# 611 Project–Wine Quality Analysis
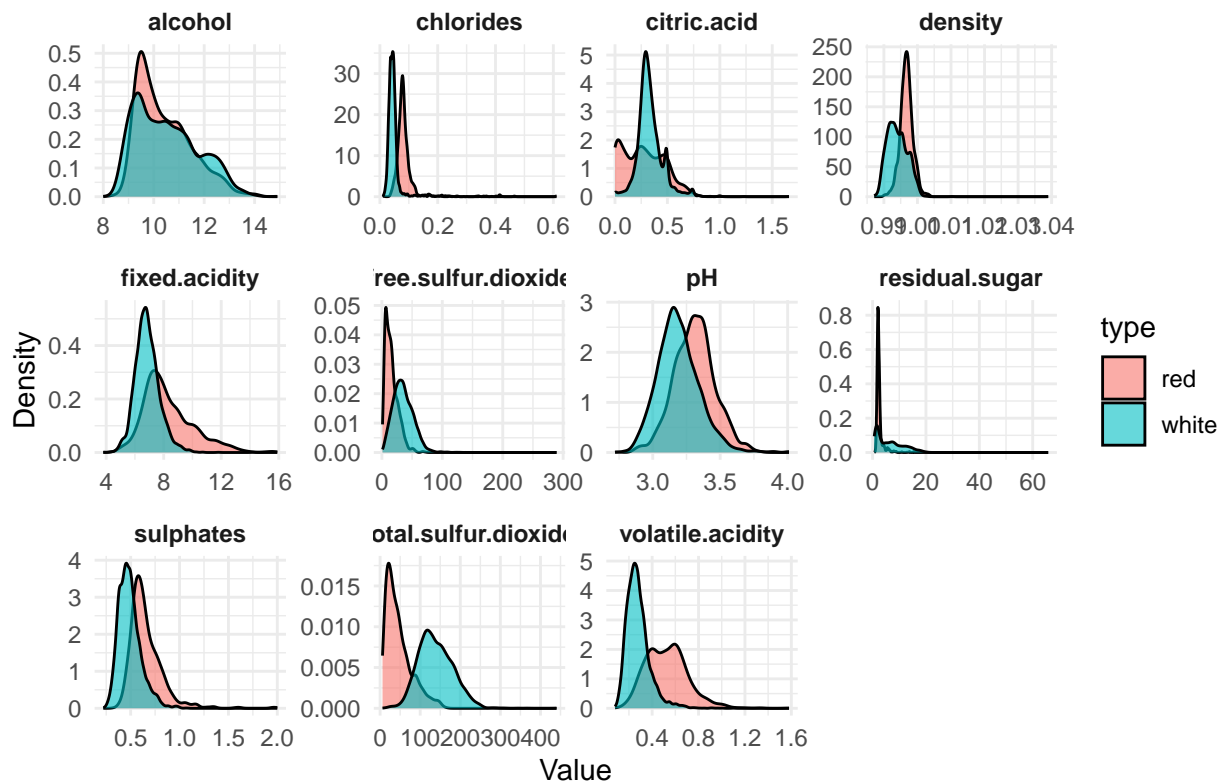
## Zoe Curtis

## 2025-09-30

This analysis will explore a concatenated data set containing chemical components of red and white variants of the Portuguese "Vinho Verde" wine from 2009. It contains chemical components including fixed acidity, volatile acidity, citric acid, residual sugar, chlorides, free sulfur dioxide, total sulfur dioxide, density, pH, sulphates, and alcohol. It also contains a quality score. I am interested in exploring whether wine type (red vs white) differs in their chemical composition, and whether there is a quality difference between the types. I am also interested in using the chemical composition to predict alcohol levels.

The data set was obtained at https://archive.ics.uci.edu/dataset/186/wine+quality. I downloaded the red and white data separately and concatenated them. Fortunately, there is no missing data and the data is fairly clean and ready for analysis. There is more white wine data than red wine data.

```r
# Reshape data
long_wine <- allwine %>%
  pivot_longer(cols = -c(type, quality), names_to = "variable", values_to = "value")

# Plot density (smoothed distribution)
densall<-ggplot(long_wine, aes(x = value, fill = type)) +
  geom_density(
    alpha = 0.6,
    position = "identity"
  ) +
  facet_wrap(~ variable, scales = "free") +
  labs(
    y = "Density",
    x = "Value",
    title = "Distribution of Wine Characteristics by Type"
  ) +
  theme_minimal() +
  theme(strip.text = element_text(face = "bold"))
densall
```

## Distribution of Wine Characteristics by Type



```
save_base_plot(densall, "density_plot.png")
```

```
## pdf
##   2
```

From this, we can see initially that red wines may have a higher pH, sulfates, and volatile acidity than white wines. White wines may have more variation in residual sugar and higher sulfur dioxide than red wines.

We will proceed with a PCA to examine the contributions of chemical components to the variation across the whole data set.

```
results<-prcomp(allwine[-1], center=TRUE, scale=TRUE)
rx <- results$x
rx_df <- as_tibble(rx) %>%
  mutate(row = row_number()) %>%
  pivot_longer(cols = -row, names_to = "pc", values_to = "value") %>%
  mutate(pc = suppressWarnings(as.numeric(gsub("^[^0-9]*", "", pc))))

summary(results)
```
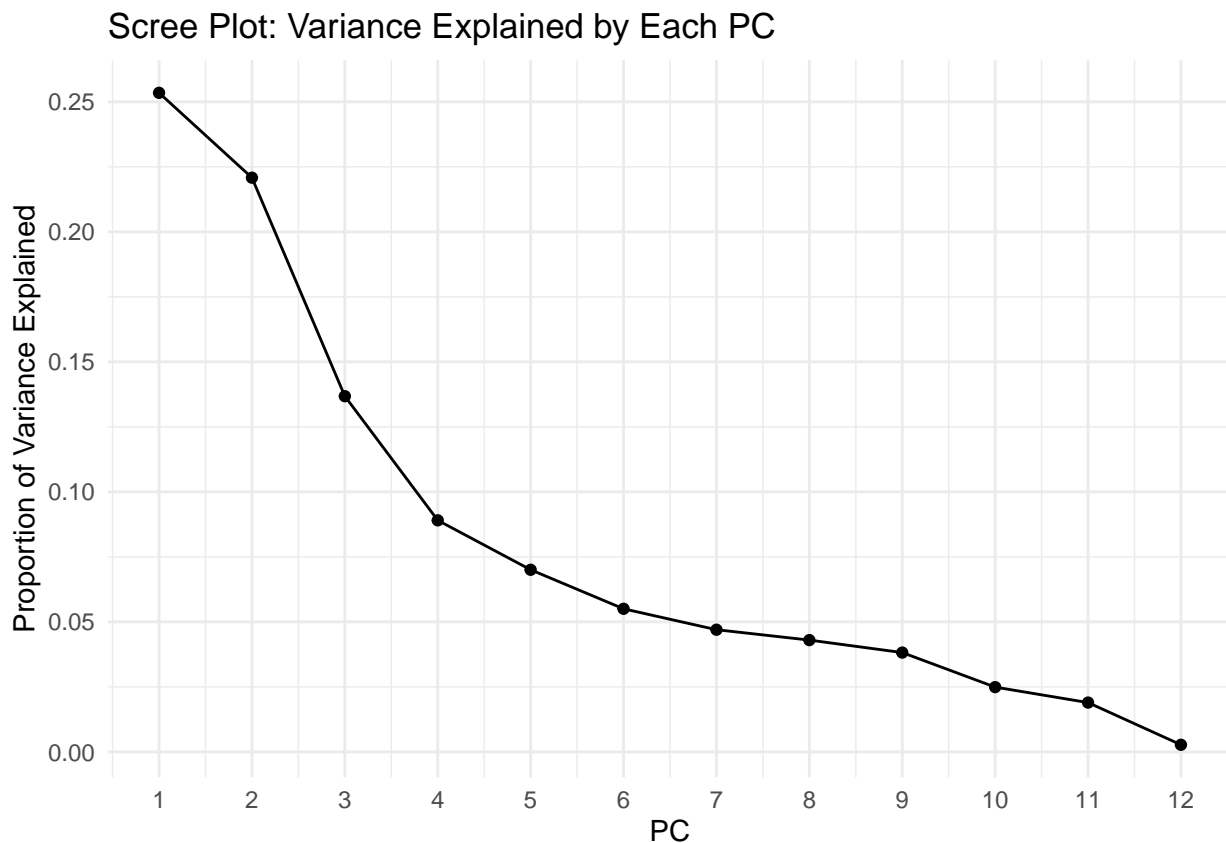
```
## Importance of components:
##                            PC1    PC2    PC3     PC4     PC5     PC6     PC7
## Standard deviation      1.7440 1.6278 1.2812 1.03374 0.91679 0.81265 0.75088
## Proportion of Variance  0.2535 0.2208 0.1368 0.08905 0.07004 0.05503 0.04699
## Cumulative Proportion   0.2535 0.4743 0.6111 0.70013 0.77017 0.82520 0.87219
##                            PC8    PC9   PC10    PC11    PC12
## Standard deviation      0.7183 0.6770 0.54682 0.47706 0.18107
## Proportion of Variance  0.0430 0.0382 0.02492 0.01897 0.00273
## Cumulative Proportion   0.9152 0.9534 0.97830 0.99727 1.00000
```

```
var_explained <- results$sdev^2 / sum(results$sdev^2)
scree_df <- data.frame(
  PC = 1:length(var_explained),
  Variance = var_explained
)

scre<-ggplot(scree_df, aes(x = PC, y = Variance)) +
  geom_line() +
  geom_point() +
  labs(
    title = "Scree Plot: Variance Explained by Each PC",
    y = "Proportion of Variance Explained"
  ) +
  scale_x_continuous(breaks = 1:ncol(allwine[-1])) +
  theme_minimal()
scre
```

## Scree Plot: Variance Explained by Each PC



```
save_base_plot(scre, "screeplot.png")
```

```
## pdf
##   2
```

```
scores <- as.data.frame(results$x[,1:2])
scores$type <- allwine$type

pcplot<-ggplot(scores, aes(x = PC1, y = PC2, color=type)) +
  geom_point(size = 1.5, alpha=0.3) +
  theme_minimal()
```

```
pcplot
save_base_plot(pcplot, "pcplot.png")
```

```
## pdf
##   2
```

I was planning to do kmeans if there was obvious clustering that appears, but as seen in this plot, there is overlap and not tight separate clusters in the data. Red and white wines are distinct in terms of PC1, but not distinct in terms of PC2. There is much more overlap on the Y axis between the types. Thus PC1 captures more of the chemical differences in the types of wines.

Let's continue by investigating the loadings of the first two PCs:

```
pcs<-as.data.frame(results$rotation)[,1:2]
pcs
```

```
##                            PC1        PC2
## fixed.acidity         0.25692873  0.2618431
## volatile.acidity      0.39493118  0.1051983
## citric.acid          -0.14646061  0.1440935
## residual.sugar       -0.31890519  0.3425850
## chlorides             0.31344994  0.2697701
## free.sulfur.dioxide  -0.42269137  0.1111788
## total.sulfur.dioxide -0.47441968  0.1439475
## density               0.09243753  0.5549205
## pH                    0.20806957 -0.1529219
## sulphates             0.29985192  0.1196342
## alcohol               0.05892408 -0.4927275
## quality              -0.08747571 -0.2966009
```
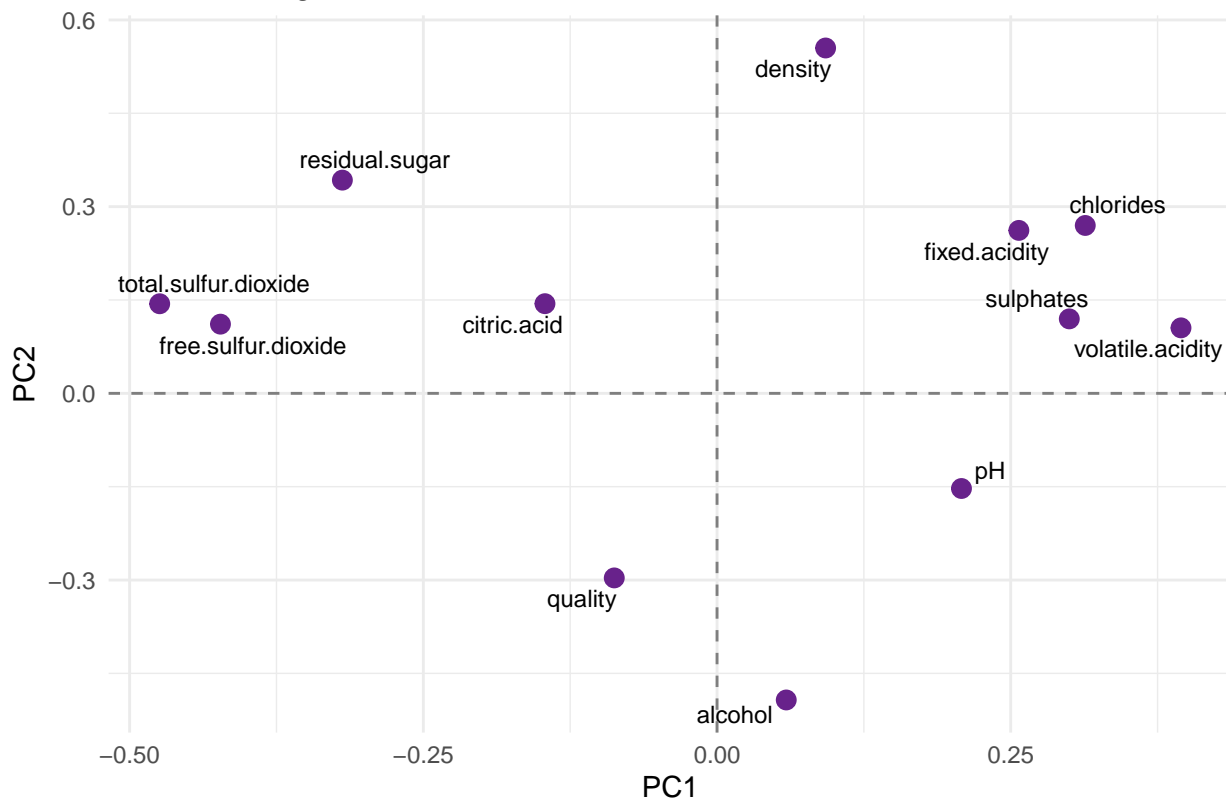
It appears that PC1 encompasses wines with higher volatile acidity, low sulfur dioxide, higher sulphates, and lower residual sugar. On the other hand, PC2 encompasses wines with higher citric acid, residual sugar, sulfur dioxide, density, and lower alcohol.

As red wines are fermented, this might make sense for them to have higher sulphates and for white wines to have higher sugars(white wines lie on the negative side, red on the positive side of the PC1 plot).

```
loadings <- as.data.frame(results$rotation[, 1:2])
loadings$variable <- rownames(loadings)

pcword<-ggplot(loadings, aes(x = PC1, y = PC2, label = variable)) +
  geom_point(color = "darkorchid4", size = 3) +
  geom_text_repel(aes(label = variable), size = 3) +
  labs(title = "PCA Loadings: Contribution of Variables to PC1 & PC2") +
  theme_minimal() +
  geom_hline(yintercept = 0, linetype = "dashed", color = "gray50") +
  geom_vline(xintercept = 0, linetype = "dashed", color = "gray50")
pcword
```

## PCA Loadings: Contribution of Variables to PC1 & PC2



```
save_base_plot(pcword, "pcword.png")
```

```
## pdf
##   2
```

This plot is a more intuitive way of understanding the contributions of chemical compounds to the two PCs. We see PC1 has positive contributions from density, fixed acidity, chlorides, sulphates, and volatile acidity. PC2 on the other hand has positive contributions from every factor except quality, alcohol, and pH.
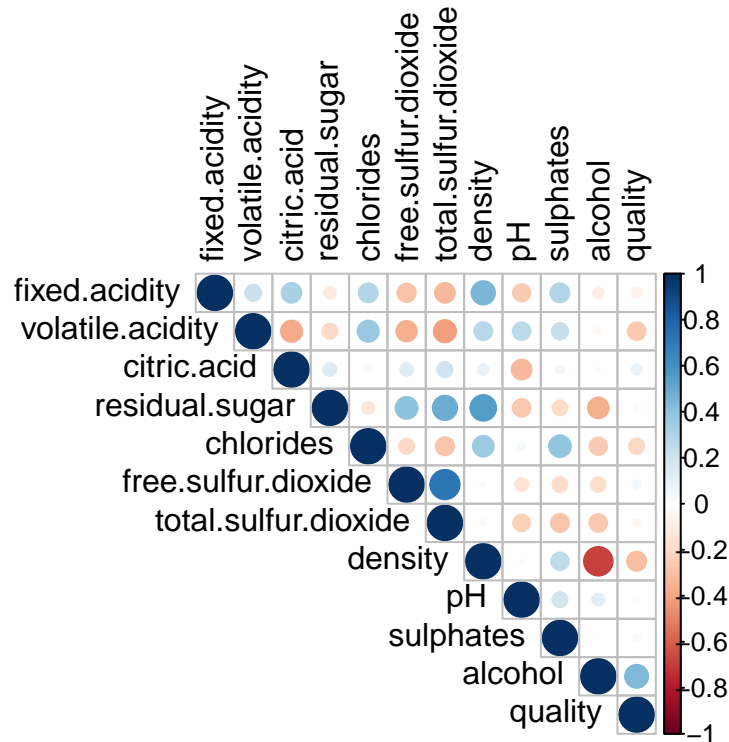
I am interested in how chemical components may differ by type, so I will look into correlation plots:

```
# Correlation matrix
mycor <- cor(allwine[,-1])
dir.create("figures", showWarnings = FALSE)
png("figures/corplotall.png", width = 8, height = 6, units = "in", res = 150)
par(oma = c(0, 0, 3, 0))
corrplot(mycor, type = "upper", tl.col = "black")
title("Correlation Matrix for All Wine", outer=TRUE, cex.main = 1.4)
dev.off()
```

```
## pdf
##   2
```

```
par(oma = c(0, 0, 3, 0))
corrplot(mycor, type = "upper", tl.col = "black")
title("Correlation Matrix for All Wine", outer=TRUE, cex.main = 1.4)
```
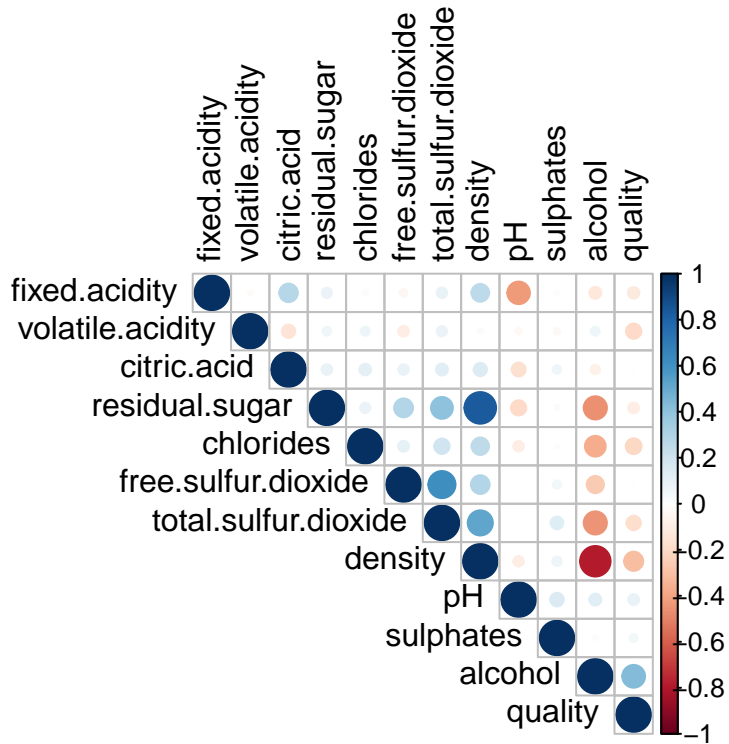
# Correlation Matrix for All Wine



```
mycor2<-cor(white[,-1])
dir.create("figures", showWarnings = FALSE)
png("figures/corplotwhite.png", width = 8, height = 6, units = "in", res = 150)
par(oma = c(0, 0, 3, 0))
corrplot(mycor2, type = "upper", tl.col = "black")
title("Correlation Matrix for White Wine", outer=TRUE, cex.main = 1.4)
dev.off()
```

```
## pdf
##   2
```

```
par(oma = c(0, 0, 3, 0))
corrplot(mycor2, type = "upper", tl.col = "black")
title("Correlation Matrix for White Wine", outer=TRUE, cex.main = 1.4)
```
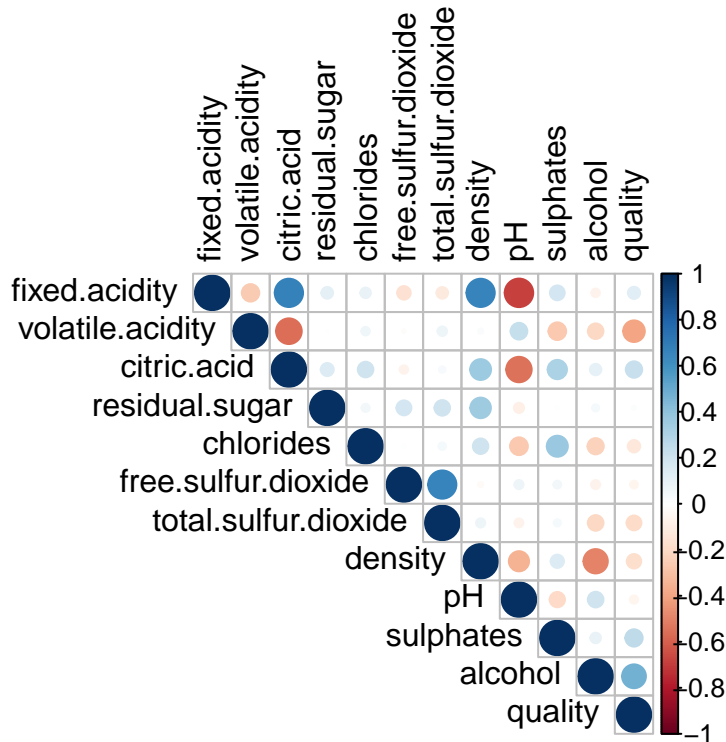
# Correlation Matrix for White Wine



```
mycor3<-cor(red[,-1])
dir.create("figures", showWarnings = FALSE)
png("figures/corplotred.png", width = 8, height = 6, units = "in", res = 150)
par(oma = c(0, 0, 3, 0))
corrplot(mycor3, type = "upper", tl.col = "black")
title("Correlation Matrix for Red Wine", outer=TRUE, cex.main = 1.4)
dev.off()
```

```
## pdf
##   2
```

```
par(oma = c(0, 0, 3, 0))
corrplot(mycor3, type = "upper", tl.col = "black")
title("Correlation Matrix for Red Wine", outer=TRUE, cex.main = 1.4)
```

# Correlation Matrix for Red Wine



In the overall dataset, we can see that alcohol and density are highly NEGATIVELY correlated, makes sense since the higher the alcohol content, the lower the density. Additionally, total sulfur dioxide and free sulfur dioxide also highly correlated, which makes sense because they are different measures of the same thing.

Yet when we split by wine type, we see some different correlations. In the white wine correlation plot, we also see that residual sugar and density are highly correlated, and that free sulfur dioxide and total sulfur dioxide and moderately positively correlated.

In the red wine correlation plot, fixed acidity and pH are highly negatively correlated, which makes sense, but it is curious that this does not appear in the white wine data set. The free and total sulfur dioxide relationship remains, yet the density-alcohol correlation is smaller, there is a stronger positive relationship between fixed acidity and citric acid, and a positive relationship between fixed acidity and density, which do not really show up in the white wine data.

This indicates that the underlying relationship between variables is not the same for red and white wine, so I will do a PCA for red and white wine separately now to see if chemical components contribute to variation differently for the two types of wine:

```
whiteresults<-prcomp(white[-1], center=TRUE, scale=TRUE)
whiterx <- whiteresults$x
whiterx_df <- as_tibble(whiterx) %>%
  mutate(row = row_number()) %>%
  pivot_longer(cols = -row, names_to = "pc", values_to = "value") %>%
  mutate(pc = suppressWarnings(as.numeric(gsub("^[^0-9]*", "", pc))))

summary(whiteresults)

## Importance of components:
##                          PC1    PC2    PC3     PC4     PC5     PC6     PC7
## Standard deviation    1.8294 1.2594 1.1710 1.04157 0.98756 0.96890 0.8771
```
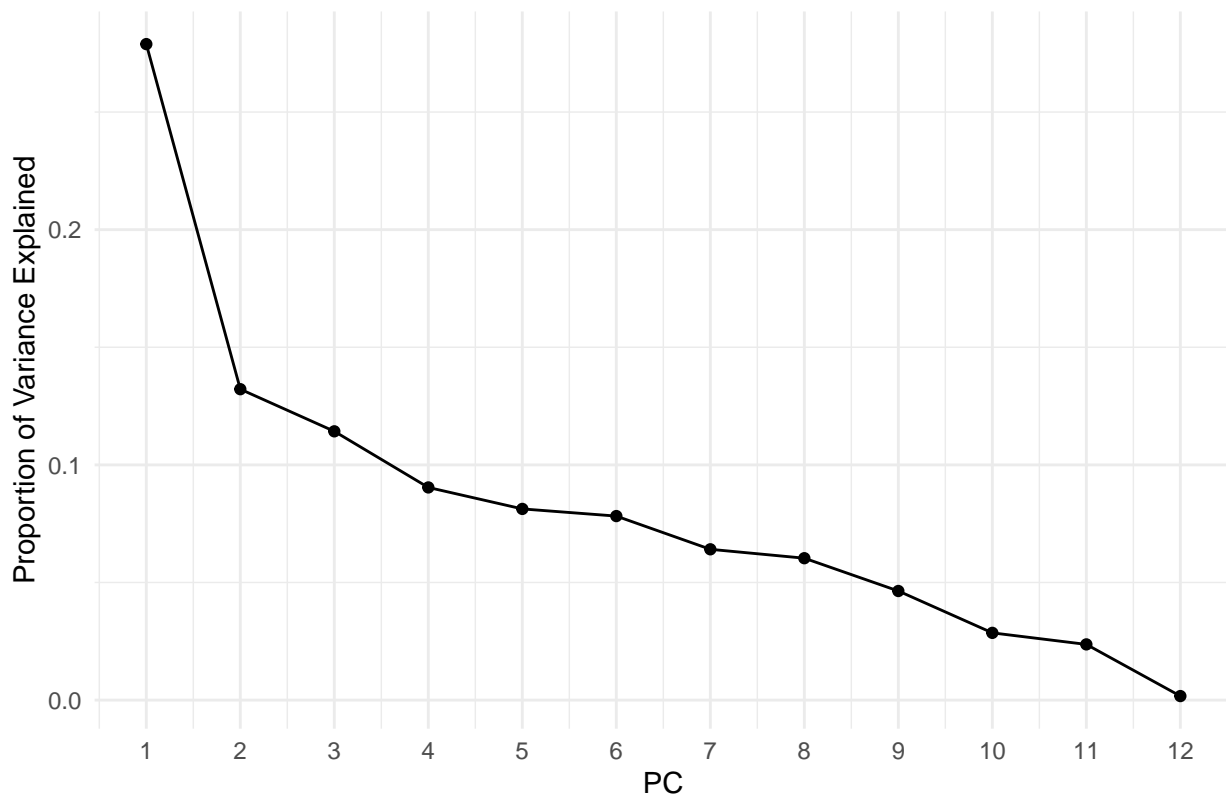
```
## Proportion of Variance 0.2789 0.1322 0.1143 0.09041 0.08127 0.07823 0.0641
## Cumulative Proportion  0.2789 0.4111 0.5253 0.61573 0.69701 0.77524 0.8393
##                                PC8     PC9    PC10    PC11    PC12
## Standard deviation     0.85082 0.74599 0.58561 0.53302 0.14307
## Proportion of Variance 0.06032 0.04638 0.02858 0.02368 0.00171
## Cumulative Proportion  0.89967 0.94604 0.97462 0.99829 1.00000
```

```r
var_explained <- whiteresults$sdev^2 / sum(whiteresults$sdev^2)
whitescree_df <- data.frame(
  PC = 1:length(var_explained),
  Variance = var_explained
)

scre<-ggplot(whitescree_df, aes(x = PC, y = Variance)) +
  geom_line() +
  geom_point() +
  labs(
    title = "White Wine Scree Plot: Variance Explained by Each PC",
    y = "Proportion of Variance Explained"
  ) +
  scale_x_continuous(breaks = 1:ncol(white[-1])) +
  theme_minimal()
scre
```



White Wine Scree Plot: Variance Explained by Each PC

```r
save_base_plot(scre, "white.screeplot.png")
```

```
## pdf
##   2
```

```r
scores <- as.data.frame(whiteresults$x[,1:2])
scores$type <- white$type

pcplot<-ggplot(scores, aes(x = PC1, y = PC2)) +
  geom_point(size = 1.5, alpha=0.3, color="#99CCFF") +
  theme_minimal()
pcplot
save_base_plot(pcplot, "white.pcplot.png")
```

```
## pdf
##   2
```

```r
whitepcs<-as.data.frame(whiteresults$rotation)[,1:2]
whitepcs
```

```
##                             PC1         PC2
## fixed.acidity        -0.15690447  0.56066866
## volatile.acidity     -0.02428722  0.01606694
## citric.acid          -0.13294430  0.28938115
## residual.sugar       -0.40605288 -0.03882402
## chlorides            -0.21754400  0.03691144
## free.sulfur.dioxide  -0.27471931 -0.34554881
## total.sulfur.dioxide -0.39044148 -0.27232605
## density              -0.50129557 -0.01773344
## pH                    0.13003701 -0.56714503
## sulphates            -0.03364168 -0.24826266
## alcohol               0.44279498  0.01698188
## quality               0.22713722 -0.14603134
```

```r
whiteloadings <- as.data.frame(whiteresults$rotation[, 1:2])
whiteloadings$variable <- rownames(whiteloadings)
```

```r
redresults<-prcomp(red[-1], center=TRUE, scale=TRUE)
redrx <- redresults$x
redrx_df <- as_tibble(rx) %>%
  mutate(row = row_number()) %>%
  pivot_longer(cols = -row, names_to = "pc", values_to = "value") %>%
  mutate(pc = suppressWarnings(as.numeric(gsub("^[^0-9]*", "", pc))))

summary(redresults)
```
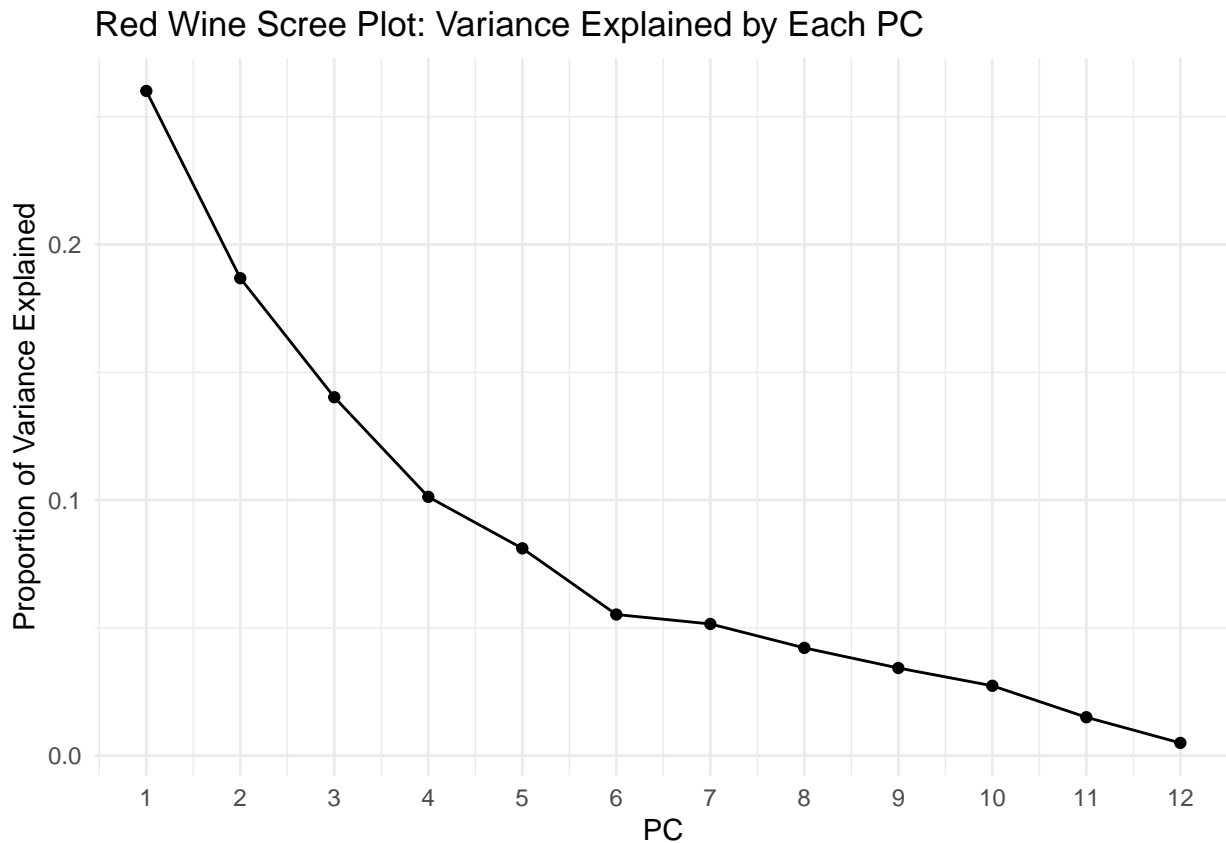
```
## Importance of components:
##                           PC1    PC2    PC3    PC4     PC5     PC6     PC7
## Standard deviation     1.7667 1.4973 1.2973 1.1023 0.98654 0.81400 0.78633
## Proportion of Variance 0.2601 0.1868 0.1402 0.1013 0.08111 0.05522 0.05153
## Cumulative Proportion  0.2601 0.4469 0.5872 0.6884 0.76952 0.82474 0.87626
##                           PC8    PC9   PC10    PC11    PC12
## Standard deviation     0.71125 0.64133 0.57264 0.42452 0.24396
## Proportion of Variance 0.04216 0.03428 0.02733 0.01502 0.00496
## Cumulative Proportion  0.91842 0.95270 0.98002 0.99504 1.00000
```

```r
var_explained <- redresults$sdev^2 / sum(redresults$sdev^2)
scree_df <- data.frame(
  PC = 1:length(var_explained),
  Variance = var_explained
)
```

```r
scre<-ggplot(scree_df, aes(x = PC, y = Variance)) +
  geom_line() +
  geom_point() +
  labs(
    title = "Red Wine Scree Plot: Variance Explained by Each PC",
    y = "Proportion of Variance Explained"
  ) +
  scale_x_continuous(breaks = 1:ncol(red[-1])) +
  theme_minimal()
scre
```



Red Wine Scree Plot: Variance Explained by Each PC

```r
save_base_plot(scre, "red.screeplot.png")
```

```
## pdf
##   2
```

```r
redscores <- as.data.frame(redresults$x[,1:2])
redscores$type <- red$type

pcplot<-ggplot(redscores, aes(x = PC1, y = PC2)) +
  geom_point(size = 1.5, alpha=0.3, color="#99000F") +
  theme_minimal()
pcplot
save_base_plot(pcplot, "red.pcplot.png")
```

```
## pdf
##   2
```

```
redpcs<-as.data.frame(redresults$rotation)[,1:2]
redpcs
```

```
##                               PC1          PC2
## fixed.acidity          0.487883358  0.004173212
## volatile.acidity      -0.265128984 -0.338967858
## citric.acid            0.473335467  0.137358104
## residual.sugar         0.139154423 -0.167736336
## chlorides              0.197426792 -0.189788185
## free.sulfur.dioxide   -0.045880713 -0.259483136
## total.sulfur.dioxide   0.004066746 -0.363971374
## density                0.370301191 -0.330780789
## pH                    -0.432720849  0.065440145
## sulphates              0.254535354  0.109333620
## alcohol               -0.073176777  0.502708647
## quality                0.112488776  0.473166214
```

```r
redloadings <- as.data.frame(redresults$rotation[, 1:2])
redloadings$variable <- rownames(loadings)
```

```r
red_loadings <- as.data.frame(redresults$rotation[, 1:2])
red_loadings$variable <- rownames(redloadings)
red_loadings$type <- "Red"

white_loadings <- as.data.frame(whiteresults$rotation[, 1:2])
white_loadings$variable <- rownames(whiteloadings)
white_loadings$type <- "White"

# combine
combined_loadings <- rbind(red_loadings, white_loadings)
overlay_plot <- ggplot(combined_loadings,
                    aes(x = PC1, y = PC2, label = variable, color = type)) +
  geom_point(size = 3) +
  geom_text_repel(size = 3, show.legend = FALSE) +
  labs(title = "PCA Loadings for Red vs White Wines",
       subtitle = "Comparison of variable contributions to PC1 & PC2") +
  theme_minimal() +
  geom_hline(yintercept = 0, linetype = "dashed", color = "gray50") +
  geom_vline(xintercept = 0, linetype = "dashed", color = "gray50") +
  scale_color_manual(values = c("Red" = "#99000F", "White" = "#99CCFF"))

overlay_plot
```
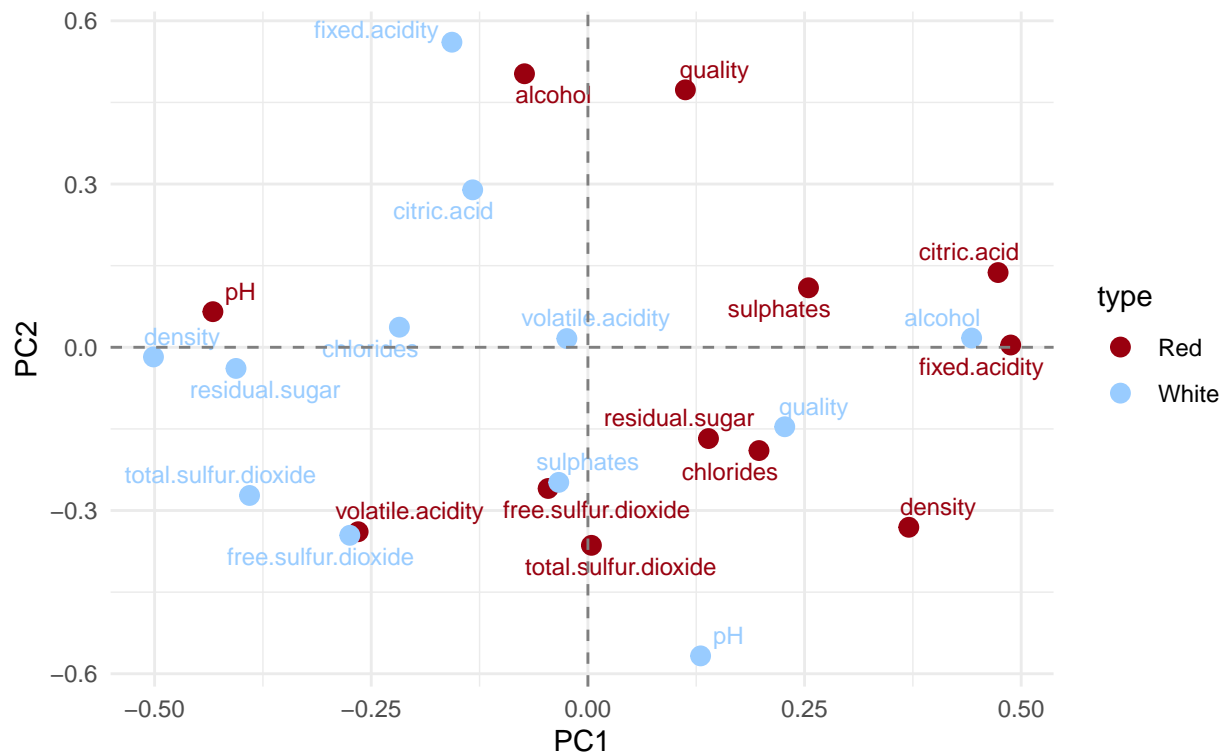
## PCA Loadings for Red vs White Wines
### Comparison of variable contributions to PC1 & PC2



```
save_base_plot(overlay_plot, "overlay_pcplot.png")
```

```
## pdf
##   2
```

I am curious if we can use the chemical properties to predict alcohol content. I will do a LASSO regression, using 70% of the data to train the model, and test it on the other 30% of the data.

```
set.seed(828)
y<-allwine$alcohol
x<-model.matrix(alcohol ~. -type -quality, data=allwine)
n<-nrow(allwine)
train_idx <- sample(1:n, 0.7 * n)
xtrain<-x[train_idx, ]
ytrain<-y[train_idx]
xtest<-x[-train_idx, ]
ytest<-y[-train_idx]

#use training data for model
cv_model<-cv.glmnet(xtrain,ytrain)
best_lambda <- cv_model$lambda.min
#make predictions on test
y_pred_test <- predict(cv_model, s = best_lambda, newx = xtest)
#calculate R^2 from test
sse <- sum((ytest - y_pred_test)^2)
sst <- sum((ytest - mean(ytest))^2)
rsq_test <- 1 - sse/sst
print(rsq_test)
```

```
## [1] 0.7192804
```

An $R^2$ value of 0.823 indicates that about 83% of the variation in alcohol content can be predicted by the other numeric chemical properties.

```r
set.seed(828)
y<-white$alcohol
x <- as.matrix(white %>% select(-alcohol, -type))
n<-nrow(white)
train_idx <- sample(1:n, 0.7 * n)
xtrain<-x[train_idx, ]
ytrain<-y[train_idx]
xtest<-x[-train_idx, ]
ytest<-y[-train_idx]

#use training data for model
cv_model<-cv.glmnet(xtrain,ytrain)
best_lambda <- cv_model$lambda.min
#make predictions on test
y_pred_test <- predict(cv_model, s = best_lambda, newx = xtest)
#calculate R^2 from test
sse <- sum((ytest - y_pred_test)^2)
sst <- sum((ytest - mean(ytest))^2)
rsq_test <- 1 - sse/sst
print(rsq_test)
```

```
## [1] 0.9050433
```

```r
set.seed(828)
y<-red$alcohol
x <- as.matrix(red %>% select(-alcohol, -type))
n<-nrow(red)
train_idx <- sample(1:n, 0.7 * n)
xtrain<-x[train_idx, ]
ytrain<-y[train_idx]
xtest<-x[-train_idx, ]
ytest<-y[-train_idx]

#use training data for model
cv_model<-cv.glmnet(xtrain,ytrain)
best_lambda <- cv_model$lambda.min
#make predictions on test
y_pred_test <- predict(cv_model, s = best_lambda, newx = xtest)
#calculate R^2 from test
sse <- sum((ytest - y_pred_test)^2)
sst <- sum((ytest - mean(ytest))^2)
rsq_test <- 1 - sse/sst
print(rsq_test)
```
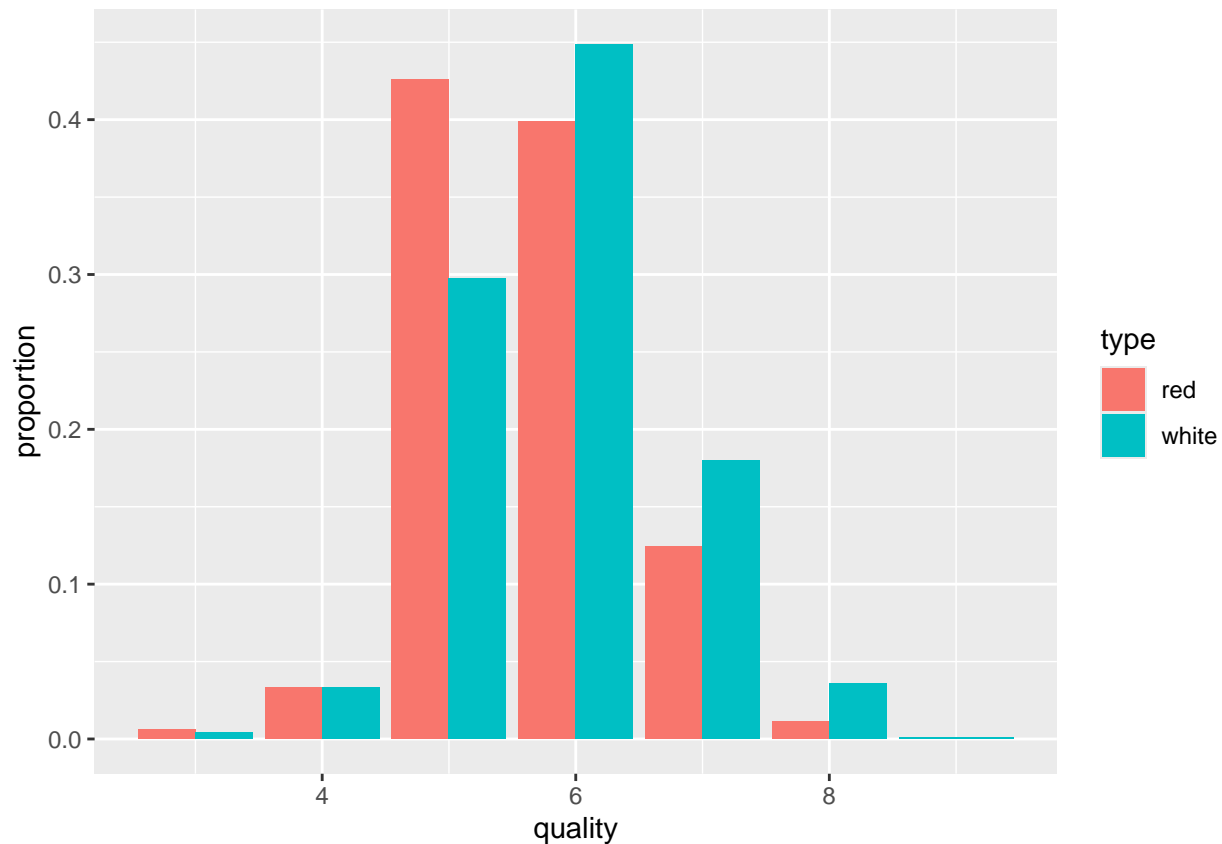
```
## [1] 0.6720744
```

In the overall data set, about 82% of the variation in alcohol data can be explained by the chemical components. However after seeing that chemical component relationships differ for type of wine and once we split by type, white wine has almost 90% of the variation in its alcohol data explained by the chemical components, but red wine has only 67% of the variation in its alcohol data explained by the chemical components.

I am now interested in investigating quality. I will look into if quality differs between red and white wines.

```r
qualityprop<-ggplot(allwine, aes(x = quality, fill = type)) +
  geom_bar(aes(y = after_stat(prop), group = type), position = "dodge")+
  ylab("proportion")
qualityprop
```



```r
save_base_plot(qualityprop, "qualityprop.png")
```

```
## pdf
##   2
```

There does not appear to be a huge difference in the quality ratings between red and white wines.

I will now do a regression (similar to what was done for alcohol) to see what chemical components are most influential on quality score. We do a LASSO to account for contributions when all predictors are included and because of multicollinearity.

```r
set.seed(828)
y<-allwine$quality
x<-model.matrix(quality ~., data=allwine)
n<-nrow(allwine)
train_idx <- sample(1:n, 0.7 * n)
xtrain<-x[train_idx, ]
ytrain<-y[train_idx]
xtest<-x[-train_idx, ]
ytest<-y[-train_idx]

#use training data for model
cv_model<-cv.glmnet(xtrain,ytrain)
```

```r
coef_best <- coef(cv_model, s = "lambda.min")
nonzero <- round(coef_best[coef_best[,1] != 0, , drop=FALSE],3)
print(nonzero)
```

```
## 13 x 1 sparse Matrix of class "dgCMatrix"
##                              s1
## (Intercept)            114.332
## typewhite               -0.407
## fixed.acidity            0.093
## volatile.acidity        -1.495
## citric.acid             -0.085
## residual.sugar           0.066
## chlorides               -0.575
## free.sulfur.dioxide      0.004
## total.sulfur.dioxide    -0.001
## density               -113.496
## pH                       0.527
## sulphates                0.598
## alcohol                  0.210
```

```r
best_lambda <- cv_model$lambda.min
#make predictions on test
y_pred_test <- predict(cv_model, s = best_lambda, newx = xtest)
#calculate adj R^2 from test
p <- length(nonzero) - 1   # exclude intercept
n_test <- length(ytest)
# predictions
y_pred_test <- predict(cv_model, s = best_lambda, newx = xtest)

# compute R^2
sse <- sum((ytest - y_pred_test)^2)
sst <- sum((ytest - mean(ytest))^2)
rsq_test <- 1 - sse/sst

rsq_adj <- 1 - (1 - rsq_test) * (n_test - 1) / (n_test - p - 1)
rsq_adj
```

```
## [1] 0.2975302
```

The model does not capture much of the variation in the quality data ($R^2$=0.298), and most predictors have small coefficients at the optimal penalty (lambda). However, density has a coefficient of -113.496, thus it has stronger negative influence on the quality response. Higher density wines thus have lower quality ratings.

This exploration gives us an insight into how chemical components of wine are related to type of wine, alcohol, and quality. Further research could be done into if clusters can be found within red and white wines, rather than considering all of the wines together. Clustering could possibly tease apart different types of Vinho Verde, and it would be interesting to see what components are influential in differentiating them.