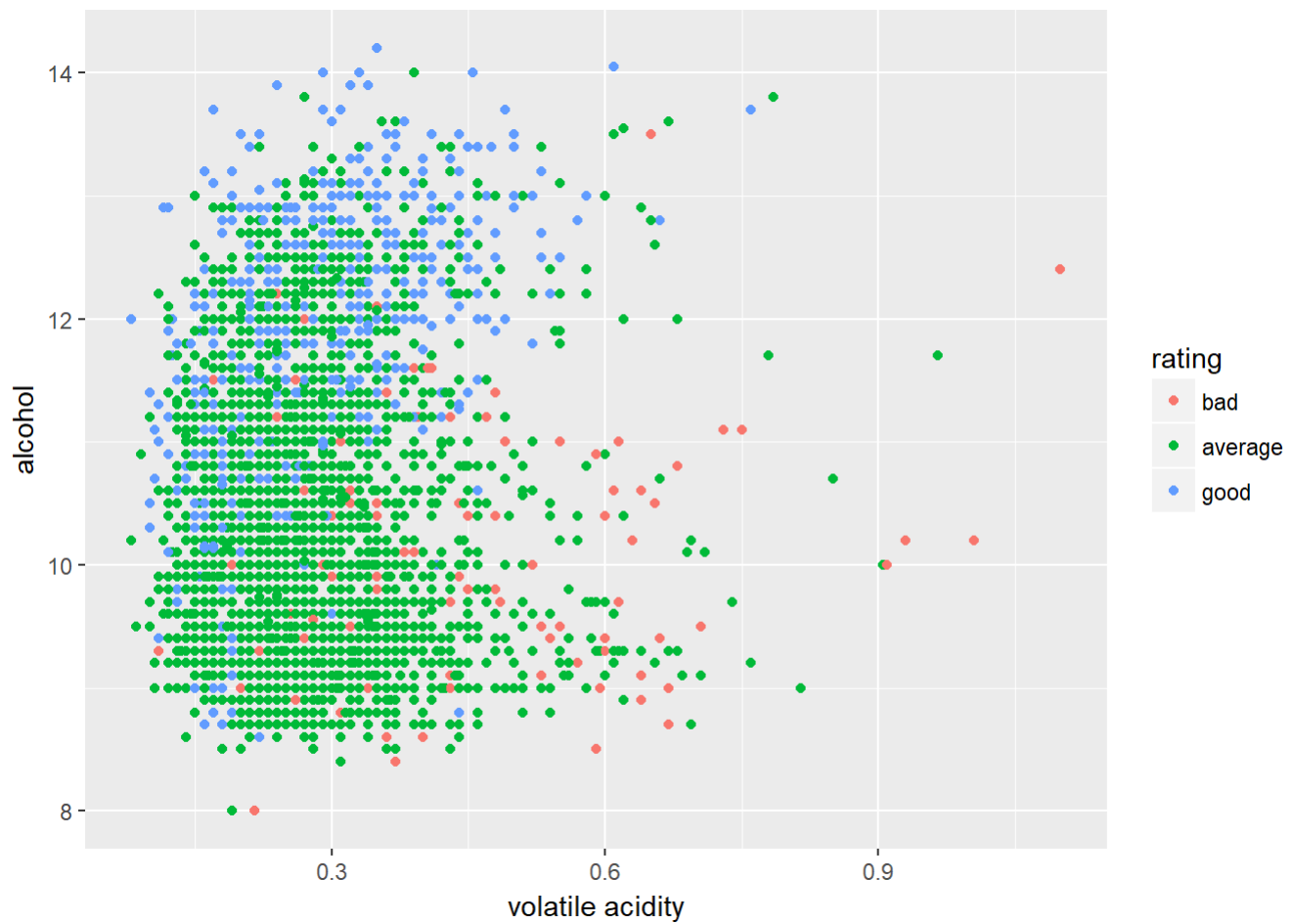
```
library(ggplot2)  
wine$rating <- ifelse(wine$quality < 5, 'bad',  
                     ifelse(wine$quality < 7, 'average', 'good'))  
wine$rating <- ordered(wine$rating,  
                      levels = c('bad', 'average', 'good'))  
  
plotRating <- function(x, y, xl, yl){  
  ggplot(data = wine, aes(x=x, y=y, color=rating))+  
    geom_point()+  
    xlab(xl)+  
    ylab(yl)  
}  
  
plotRating(wine$citric.acid, wine$alcohol, 'citric acid', 'alcohol')
```



```
plotRating(wine$sulphates, wine$alcohol, 'sulphates', 'alcohol')
```



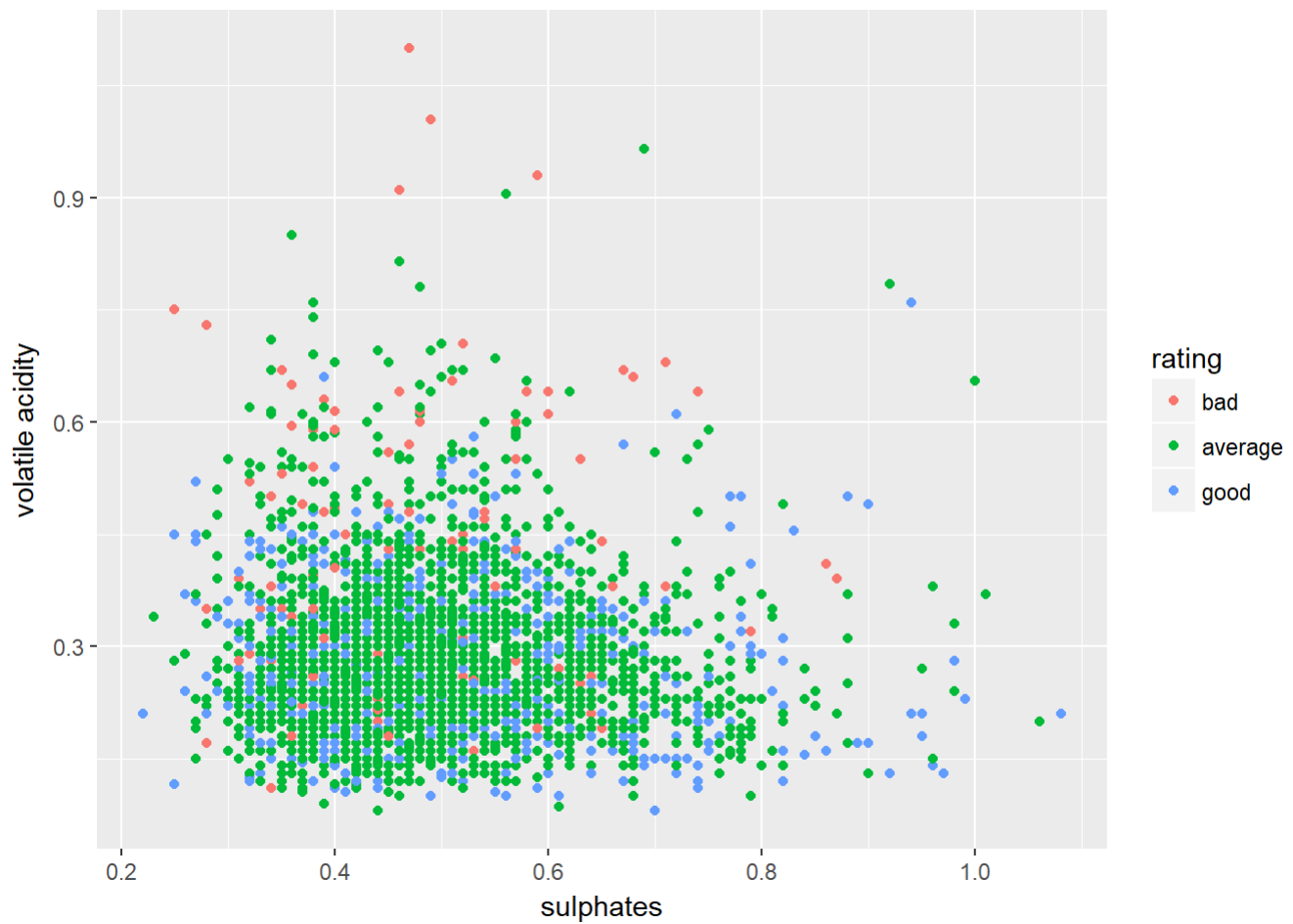
```
plotRating(wine$volatile.acidity, wine$alcohol,  
  'volatile acidity', 'alcohol')
```

```
plotRating(wine$sulphates, wine$citric.acid, 'sulphates',  
           'citric acid')
```



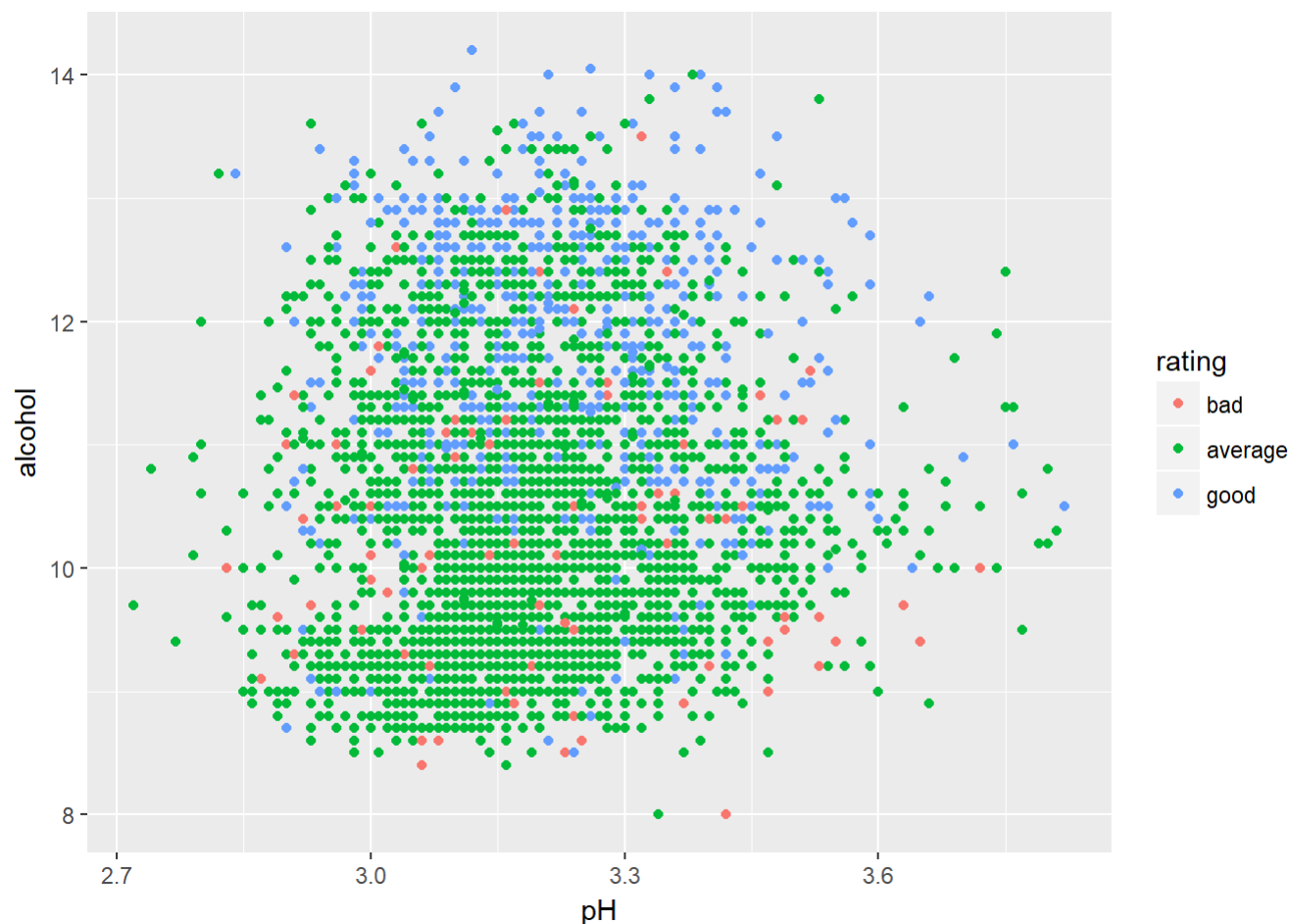
```
plotRating(wine$sulphates, wine$volatile.acidity, 'sulphates',  
           'volatile acidity')
```



```
plotRating(wine$volatile.acidity, wine$citric.acid, 'volatile acidity',  
           'citric acid')
```



```
plotRating(wine$pH, wine$alcohol, 'pH', 'alcohol')
```



Final Plots and Summary

Plot One

These plots show the effect of acidity and pH on wine quality. I observe that higher acidity (lower pH), apart from the volatile acid, is shown to yield better wines. Higher acidity and lower pH also makes sense because these variables are correlated and this is why acidity is inline with pH (making a balance for good wine).

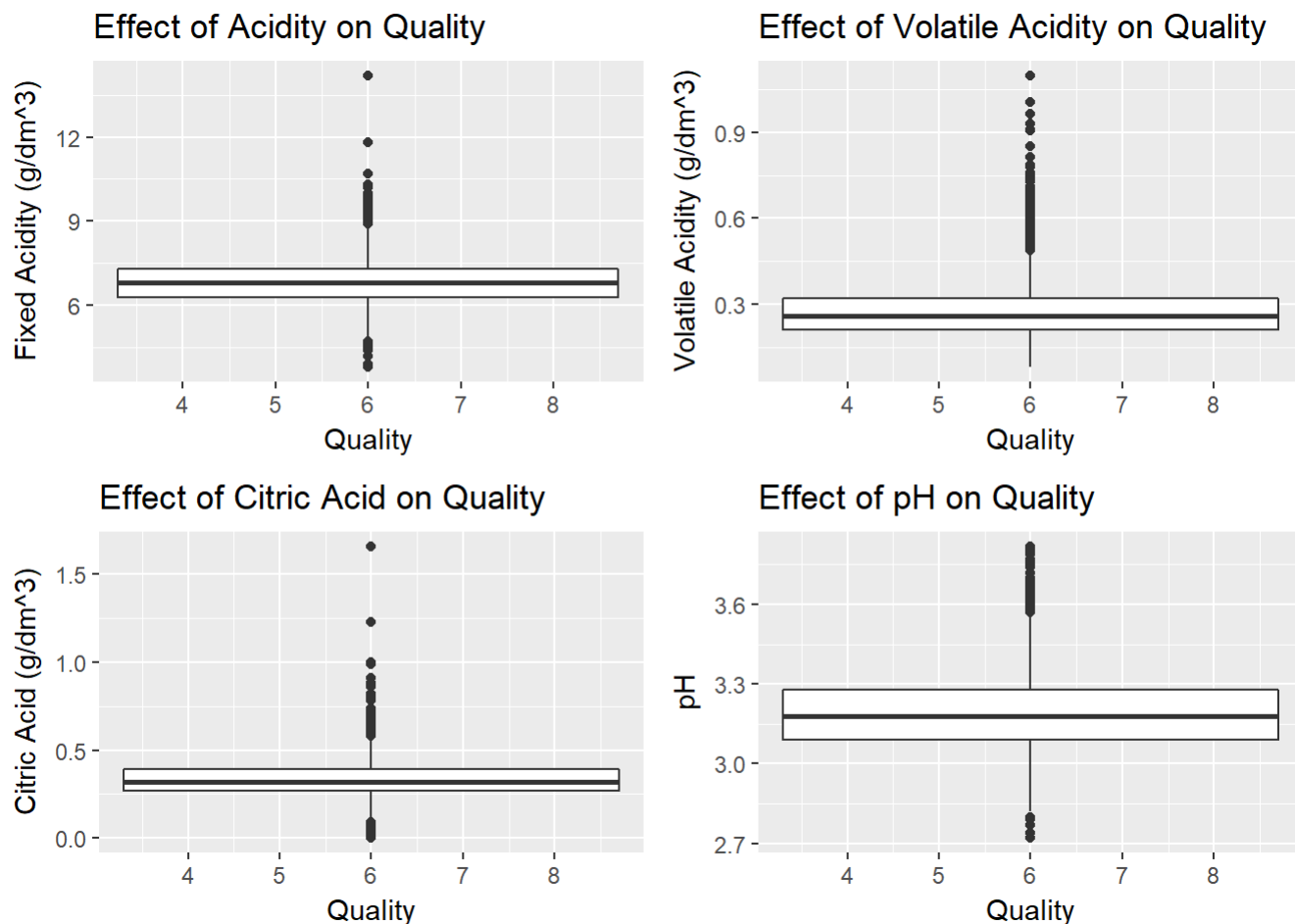
```
grid.arrange(ggplot(data=wine, aes(x=quality, y=fixed.acidity, fill=quality))+
  ggtitle('Effect of Acidity on Quality')+
  ylab('Fixed Acidity (g/dm^3)')+
  xlab('Quality')+
  geom_boxplot(),
  ggplot(data=wine, aes(x=quality, y=volatile.acidity, fill=quality))+
  ggtitle('Effect of Volatile Acidity on Quality')+
  ylab('Volatile Acidity (g/dm^3)')+
  xlab('Quality')+
  geom_boxplot(),
  ggplot(data=wine, aes(x=quality, y=citric.acid, fill=quality)) +
  ggtitle('Effect of Citric Acid on Quality')+
  ylab('Citric Acid (g/dm^3)')+
  xlab('Quality')+
  geom_boxplot(),
  ggplot(data=wine, aes(x=quality, y=pH, fill=quality))+
  ggtitle('Effect of pH on Quality')+
  ylab('pH')+
  xlab('Quality')+
  geom_boxplot()
)
```

```
## Warning: Continuous x aesthetic -- did you forget aes(group=...)?
```

```
## Warning: Continuous x aesthetic -- did you forget aes(group=...)?
```

```
## Warning: Continuous x aesthetic -- did you forget aes(group=...)?
```

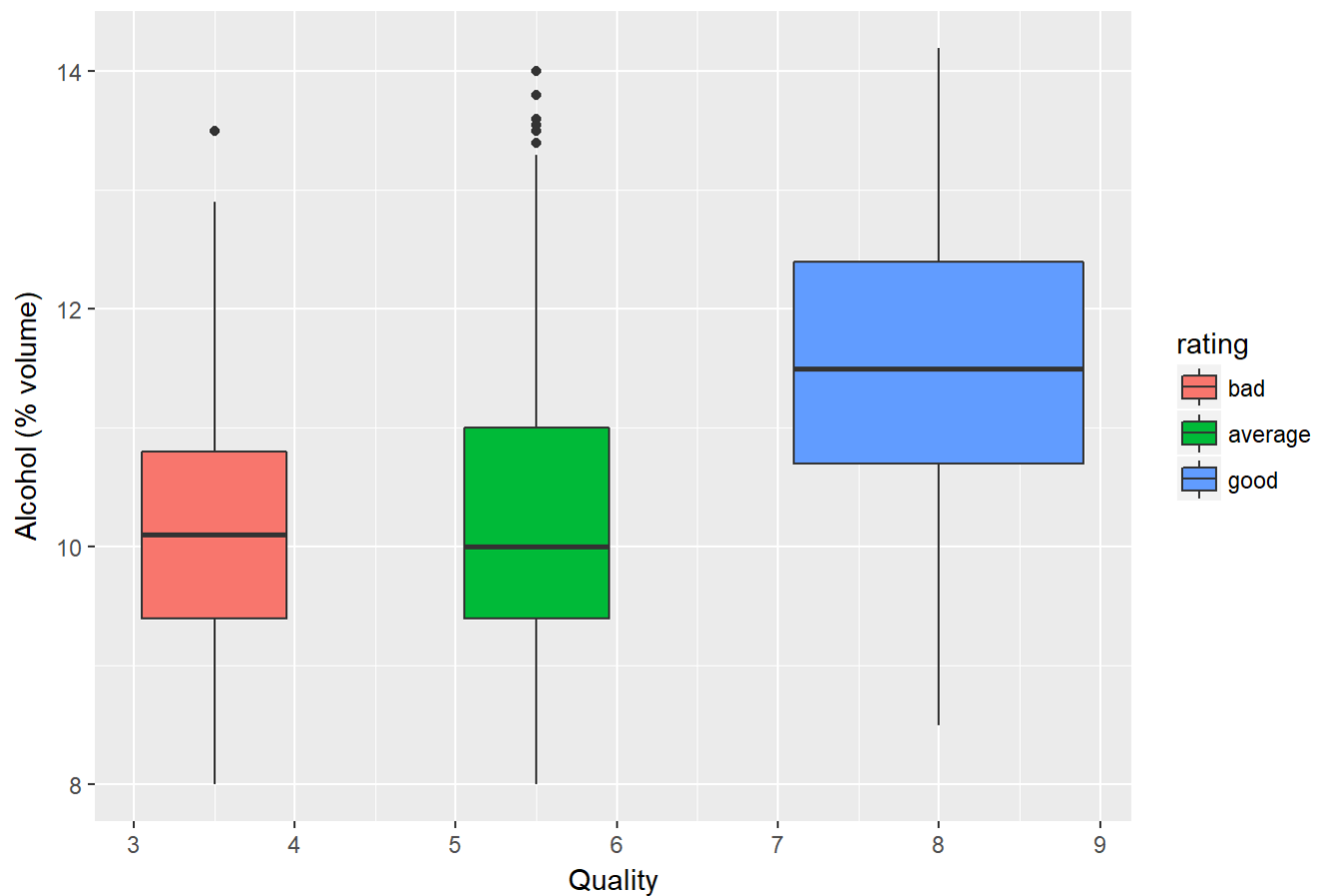
```
## Warning: Continuous x aesthetic -- did you forget aes(group=...)?
```



Plot Two These plots show the effect of alcohol on wine quality. On average, we observe that the wine quality increases with alcohol, except for wine quality 5. To further understand the difference between quality 4 and 5, I computed the stats considering only wines with that rating.

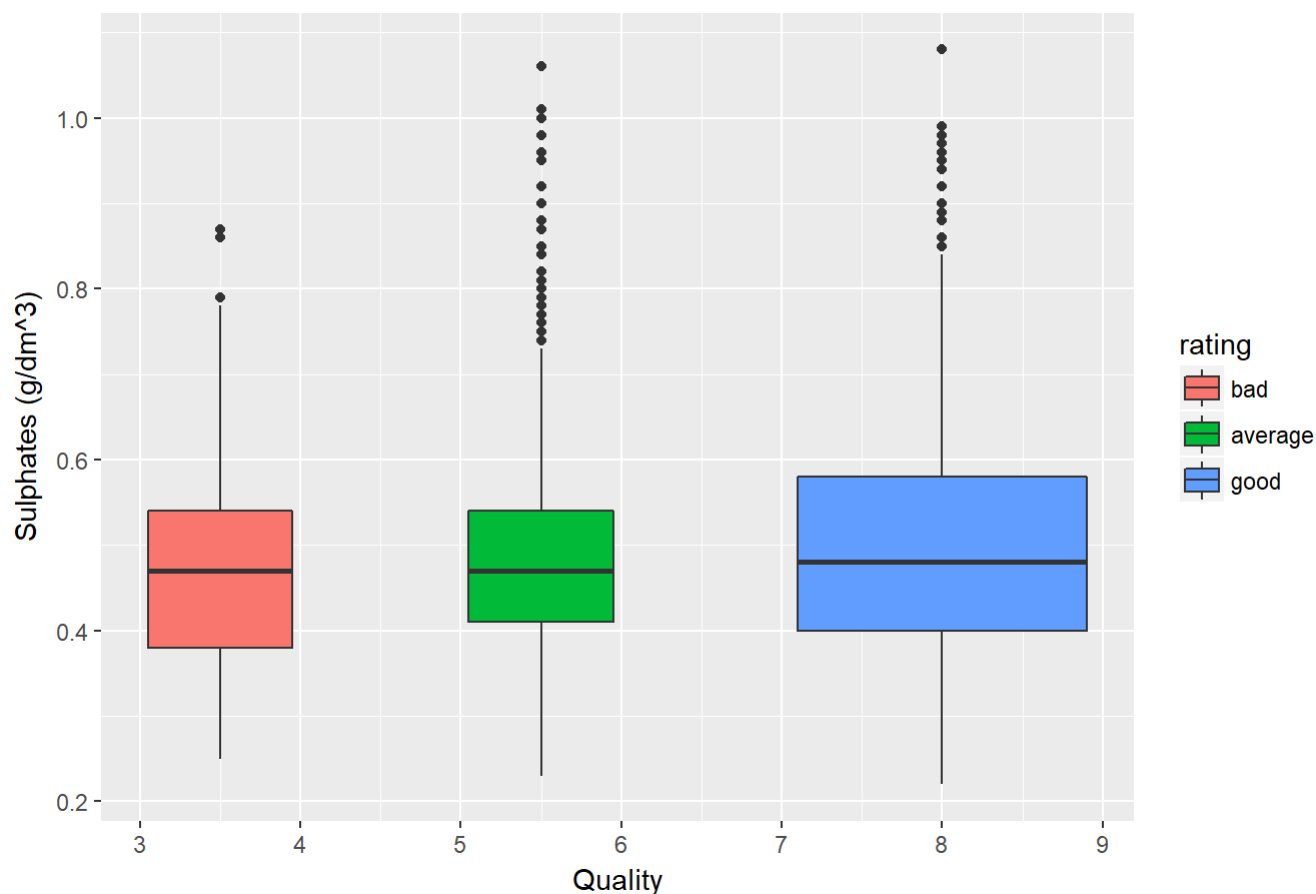
```
ggplot(data=wine, aes(x=quality, y=alcohol, fill=rating))+
  geom_boxplot()+
  ggtitle('Alcohol Levels in Different Wine Qualities')+
  xlab('Quality')+
  ylab('Alcohol (% volume)')
```

Alcohol Levels in Different Wine Qualities



```
ggplot(data=wine, aes(x=quality, y=sulphates, fill=rating))+  
  geom_boxplot()+  
  ggtitle('Sulphates Levels in Different Wine Qualities')+  
  xlab('Quality')+  
  ylab('Sulphates (g/dm^3)')
```


Sulphates Levels in Different Wine Qualities



```
print("Alcohol summaries for quality = 4")
```

```
## [1] "Alcohol summaries for quality = 4"
```

```
summary(subset(wine$alcohol, wine$quality==4))
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      8.40   9.40   10.10   10.15  10.75   13.50
```

```
print("Alcohol summaries for quality = 5")
```

```
## [1] "Alcohol summaries for quality = 5"
```

```
summary(subset(wine$alcohol, wine$quality==5))
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      8.000   9.200   9.500   9.809  10.300  13.600
```

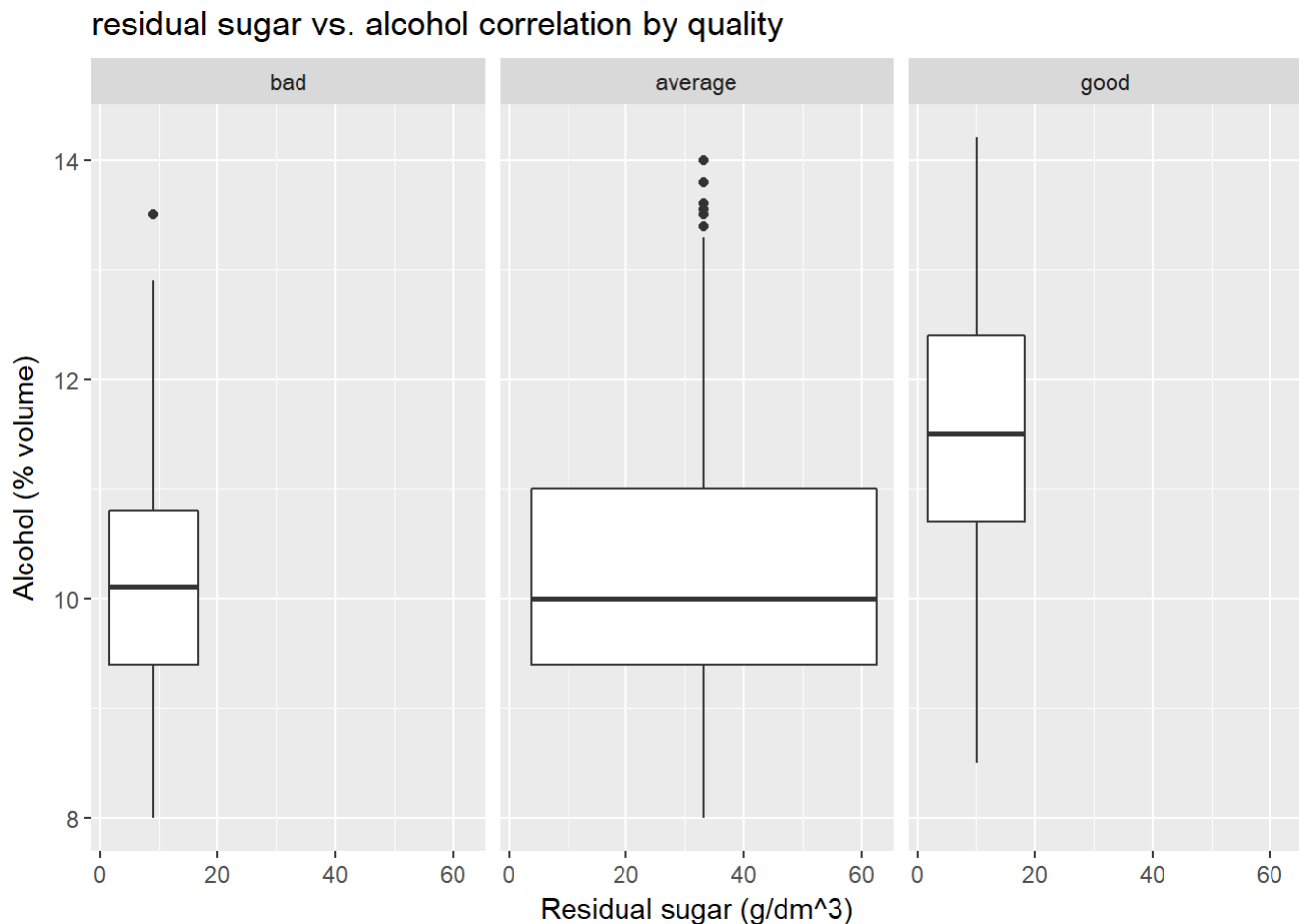
I noticed tha the wine quality increases with sulphates. However, having more of it will not help much sulphates are only good for quality upto a certain level.

Plot Three

This boxplot shows the correlation between quality and residual sugar vs. alcohol. I observe that lower residual sugars and higher alcohol leads to better wine quality.

```
ggplot(data=wine, aes(x=residual.sugar, y=alcohol))+
  facet_wrap(~rating) +
  geom_boxplot() +
  xlab('Residual sugar (g/dm^3)')+
  ylab('Alcohol (% volume)') +
  ggtitle('residual sugar vs. alcohol correlation by quality')
```

```
## Warning: Continuous x aesthetic -- did you forget aes(group=...)?
```



Reflections

This study about white wine reveals following:

- Most wines in the dataset are rated as 'average' just small number is rated as bad.
- The data about citric acid is not accurate because many entries are set to 0.
- Higher alcohol level impacts the wine rating. That is, higher alcohol means better wine.
- More alcohol and low residual sugars yield to good red wines.
- Residual sugar and chlorides do not seem to have much impact on wine quality.
- The lower volatile acidity the better wine (as there is negative correlation).

Given that only alcohol correlated with quality of wine, it suggests that alcohol has impact on wine quality. However there are other factors like harvest year, vineyard location, temperature of harvest missing from the dataset. It will be interesting to do further exploration on these factors and see how it affects the quality of wine.

My major struggle in this project was to find more information about the wine that would have been important in this study. I spent sometime online researching about information on wines and this project increased my knowledge about wines. I think next time while drinking wine I will recall this project and check if my rating matches with what I have discovered.

Thank you