

Red Wine Quality

Made by majafoi

UPDATED: 1/8/2021

MADE IN R

License – Include the citation of authors of the dataset, but mine also. And before you use this code, please contact me.

```
library(readr)
RWQ<-read.csv("winequality-red.csv")

## Parsed with column specification:
## cols(
##   `fixed acidity` = col_double(),
##   `volatile acidity` = col_double(),
##   `citric acid` = col_double(),
##   `residual sugar` = col_double(),
##   chlorides = col_double(),
##   `free sulfur dioxide` = col_double(),
##   `total sulfur dioxide` = col_double(),
##   density = col_double(),
##   pH = col_double(),
##   sulphates = col_double(),
##   alcohol = col_double(),
##   quality = col_double()
## )

cs = complete.cases(RWQ) # since there is a possibility that this original dataset has some missing data or NAs (not available data), I will use the function complete.cases in order to remove that part of the dataset.
RWQ= RWQ [cs,]
attach(RWQ) # I will only work with a complete dataset, or dataset without missing data and NAs.

Dfrm <- as.data.frame(RWQ)

colnames(dfrm) <- c("fixed.acidity","volatile.acidity","citric.acid","residual.sugar","chlorides","free.sulfur.dioxide","total.sulfur.dioxide","density","pH","sulphates","alcohol","quality")
# it is easier to work with variables without separation between the names

attach(dfrm) # this function attaches the new data table with all changes made above.
```

Hello, all!

In this R notebook, I will show you all the details about statistical analysis I performed based on a dataset called Red Wine Quality, and the dataset is found on the following link:

<https://archive.ics.uci.edu/ml/datasets/wine+quality> (it is also on Kaggle website).

Citation info: P. Cortez, A. Cerdeira, F. Almeida, T. Matos and J. Reis. Modeling wine preferences by data mining from physicochemical properties. In Decision Support Systems, Elsevier, 47(4):547-553, 2009.

As you can see, I do a lot of explaining, as I can't know who will be reading this analysis, and I want to make the results as clear as possible to the non-data scientist people.

The variables are as it follows:

1. **Fixed acidity** – most acids involved with wine are fixed or nonvolatile (do not evaporate readily),
2. **Volatile acidity** – the amount of acetic acid in wine, which at too high levels can lead to an unpleasant, vinegar taste of the wine,
3. **Citric acid** – found in small quantities, citric acid can add „freshness“ and flavor to wine,
4. **Residual sugar** – the amount of sugar remaining after fermentation stops, it is rare to find wines with less than 1 g/L and wines with greater than 45g/L are considered sweet,
5. **Chlorides** – the amount of salt in the wine,
6. **Free sulfur dioxide** – the free form of SO₂ exists in equilibrium between molecular SO₂ (as a dissolved gas) and bisulfite ion; it prevents microbial growth and the oxidation of wine,
7. **Total sulfur dioxide** - amount of free and bound forms of SO₂; in low concentrations, SO₂ is mostly undetectable in wine, but at free SO₂ concentrations over 50 ppm, SO₂ becomes evident in the nose and taste of wine,
8. **Density** - the density of water is close to that of water depending on the percent alcohol and sugar content,
9. **ph** - describes how acidic or basic the wine is, on a scale from 0 (very acidic) to 14 (very basic); most wines are between 3-4 on the pH scale,
10. **Sulphates** - a wine additive which can contribute to sulfur dioxide gas (SO₂) levels, which acts as an antimicrobial and antioxidant,
11. **Alcohol** - the percent alcohol content of the wine,
12. **Quality** - output variable (based on sensory data, score between 0 and 10).

- The main goal of this analysis is to perform an EDA (Exploratory Data Analysis), and to try to create a model in order to predict the Quality variable of the red wines in the sample.

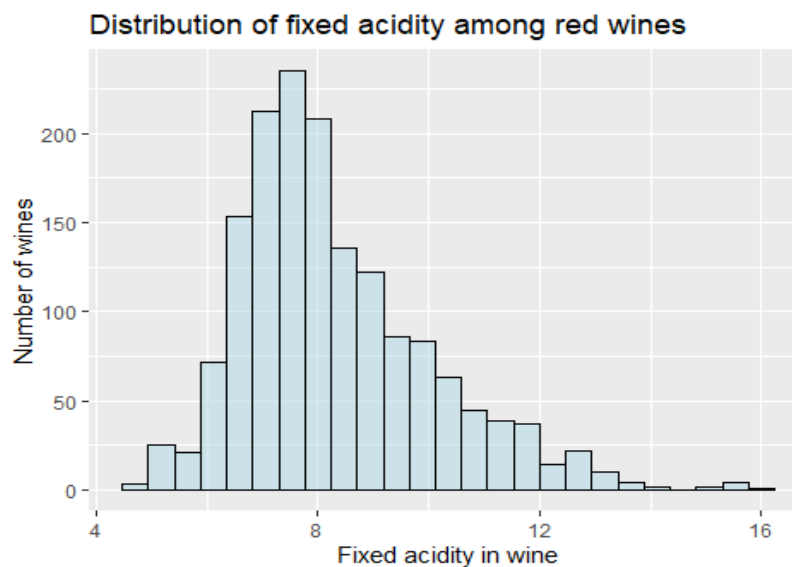
Summaries of the variables

In this chapter, I am going to create summaries of the variables in the sample. The sample itself has 1599 observations or wines, and 12 variables which have been presented and explained above.

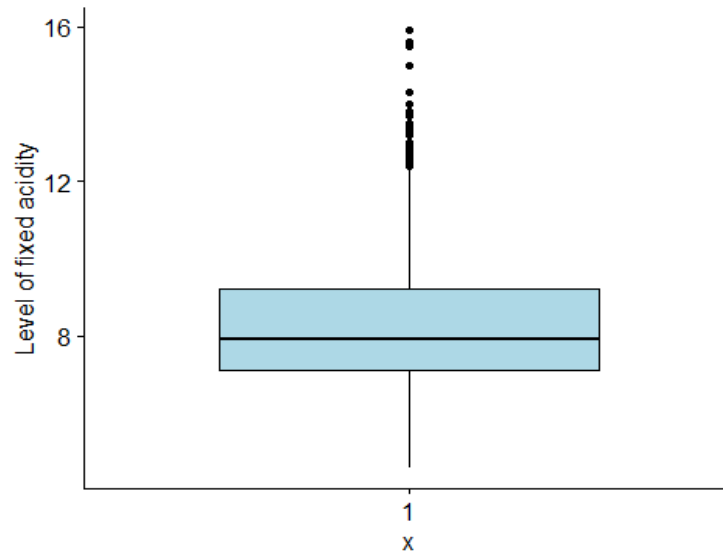
`summary` (fixed.acidity)

```
##  Min. 1st Qu.  Median    Mean 3rd Qu.   Max.
##  4.60  7.10   7.90   8.32   9.20  15.90
```

`gf_histogram` (~fixed.acidity, fill="lightblue", color="black", xlab="Fixed acidity in wine", ylab="Number of wines", title="Distribution of fixed acidity among red wines")



`ggboxplot` (fixed.acidity, fill="lightblue", color="black", ylab="Level of fixed acidity")



shapiro.test (fixed.acidity)

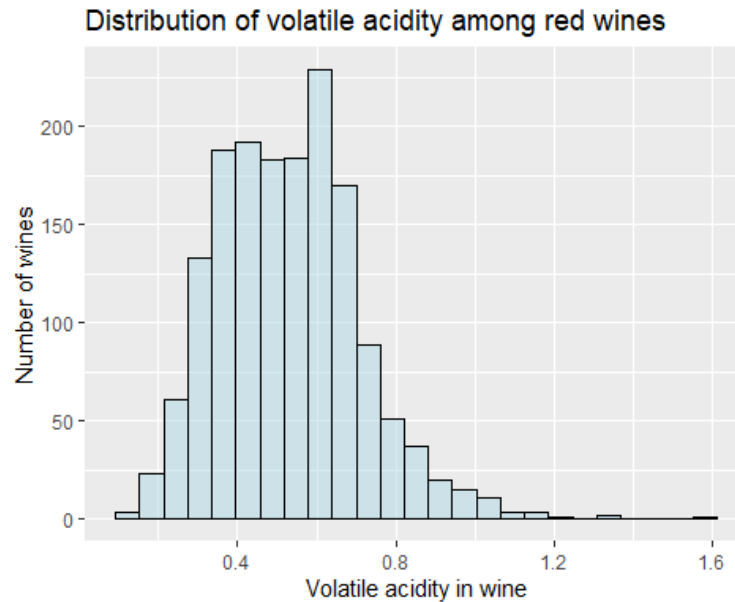
```
##
## Shapiro-Wilk normality test
##
## data: fixed.acidity
## W = 0.94203, p-value < 2.2e-16
```

The average fixed acidity level in the sample's red wines is 8.32, but the distribution has positive skewness, as it is last a longer tail on its right side (which means that there are wines in the sample with fixed acidity level greater than 12). Those are being considered as outliers (different than others), which are also shown as black dots on boxplot. Right away, I performed Shapiro test which shows that this variable distribution hasn't come from a normal population distribution, or that the population itself has outliers too. Shapiro test is a statistical test which is used to determine whether a data sample is from a normally distributed population, as that is a prerequisite for many statistical tests. If a p value is lower than 0.05, that indicates that the population is not normally distributed, whereas a p value greater than 0.05 provides no such evidence.

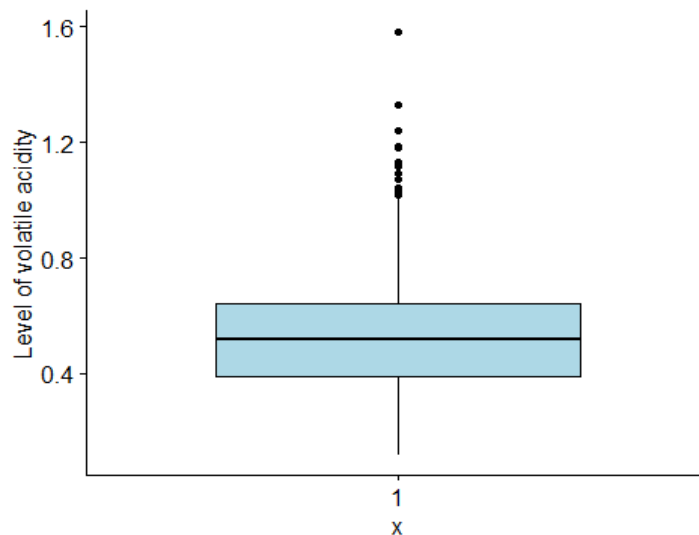
summary (volatile.acidity)

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
## 0.1200 0.3900 0.5200 0.5278 0.6400 1.5800
```

gf_histogram (~volatile.acidity, fill="lightblue", color="black", xlab="Volatile acidity in wine", ylab="Number of wines", title="Distribution of volatile acidity among red wines")



```
ggboxplot (volatile.acidity, fill="lightblue", color="black", ylab="Level of volatile acidity")
```



```
shapiro.test (volatile.acidity)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: volatile.acidity  
## W = 0.97434, p-value = 2.693e-16
```

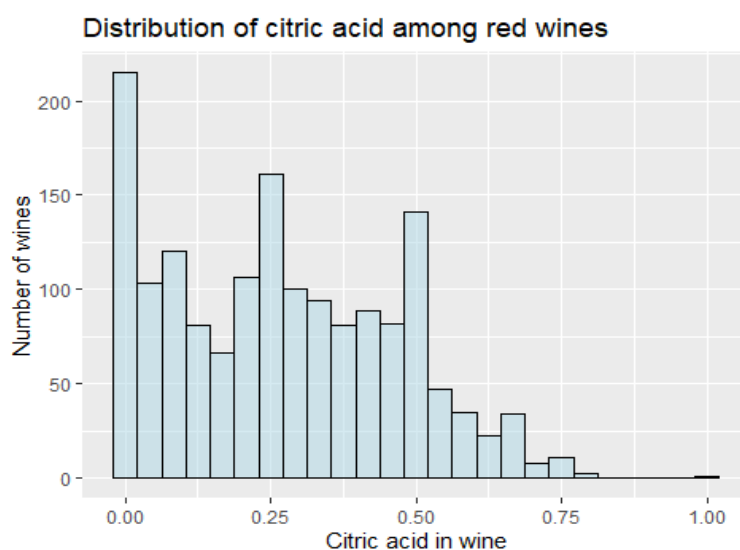
Volatile acidity is the amount of acetic acid in wine which, at too high of levels, can lead to an unpleasant and vinegar taste. The average volatile acidity of red wines in the sample is 0.52, but there is a big difference

between minimum and maximum value (points that some wines have too high level of acetic acid in it). At histogram, we can see that the distribution has positive skewness and outliers on boxplot (points on the same things and conclusion as for the first variable). Also, one outlier on the boxplot is away from other outliers (point at 1.6). Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

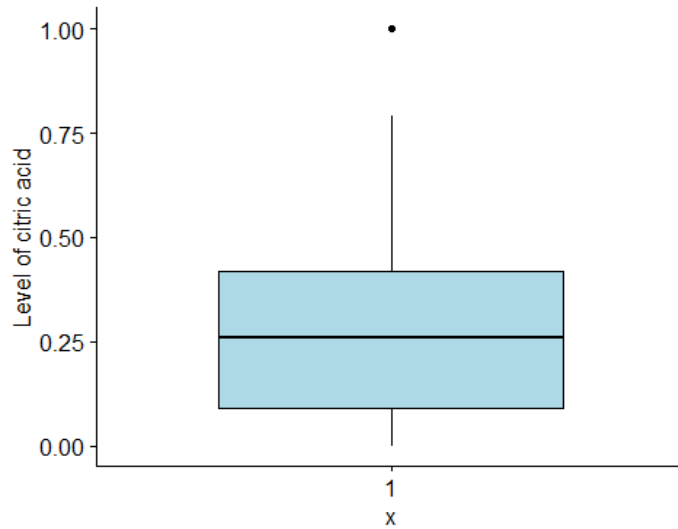
```
summary (citric.acid)
```

```
##  Min. 1st Qu. Median  Mean 3rd Qu.  Max.
##  0.000  0.090  0.260  0.271  0.420  1.000
```

```
gf_histogram (~citric.acid, fill="lightblue", color="black", xlab="Citric acid in wine", ylab="Number of wines", title="Distribution of citric acid among red wines")
```



```
ggboxplot (citric.acid, fill="lightblue", color="black", ylab="Level of citric acid")
```



Citric acid is found in small quantities in wines, adding “freshness” and flavor to wines. So, we can obviously see that the levels of this acid are much smaller than the amounts of acetic acid. Most red wines in the sample have a level of citric acid between 0 and 0.5, with the average value of 0.27. On the boxplot it is visible that only one red wine is considered as outlier, with the level of citric acid around 1. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn’t have a normal distribution.

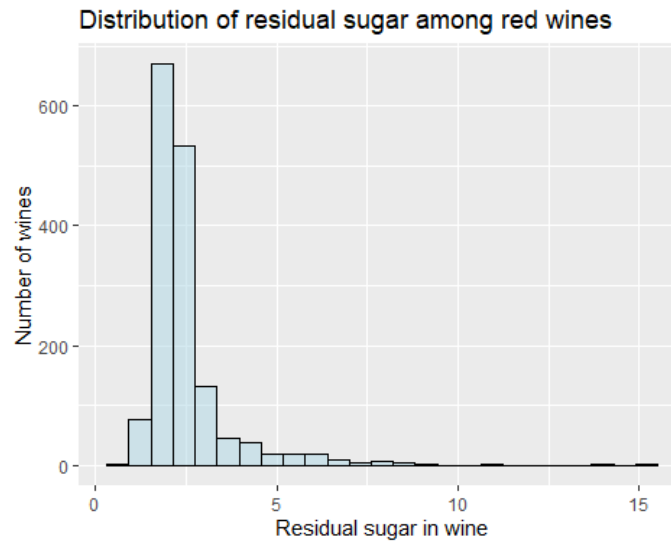
```
summary(residual.sugar)
```

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.   Max.
##  0.900  1.900  2.200  2.539  2.600 15.500
```

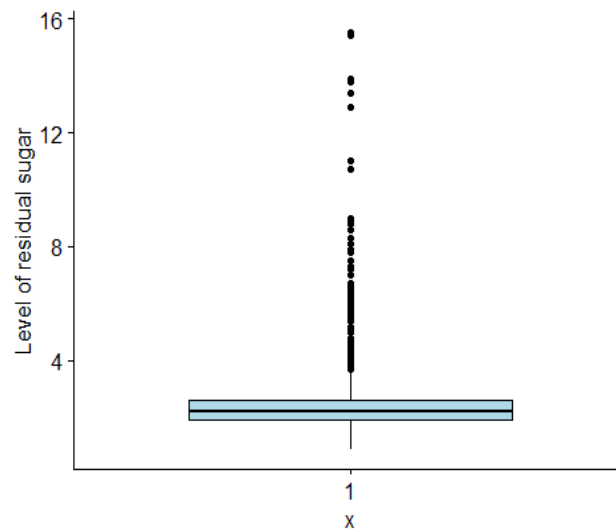
```
sd(residual.sugar)
```

```
## [1] 1.409928
```

```
gf_histogram(~residual.sugar, fill="lightblue", color="black", xlab="Residual sugar in wine", ylab="Number of wines", title="Distribution of residual sugar among red wines")
```



```
ggboxplot (residual.sugar, fill="lightblue", color="black", ylab="Level of residual sugar")
```



```
shapiro.test (residual.sugar)
```

```
##
## Shapiro-Wilk normality test
##
## data: residual.sugar
## W = 0.56608, p-value < 2.2e-16
```

Residual sugar is the amount of sugar remaining after fermentation stops, it is very rare to find wines with less than 1g/L and wines with greater than 45g/L (those are considered sweet). In this sample's distribution, we can see that the minimum value is 0.9, which makes that wine – rare, when speaking of residual sugar, but the maximum value is below 45g/L - only 15.5. The average amount of residual sugar in these red wines is 2.53,

with standard deviation of 1.4, which means that outliers haven't had much influence on the distribution. Looking at the histogram distribution, that is not expected because we can see a huge positive skewness. But, as we can see on boxplot, there are a lot of red wine's residual sugar which are outliers. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

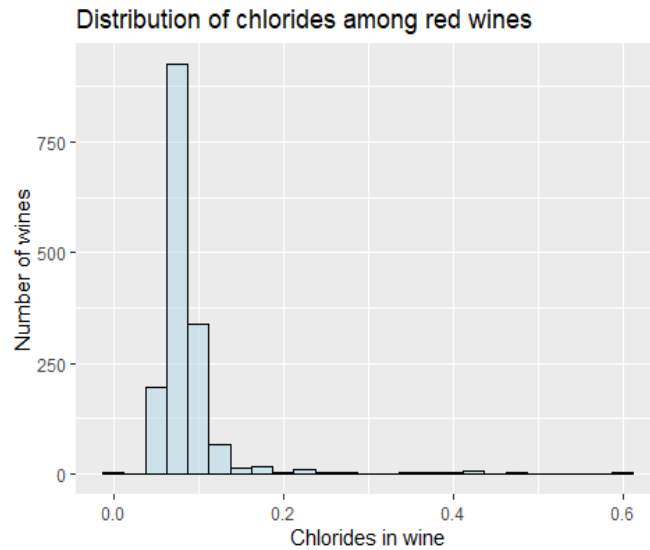
```
summary(chlorides)
```

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.   Max.
## 0.01200 0.07000 0.07900 0.08747 0.09000 0.61100
```

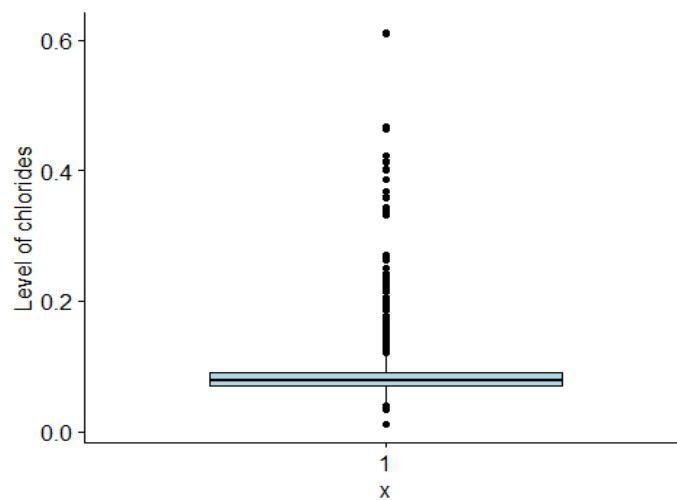
```
sd(chlorides)
```

```
## [1] 0.0470653
```

```
gf_histogram(~chlorides, fill="lightblue", color="black", xlab="Chlorides in wine", ylab="Number of wines",
title="Distribution of chlorides among red wines")
```



```
ggboxplot (chlorides, fill="lightblue", color="black", ylab="Level of chlorides")
```



```
shapiro.test(chlorides)
```

```
##  
## Shapiro-Wilk normality test  
##  
## data: chlorides  
## W = 0.48425, p-value < 2.2e-16
```

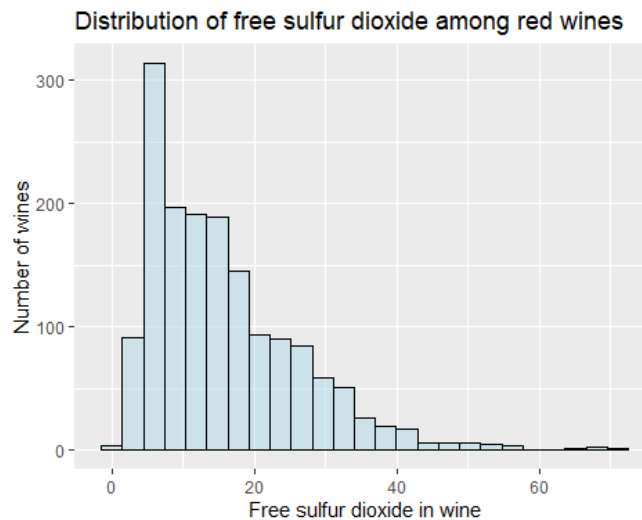
Chlorides represent the amount of salt in the wine. The average amount of chlorides in these red wines is 0.08, but again - the minimum and maximum value differentiate largely. This time, we also have a positive skewness

of the distribution, but also outliers on both sides of the medium and average value. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

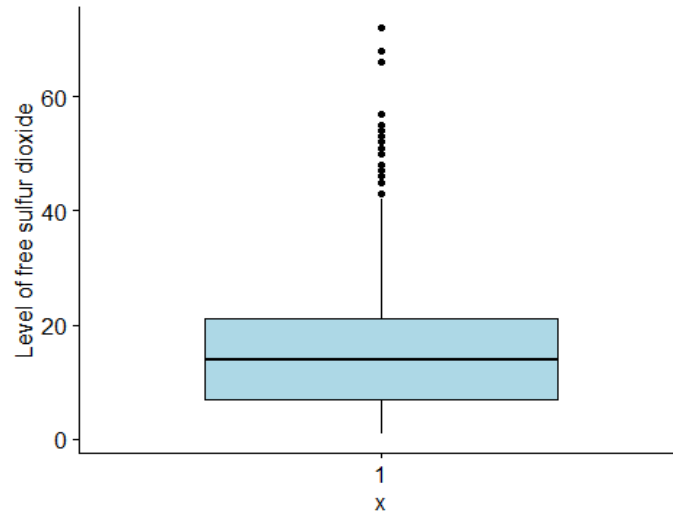
```
summary (free.sulfur.dioxide)
```

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.   Max.
##   1.00   7.00  14.00  15.87  21.00  72.00
```

```
gf_histogram (~free.sulfur.dioxide, fill="lightblue", color="black", xlab="Free sulfur dioxide in wine", ylab="Number of wines", title="Distribution of free sulfur dioxide among red wines")
```



```
ggboxplot (free.sulfur.dioxide, fill="lightblue", color="black", ylab="Level of free sulfur dioxide")
```



shapiro.test (free.sulfur.dioxide)

```
##  
## Shapiro-Wilk normality test  
##  
## data: free.sulfur.dioxide  
## W = 0.90184, p-value < 2.2e-16
```

Free sulfur dioxide is the free form of SO₂ which exists in the equilibrium between molecular SO₂ (as a dissolved gas) and bisulfite ion, it also prevents microbial growth and the oxidation of the wine. The minimum and maximum value differentiate highly here too, whereas the average value of free sulfur dioxide is 15.87. This distribution has a hard positive skewness, with outliers of some wines above 40 amount of free sulfur dioxide. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

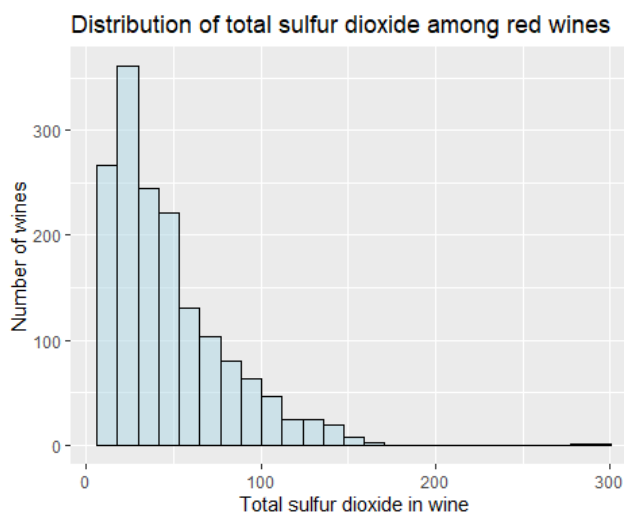
summary (total.sulfur.dioxide)

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   6.00  22.00  38.00  46.47  62.00 289.00
```

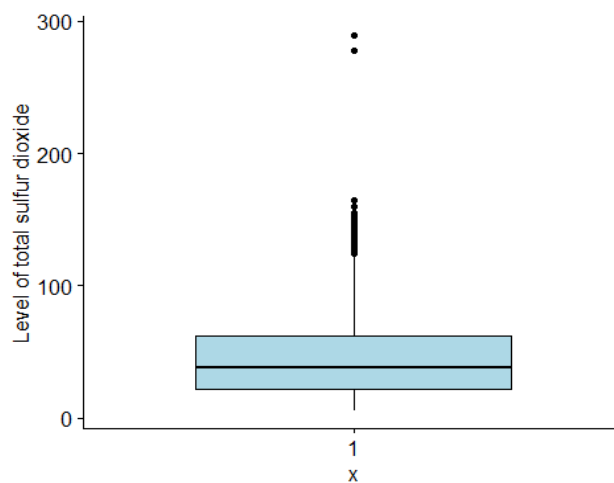
```
sd (total.sulfur.dioxide)
```

```
## [1] 32.89532
```

```
gf_histogram (~total.sulfur.dioxide, fill="lightblue", color="black", xlab="Total sulfur dioxide in wine", ylab="Number of wines", title="Distribution of total sulfur dioxide among red wines")
```



```
ggboxplot (total.sulfur.dioxide, fill="lightblue", color="black", ylab="Level of total sulfur dioxide")
```



```
shapiro.test (total.sulfur.dioxide)
```

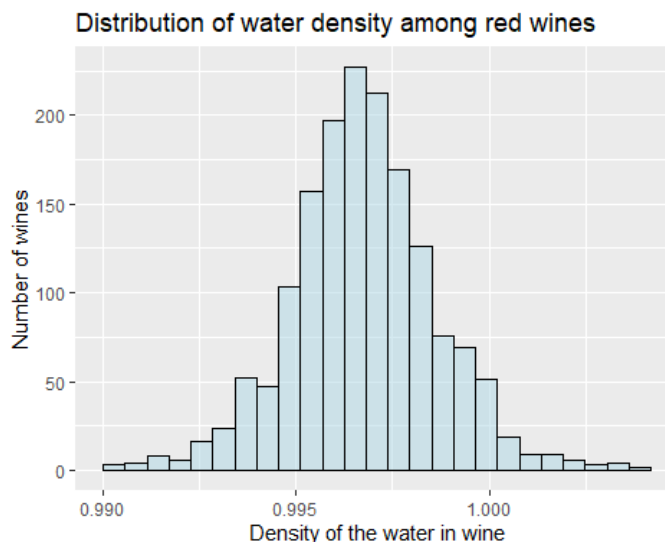
```
##
## Shapiro-Wilk normality test
##
## data: total.sulfur.dioxide
## W = 0.87322, p-value < 2.2e-16
```

Total sulfur dioxide is the amount of free and bound forms of SO₂. In low concentrations, SO₂ is mostly undetectable in wine, but at free SO₂ concentrations over 50 ppm, SO₂ becomes evident in the nose of people and taste of wine. Here, some wines have that concentration above 50 ppm, even the average value is pretty close. Maximum value is astonishing 289 ppm! This distribution has a big positive skewness, with outliers with values around 120 ppm, but there are some extreme outliers around 300 ppm. We should take care of this variable's distribution, as it has the biggest outliers so far, and that can distort future statistical tests or make those tests insignificant. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

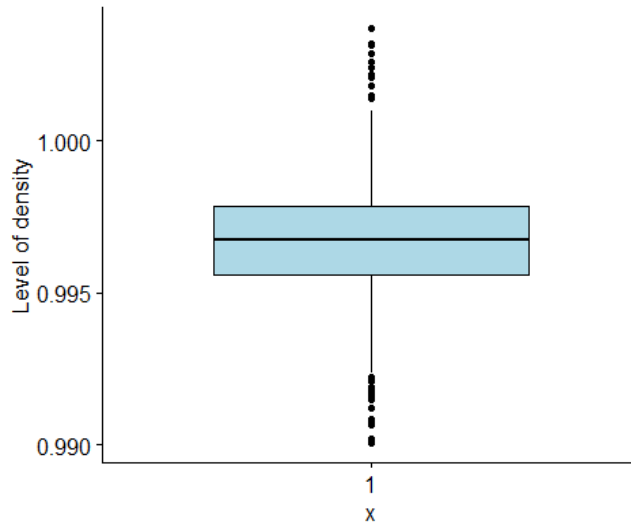
summary (density)

```
## Min. 1st Qu. Median Mean 3rd Qu. Max.
## 0.9901 0.9956 0.9968 0.9967 0.9978 1.0037
```

gf_histogram (~density, fill="lightblue", color="black", xlab="Density of the water in wine", ylab="Number of wines", title="Distribution of water density among red wines")



ggboxplot (density, fill="lightblue", color="black", ylab="Level of density")



shapiro.test (density)

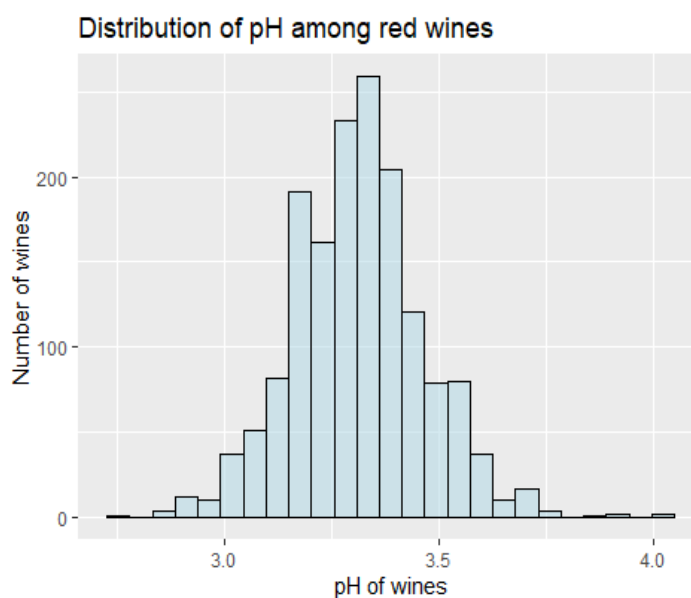
```
##  
## Shapiro-Wilk normality test  
##  
## data: density  
## W = 0.99087, p-value = 1.936e-08
```

Density represents density of water, which is close to that of the regular water itself, depending on the percent alcohol and sugar content. The mean and the median almost are the same, and the histogram is showing normal distribution. The boxplot shows that there are some outliers on both sides, but those aren't outliers visible on the histogram, and the summary proves it. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution, whereas we can say that this variable's distribution in the sample is normal. But this difference between distributions in the sample and population can happen when you are extracting numerous samples from one big population. Some samples tend to have many outliers, whereas some samples are normally distributed.

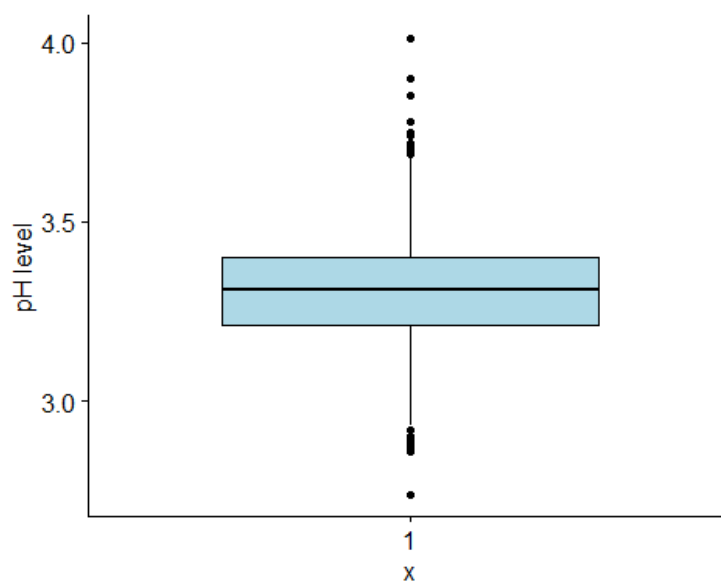
summary (pH)

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##  2.740  3.210   3.310   3.311  3.400   4.010
```

```
gf_histogram (~pH, fill="lightblue", color="black", xlab="pH of wines", ylab="Number of wines", title="Dist  
ribution of pH among red wines", xlim=c(40,95), ylim=c(0,50))
```



```
ggboxplot (pH, fill="lightblue", color="black", ylab="pH level")
```



```
shapiro.test (pH)
```



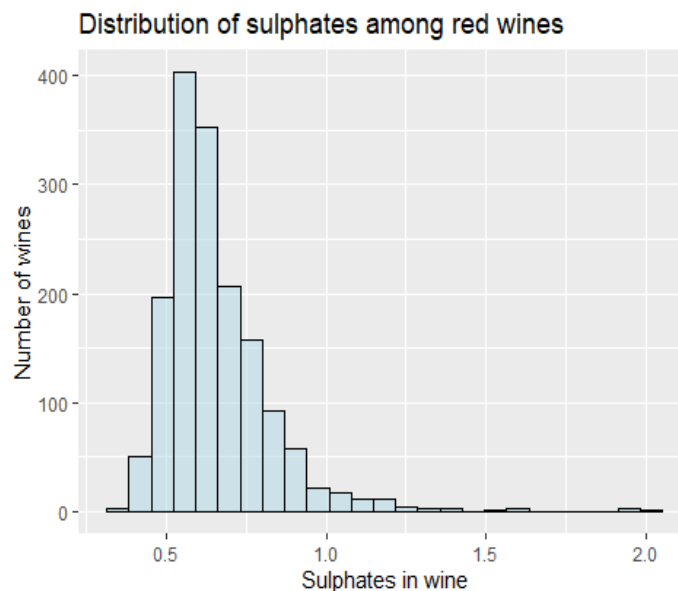
```
##  
## Shapiro-Wilk normality test  
##  
## data: pH  
## W = 0.99349, p-value = 1.712e-06
```

pH describes how acidic or basic a wine is on a scale from 0 (very acidic) to 14 (very basic), and most wines are between 3-4 on the pH scale. That is also visible here on the histogram, which almost has a normal looking distribution too, as a variable density before. The mean and median do differentiate, but just slightly. Boxplot does show that there are outliers on both sides of the distribution, but the histogram and the summary shows us that those aren't real outliers. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

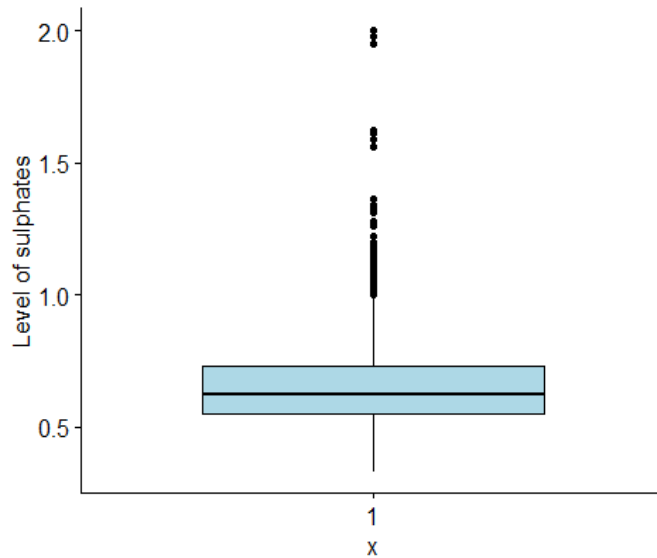
summary (sulphates)

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
## 0.3300 0.5500 0.6200 0.6581 0.7300 2.0000
```

```
gf_histogram (~sulphates, fill="lightblue", color="black", xlab="Sulphates in wine", ylab="Number of wines", title="Distribution of sulphates among red wines", xlim=c(40,95), ylim=c(0,50))
```



```
ggboxplot (sulphates, fill="lightblue", color="black", ylab="Level of sulphates")
```



`shapiro.test` (sulphates)

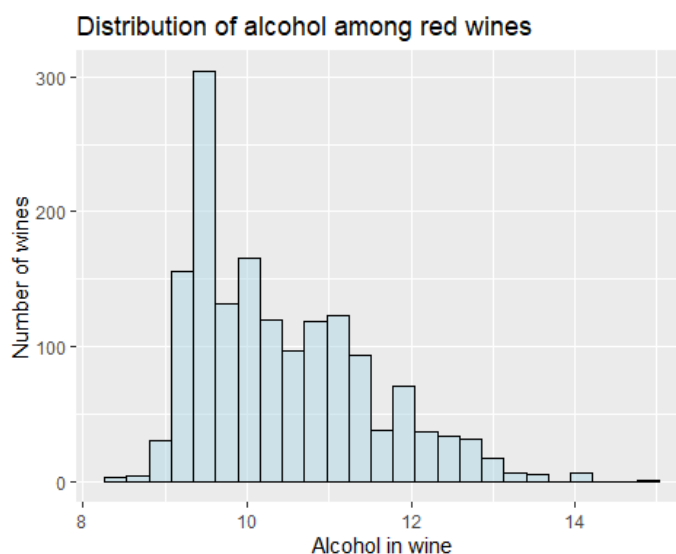
```
##  
## Shapiro-Wilk normality test  
##  
## data: sulphates  
## W = 0.83304, p-value < 2.2e-16
```

Sulphates are a wine additive which can contribute to sulfur dioxide gas (SO₂) levels, which acts as an antimicrobial and antioxidant. Mean and the median are close together, but not enough to say it is a normal distribution of the variable like for the last two variables. Also, there is a big difference between minimum and maximum value and if we look at the histogram, we can see that the distribution has a negative skewness. On boxplot, we can also see that there are outliers from level 1 and beyond. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

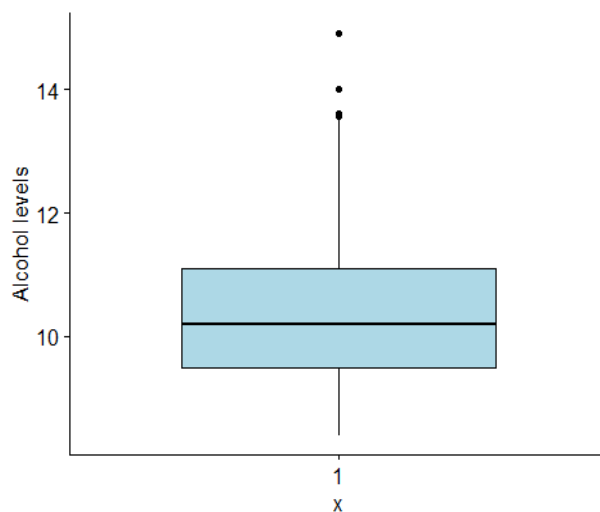
`summary` (alcohol)

```
##   Min. 1st Qu.  Median    Mean 3rd Qu.   Max.
##   8.40   9.50   10.20   10.42   11.10   14.90
```

```
gf_histogram(~alcohol, fill="lightblue", color="black", xlab="Alcohol in wine", ylab="Number of wines", title="Distribution of alcohol among red wines", xlim=c(40,95), ylim=c(0,50))
```



```
ggboxplot(alcohol, fill="lightblue", color="black", ylab="Alcohol levels")
```



```
shapiro.test(alcohol)
```

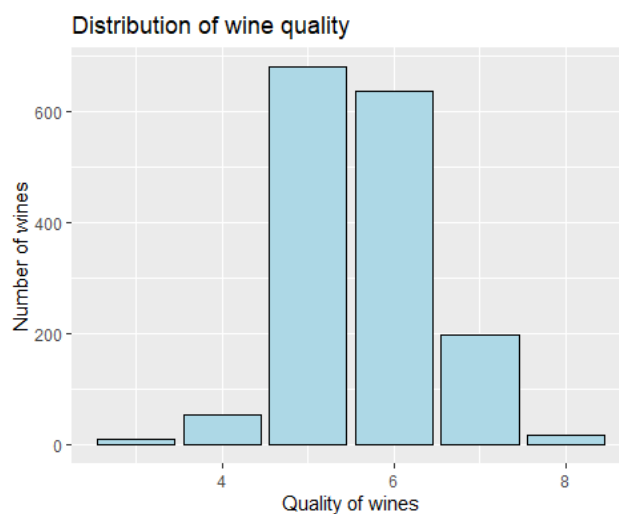
```
##
## Shapiro-Wilk normality test
##
## data: alcohol
## W = 0.92884, p-value < 2.2e-16
```

Alcohol variable represents the percent alcohol of the wine. The average percent of alcohol in the wines in the sample is 10.42, but the median isn't far from the average value/mean. But, the difference between minimum and maximum value is high, which means that there is some (positive) skewness visible on histogram. The boxplot shows that there are some outliers after a value of approximate 13.5%. Again, Shapiro test shows that the distribution of this variable has come from a population which also doesn't have a normal distribution.

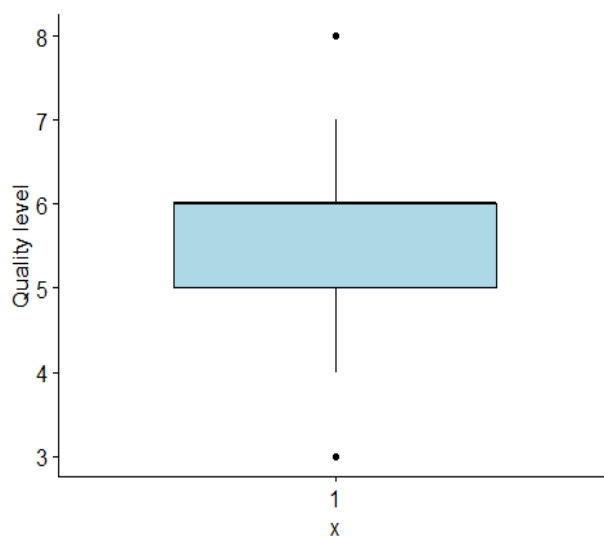
`summary` (quality)

```
##  Min. 1st Qu.  Median    Mean 3rd Qu.   Max.
## 3.000  5.000  6.000  5.636  6.000  8.000
```

`gf_bar` (~quality, fill="lightblue", color="black", xlab="Quality of wines", ylab="Number of wines", title="Distribution of wine quality", xlim=c(40,95), ylim=c(0,50))



`ggboxplot` (quality, fill="lightblue", color="black", ylab="Quality level")



The quality variable is an output variable, based on sensory data, and it scores between 0 and 10. Unfortunately, the average quality grade is 5.63. The difference is existent between minimum and maximum value, but the histogram seems to show pretty normal distribution. Boxplot shows outliers on both sides, which, again, shouldn't be real outliers.

With this, I have finished summaries of all variables in the sample. Since we want to predict the quality variable, this is going to be our y-variable or dependent variable, whereas all other variables will be independent variables (x-variables). More info on that in the next chapter!

Regression analysis and correlation

In the prior chapter, I have shown summaries and basic, simple graphics of each variable. The main goal of this dataset analysis is that we determine which variables make a wine good, or which variables give higher quality of red wines in the sample, but also in the population as well. For that, I am going to calculate correlations, to see what the relationship between variables is. After that, I am going to calculate linear regression models where the y-variable (dependent variable) is quality variable, and all other variables will be independent variables or x-variables. So, let us start with correlation first. Since there are 66 combination of correlation, I am making a correlation matrix, and then proceed from there to extract only significant correlations (above 0.5) and make a regression model. Since the output correlation matrix is very strange-looking in Word output, I will skip it here.

Correlation is a statistical measure of a relationship between two variables. It can go between -1 and +1. If it is positive, that means that if a value of one variable goes up, the value of the other variable goes up too. The same logic goes for negative values (it goes down). The best values for correlation is above 0.4, as all below that aren't that significant.

From the correlation matrix, I see 6 correlation coefficients higher than 0.50. Those are:

- 1) Citric acid and fixed acidity = 0.672
- 2) Density and fixed acidity = 0.668
- 3) pH and fixed acidity = -0.683
- 4) Citric acid and volatile acidity = -0.55
- 5) pH and citric acid = -0.54
- 6) total sulfur dioxide and free sulfur dioxide = 0.66

Earlier, during the summary part, I saw what could be a higher correlation coefficient, between density, residual sugar and alcohol. Unfortunately, coefficient between density and residual sugar variables, and coefficient between residual sugar and alcohol variables aren't high enough. Only the coefficient between Alcohol and density is -0.496, which is pretty close to our limit (0.5), so I'll take that one in consideration also. Now, at last, we have 7 significant correlations to check via regression model. Most of them aren't really significant, when looking at *adjusted R-squared which explains how much of the variability is being explained by the model. To be significant, that R squared should be around 60-70%.*

When talking about regression model, think about it as linear equations you learn in school. Mostly, they start with $y = ax + b$, just like the regression equation which is linear. But the equation can be not-linear, which doesn't have a linear look on the graph, but as a polynomial. Those are harder to explain and not used in normal examples as much as a linear equation.

First regression model to check is $\text{citric.acid} = -0.354270 + 0.075153 \cdot \text{fixed.acidity}$.

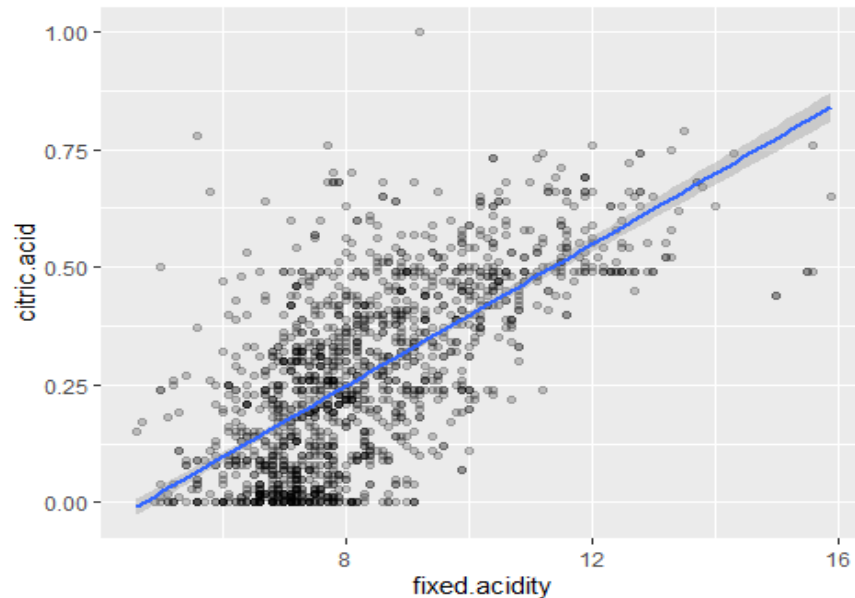
As we can see, the p-values are all less than 0,05, so the means of the model and the variables in are significant, but the adjusted R squared is less than 0.50 (it is actually 0.4508). Still, that level of R-squared is still semi-significant and useful.

```
reg1<-lm(citric.acid~fixed.acidity)
print(summary(reg1))

##
## Call:
## lm(formula = citric.acid ~ fixed.acidity)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.33302 -0.11017 -0.00938  0.09317  0.71341
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.354270   0.017629  -20.09  <2e-16 ***
## fixed.acidity  0.075153   0.002074   36.23  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1444 on 1597 degrees of freedom
## Multiple R-squared:  0.4512, Adjusted R-squared:  0.4508
## F-statistic: 1313 on 1 and 1597 DF, p-value: < 2.2e-16

ggplot (dfrm, aes(x=fixed.acidity, y=citric.acid))+ geom_point (alpha=.2) + stat_smooth (
method=lm, level=0.95)

## `geom_smooth()` using formula 'y ~ x'
```

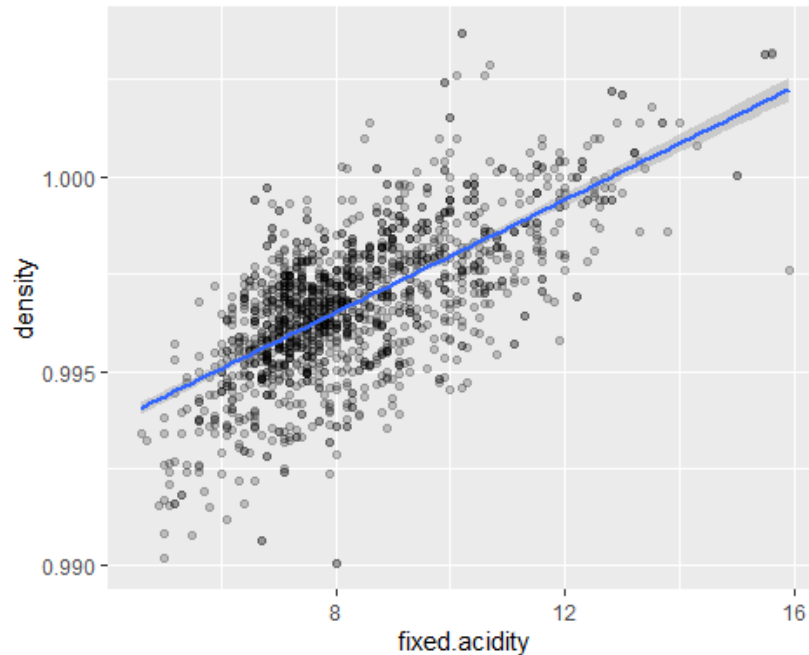
The next regression model is $\text{density} = 9.907\text{e-}01 + 7.242\text{e-}04 \times \text{fixed acidity}$. Here we are really talking about big decimal numbers after a point, and both p-values are lower than 0.05, which makes the model and the variables in it significant. The R squared value is 0.4459, which isn't particularly high.

```
reg2<-lm (density~fixed.acidity)
print (summary(reg2))

##
## Call:
## lm(formula = density ~ fixed.acidity)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0064452 -0.0007700  0.0000738  0.0009434  0.0055816
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  9.907e-01  1.716e-04  5774.70  <2e-16 ***
## fixed.acidity 7.242e-04  2.018e-05   35.88  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.001405 on 1597 degrees of freedom
## Multiple R-squared:  0.4463, Adjusted R-squared:  0.4459
## F-statistic: 1287 on 1 and 1597 DF, p-value: < 2.2e-16
```

```
ggplot (dfrm, aes(x=fixed.acidity, y=density)) + geom_point (alpha=.2) + stat_smooth(method=lm, level=0.95)
```

```
## `geom_smooth()` using formula 'y ~ x'
```



The next regression model is $\text{pH} = 3.81 - 0.06 \times \text{fixed.acidity}$. Again, p-values are lower than 0.05, which makes the model and the variables in it significant. The R squared value is 0.4661, which again isn't that particularly high to be seen significant.

```
reg3<-lm (pH~fixed.acidity)
```

```
print (summary(reg3))
```

```
##
```

```
## Call:
```

```
## lm(formula = pH ~ fixed.acidity)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
```

```
## -0.51780 -0.06547  0.00164  0.06488  0.52207
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)   3.814959   0.013776  276.93  <2e-16 ***
```

```
## fixed.acidity -0.060561   0.001621  -37.37  <2e-16 ***
```

```
## ---
```

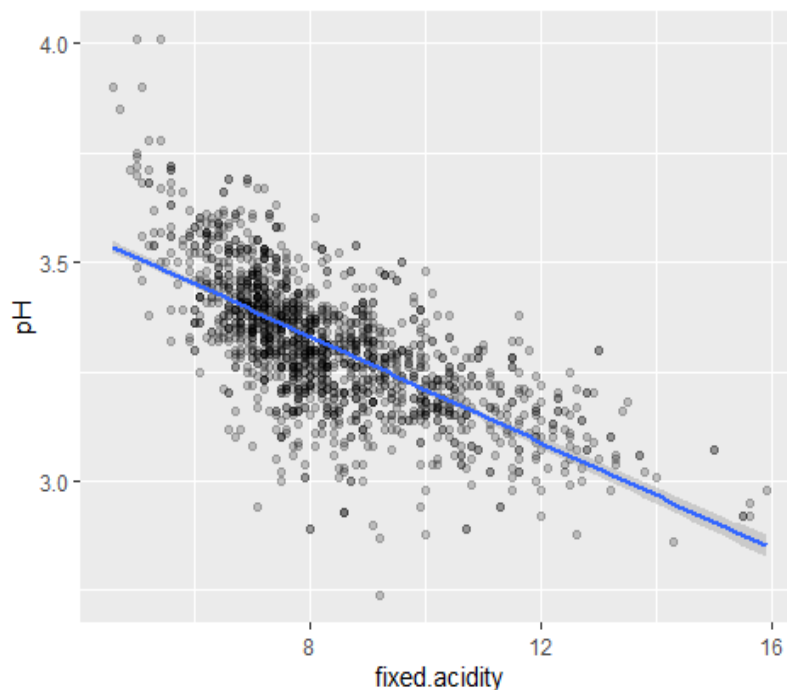
```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## Residual standard error: 0.1128 on 1597 degrees of freedom
## Multiple R-squared:  0.4665, Adjusted R-squared:  0.4661
## F-statistic: 1396 on 1 and 1597 DF, p-value: < 2.2e-16
```

```
ggplot (dfrm, aes(x=fixed.acidity, y=pH)) + geom_point (alpha=.2) + stat_smooth (method
=lm, level=0.95)
```

```
## `geom_smooth()` using formula 'y ~ x'
```



The next regression model is $\text{citric.acid} = 0.58 - 0.6 \cdot \text{volatile.acidity}$. Again, p-values are lower than 0.05, which makes the model and the variables in it significant. The R squared value is 0.3048, which again isn't that particularly high to be seen significant.

```
reg4<-lm (citric.acid~volatile.acidity)
```

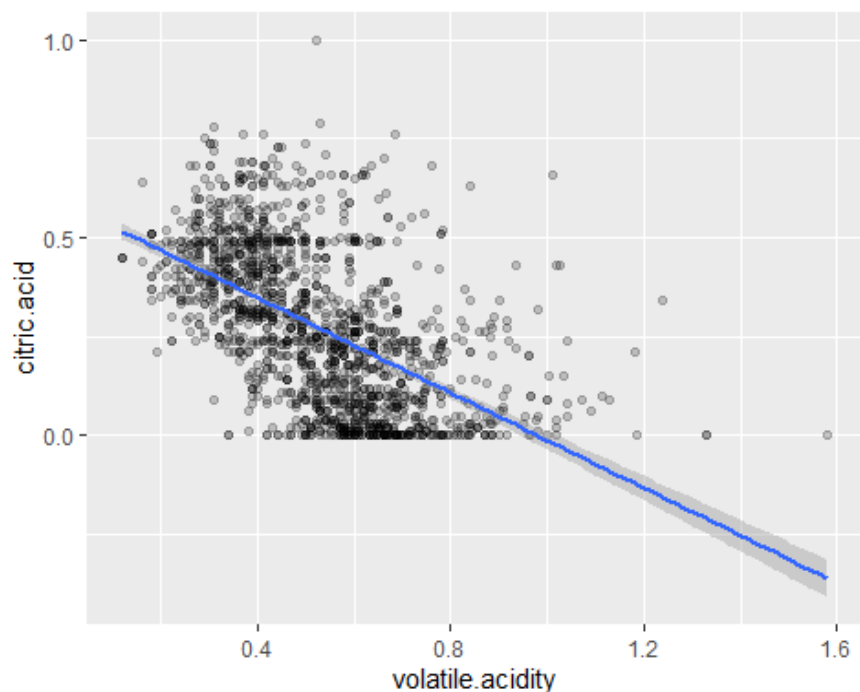
```
print (summary(reg4))
```

```
##
## Call:
## lm(formula = citric.acid ~ volatile.acidity)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.38387 -0.12073 -0.01748  0.09632  0.72432
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)    0.58823    0.01265   46.51 <2e-16 ***
## volatile.acidity -0.60107    0.02269  -26.49 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1624 on 1597 degrees of freedom
## Multiple R-squared:  0.3053, Adjusted R-squared:  0.3048
## F-statistic: 701.7 on 1 and 1597 DF, p-value: < 2.2e-16

ggplot(dfrm, aes(x=volatile.acidity, y=citric.acid)) + geom_point(alpha=.2) + stat_smooth(
  method=lm, level=0.95)

## `geom_smooth()` using formula 'y ~ x'
```



The next regression model is $\text{pH} = 3.427 - 0.429 \times \text{citric acid}$. Even though the scatter plot looks like the R squared could be higher, it is only 0.2932. The R squared is higher when the dots are close to the line.

```

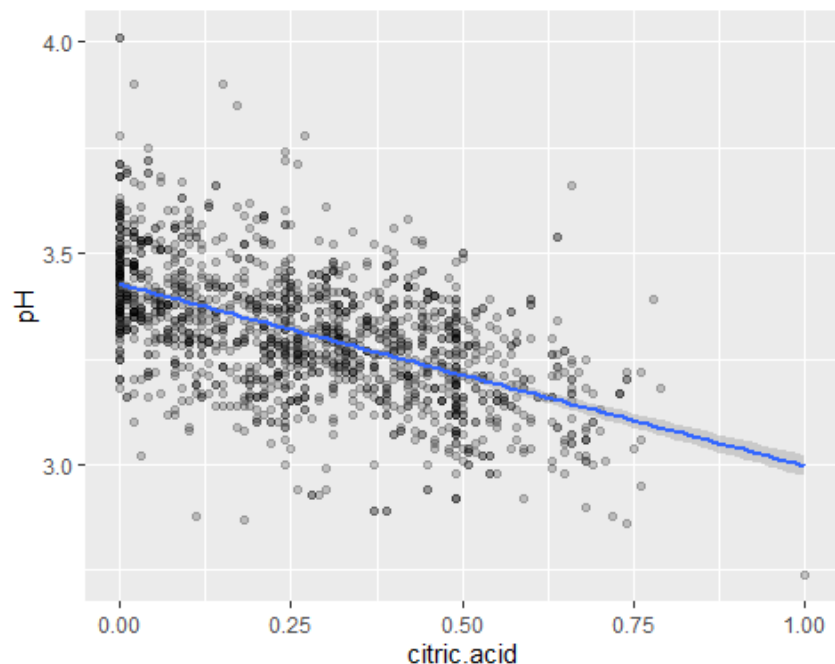
reg5<-lm (pH~citric.acid)
print (summary(reg5))

##
## Call:
## lm(formula = pH ~ citric.acid)
##
## Residuals:
##   Min     1Q   Median     3Q      Max
## -0.50025 -0.07733 -0.00570  0.08251  0.58251
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  3.427491   0.005562  616.25  <2e-16 ***
## citric.acid -0.429477   0.016668  -25.77  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1298 on 1597 degrees of freedom
## Multiple R-squared:  0.2937, Adjusted R-squared:  0.2932
## F-statistic:  664 on 1 and 1597 DF, p-value: < 2.2e-16

ggplot (dfrm, aes(x=citric.acid, y=pH)) + geom_point (alpha=.2) + stat_smooth(method=l
m, level=0.95)

## `geom_smooth()` using formula 'y ~ x'

```



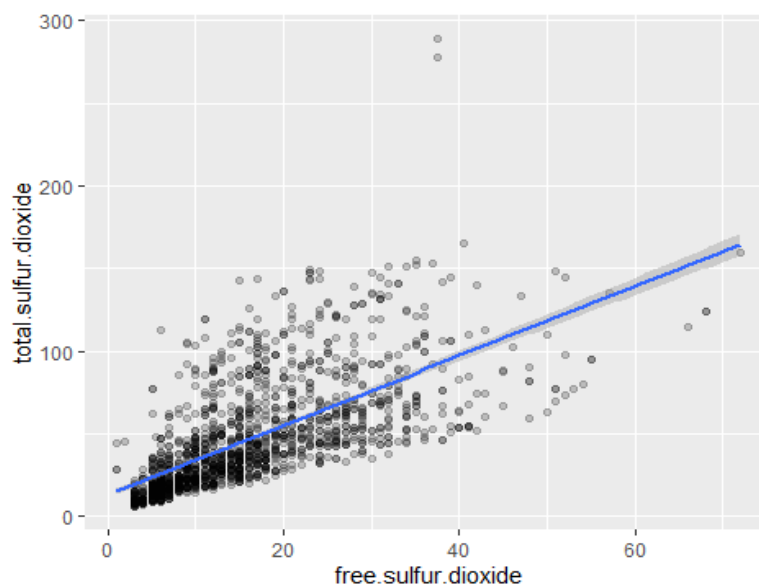
The sixth regression model is $\text{total.sulfur.dioxide} = 13.13 + 2.09 \cdot \text{free.sulfur.dioxide}$. P-values are lower than 0.05, and the R squared is 0.4454.

```
reg6<-lm (total.sulfur.dioxide~free.sulfur.dioxide)
print (summary(reg6))

##
## Call:
## lm(formula = total.sulfur.dioxide ~ free.sulfur.dioxide)
##
## Residuals:
##   Min     1Q   Median     3Q    Max
## -55.120 -13.534  -7.325   7.570 197.126
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   13.13535   1.11367   11.79 <2e-16 ***
## free.sulfur.dioxide 2.09969   0.05858   35.84 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 24.5 on 1597 degrees of freedom
## Multiple R-squared:  0.4458, Adjusted R-squared:  0.4454
## F-statistic: 1285 on 1 and 1597 DF, p-value: < 2.2e-16

ggplot (dfrm, aes(x=free.sulfur.dioxide, y=total.sulfur.dioxide))+ geom_point(alpha=.2) +
stat_smooth (method=lm, level=0.95)

## `geom_smooth()` using formula 'y ~ x'
```



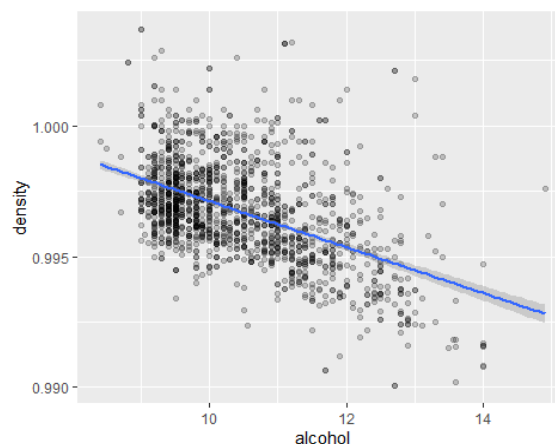
The last regression model is $\text{density} = 1.006e+00 - 8.787e-04 \cdot \text{alcohol}$. So again - big decimal numbers after a point. I investigated this regression model because it is said that density of the wine depends on the percent of alcohol in it, among other things. It turns out that R squared is very low here, only 0.2457 and we can see the dots are being scattered around the line away.

```
reg7<-lm (density~alcohol)
print (summary(reg7))

##
## Call:
## lm(formula = density ~ alcohol)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.0049845 -0.0010951 -0.0002456  0.0008467  0.0073543
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  1.006e+00  4.031e-04  2495.19  <2e-16 ***
## alcohol    -8.787e-04  3.848e-05  -22.84   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.001639 on 1597 degrees of freedom
## Multiple R-squared:  0.2462, Adjusted R-squared:  0.2457
## F-statistic: 521.6 on 1 and 1597 DF, p-value: < 2.2e-16

ggplot (dfrm, aes(x=alcohol, y=density))+ geom_point (alpha=.2) + stat_smooth (method
=lm, level=0.95)

## `geom_smooth()` using formula 'y ~ x'
```



With this, I am done with individual regression models. The last thing I want to see and check is how much variability is accounted by all of the variables together, versus Quality variable. In order to do that, I am going to create a multiple regression model.

Now, in statistics you have a linear regression model and a multiple regression model. Linear model takes in account just one x-variable, whereas the multiple regression model has more x-variables.

```
reg8<-lm (quality~fixed.acidity+volatile.acidity+citric.acid+residual.sugar+chlorides+fre
e.sulfur.dioxide+total.sulfur.dioxide+density+pH+sulphates+alcohol)
print (summary(reg8))

##
## Call:
## lm(formula = quality ~ fixed.acidity + volatile.acidity + citric.acid +
##   residual.sugar + chlorides + free.sulfur.dioxide + total.sulfur.dioxide +
##   density + pH + sulphates + alcohol)
##
## Residuals:
##   Min     1Q   Median     3Q      Max
## -2.68911 -0.36652 -0.04699  0.45202  2.02498
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.197e+01 2.119e+01  1.036  0.3002
## fixed.acidity    2.499e-02 2.595e-02  0.963  0.3357
## volatile.acidity -1.084e+00 1.211e-01 -8.948 < 2e-16 ***
## citric.acid     -1.826e-01 1.472e-01 -1.240  0.2150
## residual.sugar   1.633e-02 1.500e-02  1.089  0.2765
## chlorides       -1.874e+00 4.193e-01 -4.470 8.37e-06 ***
## free.sulfur.dioxide 4.361e-03 2.171e-03  2.009  0.0447 *
## total.sulfur.dioxide -3.265e-03 7.287e-04 -4.480 8.00e-06 ***
## density        -1.788e+01 2.163e+01 -0.827  0.4086
## pH             -4.137e-01 1.916e-01 -2.159  0.0310 *
## sulphates       9.163e-01 1.143e-01  8.014 2.13e-15 ***
## alcohol        2.762e-01 2.648e-02 10.429 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.648 on 1587 degrees of freedom
## Multiple R-squared:  0.3606, Adjusted R-squared:  0.3561
## F-statistic: 81.35 on 11 and 1587 DF, p-value: < 2.2e-16
```


As you can see, at first, R squared isn't particularly high - only 0.3561, which means that there are other factors and variables in real life that are having an effect on quality of red wines, which aren't a part of this sample. Also, not all p-values in the models are lower than 0.05. For example, we would say that variables like Fixed acidity, citric acid, residual sugar and density aren't significant in the model itself. The proposition to this sample is that we find other variables that are found in the red wine production, in order to better the whole model and raise the R squared value.

If I were to delete those variables from the model and only take those which are significant, the R squared value doesn't change almost anything. The only thing we can do, with the data in this sample, is to make a decision tree. But since there is a lot of variables, it doesn't look good and useful.