

# SDS 323: Exercises 3

Aaron Grubbs

Kaushik Koirala

Khue Tran

Matthew Tran

## Problem 1: Predictive Model Building

### Problem Overview:

Given the data on commercial rental properties, this report aims to generate the best predictive model for price. Using the generated predictive model, this report will additionally aim to determine the change in rental income per square foot given a building's green certification while holding the other features of the building constant.

### Data Analysis and Process

In order to create the best model, this report aimed at starting with a simple intuitive linear model and then doing a stepwise selection process to select the proper features and interaction relations between the features.

Using the green\_rating variable for the green certification requirement makes creating a predictive model simpler instead of dealing with the 2 categories. The other variables were then separated by their significance in predicting rent when paired with green\_rating. Afterwards the variables were ranked in order of r2 to establish the best variables to use in the stepwise process of model improvement. The cluster rent category was removed since that data is based on the rent data. Amenities and Electricity costs both had r2 values in the tenths place and subsequently improved the model based on the decrease in rsme value. Systematically other variables were added and removed, but no other variables were able to reduce the rsme value further.

In code that looks like this:

```
lm_medium = lm(Rent ~ green_rating + amenities + Electricity_Costs, data=greenbuildings)
```

For the stepwise function other variables were included and experimented with such as size and leasing rate. In code that looks like this:

```
lm_step = step(lm_medium, scope=~(. + cd_total_07 + size + class_a + leasing_rate + class_b)^3)
```

### Results

Here is the output of the stepwise selection:

```
## lm(formula = Rent ~ green_rating + amenities + Electricity_Costs +
##       cd_total_07 + class_a + leasing_rate + size + class_b + Electricity_Costs:cd_total_07 +
##       Electricity_Costs:size + cd_total_07:size + leasing_rate:size +
##       amenities:class_b + Electricity_Costs:class_a + class_a:size +
##       amenities:cd_total_07 + Electricity_Costs:leasing_rate +
##       amenities:leasing_rate + green_rating:size + size:class_b +
##       leasing_rate:class_b + Electricity_Costs:class_a:size + Electricity_Costs:leasing_rate:size +
##       amenities:leasing_rate:class_b, data = greenbuildings)
```

Ultimately, it selected 25 variables for use in the optimized model.

Here are the RMSE values for the original linear model (on the left) and the output of the stepwise selection (on the right), calculated across 100 different train test splits.

```
##      V1      V2
## 16.42058 17.65229
```

## Conclusion

The initial simpler model had the lower RMSE than the result of the stepwise feature selection. It seemed that adding more features and interactions worsened the model and its accuracy as the first three variables were sufficient enough. The coefficients of the first simple model are:

```
##      (Intercept)    green_rating    amenities Electricity_Costs
## 3.8199654     0.4054184    3.5906206    735.3183503
```

Interpreting the green\_rating coefficient, “green” certified properties seemed to improve rent by around 0.4054184 dollars per square foot.

## Problem 2: What causes what?

1. The question that the researchers were looking at was whether increasing the number of police would reduce the rate of crime in a given city. Intuitively, it might make sense to approach the problem as an independent probability by varying the number of cops and observing fluctuations in crime rates. However, an experimental construct that randomly changes the number of cops on random days is not practical. Another problem arises from the fact that crime rates also affect the number of police. So places with higher crime rates would naturally have more cops on the street at a given time, making it difficult to see the isolated effect of just increasing cops. A regression model that takes in crime rate and number of cops would not be able to account for such interactions.

### EFFECT OF POLICE ON CRIME

TABLE 2  
TOTAL DAILY CRIME DECREASES ON HIGH-ALERT DAYS

	(1)	(2)
High Alert	-7.316* (2.877)	-6.046* (2.537)
Log(midday ridership)		17.341** (5.309)
R <sup>2</sup>	.14	.17

Figure 1: The dependent variable is the daily total number of crimes in D.C. This table present the estimated coefficients and their standard errors in parenthesis. The first column refers to a model where the only variable used in the High Alert dummy whereas the model in column (2) controls form the METRO ridership. \* refers to a significant coefficient at the 5% level, \*\* at the 1% level.

2. So the best setting for looking at the correlation between police and crime would control for the positive feedback loop. The researchers accomplished this by collecting data when the terrorist alert system levels are high in D.C., when more cops are put on the street for reasons unrelated to street crime. As summarized in table 2, two models were fitted: the one in the first column only has the dummy variable

High Alert while the model in the second column included a term for ridership. For both models, the High Alert variable have negative coefficients with similar standard errors and is significant at the 5% level. The coefficient from the first model implies that on a high alert day, where there are more cops on the streets, the daily number of crimes in D.C. decreases by around 7. Similarly, the second model coefficients imply that when Metro ridership increases by 10 (on a log scale), about 17 more crimes are committed, and 6 less crimes are committed on a high alert day. So from the results of table 2, increased number of cops on high alert days does decrease the number of crimes. The R-squared value for both models were relative low (0.14 and 0.17 respectively), this could indicate that a linear fit might not be the best way to model the correlation bewteen police and crime, but there was still significant decrease in crime with increased number of cops.

3. It was unknown if the lower crime rates when the threat level was orange was due to the increase in cops on the street or because there was an increase terrorist threat. The reduce crime could be caused by victims and perpetrators being more cautious about being outside during increased terrorist threats. Measuring the Metro ridership measures the general street traffic in DC during those times, which was shown to be relatively unaffected. This was another way to ensure that the correlations observed were just from the number of police.

TABLE 4  
REDUCTION IN CRIME ON HIGH-ALERT DAYS: CONCENTRATION ON THE NATIONAL MALL

	Coefficient (Robust)	Coefficient (HAC)	Coefficient (Clustered by Alert Status and Week)
High Alert × District 1	-2.621** (.044)	-2.621* (1.19)	-2.621* (1.225)
High Alert × Other Districts	-.571 (.455)	-.571 (.366)	-.571 (.364)
Log(midday ridership)	2.477* (.364)	2.477** (.522)	2.477** (.527)
Constant	-11.058** (4.211)	-11.058 (5.87)	-11.058 <sup>+</sup> (5.923)

Figure 2: The dependent variable is the daily total number of crimes in D.C. District 1 refers to a dummy variable associated with crime incidents in the first police district area. This table present the estimated coefficients and their standard errors in parenthesis.\* refers to a significant coefficient at the 5% level, \*\* at the 1% level.

4. The analysis is further explored in table 4, looking at the effects of dummy variable High Alert and log(ridership) in different districts of D.C. by introducing interaction terms. The model from the first column has negative coefficients for interactions between High Alert and District 1, and High Alert and Other Districts, but only the first was significant. So the effect of increasing the number of police (as seen on high alert days) significantly decreases number of total crimes in District 1. The interaction coefficients tell us that on a high terrorist alert day, about 10 less crimes are committed in District 1 (-7.316 - 2.621) and for other districts, about 7 less crimes are committed. As for the ridership coefficient, the opposite correlation is seen where 2 more crimes are recorded when there are 10 times more rider, consistent with the results of table 2.

## Problem 3: Clustering and PCA

### Problem Overview:

The data for this problem contains information on 11 different chemical properties of 6500 different bottles of wine, as well as two additional classification variables of color (red or white) and quality (judged on a 1-10 scale).

### Data and Analysis Process:

Our goal in this problem is to perform dimensionality reduction on the data: principal components analysis (PCA) and clustering. From there, we must summarize and see if our results can distinguish the red and white wines as well as the quality of the wines from our own intuition.

In order to solve this problem, we first need to perform PCA, look at the loadings and principal components of each summary, and form a conclusion. We also need to perform k-means clustering, and check if the clusters can distinctly distinguish our data in terms of quality and color.

After analysis has been done, we compare our models to see which one makes the most sense to us, and see if either analysis did a better job in distinguishing the reds and whites as well as the quality.

```
wine = read.csv("https://raw.githubusercontent.com/jgscott/SDS323/master/data/wine.csv", header=TRUE)
# summary of data
summary(wine)

## fixed.acidity      volatile.acidity    citric.acid      residual.sugar
##  Min.   : 3.800   Min.   :0.0800   Min.   :0.0000   Min.   : 0.600
##  1st Qu.: 6.400   1st Qu.:0.2300   1st Qu.:0.2500   1st Qu.: 1.800
##  Median : 7.000   Median :0.2900   Median :0.3100   Median : 3.000
##  Mean   : 7.215   Mean   :0.3397   Mean   :0.3186   Mean   : 5.443
##  3rd Qu.: 7.700   3rd Qu.:0.4000   3rd Qu.:0.3900   3rd Qu.: 8.100
##  Max.   :15.900   Max.   :1.5800   Max.   :1.6600   Max.   :65.800
## chlorides          free.sulfur.dioxide total.sulfur.dioxide      density
##  Min.   :0.00900   Min.   : 1.00     Min.   : 6.0      Min.   :0.9871
##  1st Qu.:0.03800   1st Qu.: 17.00    1st Qu.: 77.0     1st Qu.:0.9923
##  Median :0.04700   Median : 29.00    Median :118.0     Median :0.9949
##  Mean   :0.05603   Mean   : 30.53    Mean   :115.7     Mean   :0.9947
##  3rd Qu.:0.06500   3rd Qu.: 41.00    3rd Qu.:156.0     3rd Qu.:0.9970
##  Max.   :0.61100   Max.   :289.00    Max.   :440.0     Max.   :1.0390
## pH                 sulphates        alcohol       quality      color
##  Min.   :2.720     Min.   :0.2200   Min.   : 8.00   Min.   :3.000   red   :1599
##  1st Qu.:3.110     1st Qu.:0.4300   1st Qu.: 9.50   1st Qu.:5.000   white:4898
##  Median :3.210     Median :0.5100   Median :10.30   Median :6.000
##  Mean   :3.219     Mean   :0.5313   Mean   :10.49   Mean   :5.818
##  3rd Qu.:3.320     3rd Qu.:0.6000   3rd Qu.:11.30   3rd Qu.:6.000
##  Max.   :4.010     Max.   :2.0000   Max.   :14.90   Max.   :9.000

str(wine)

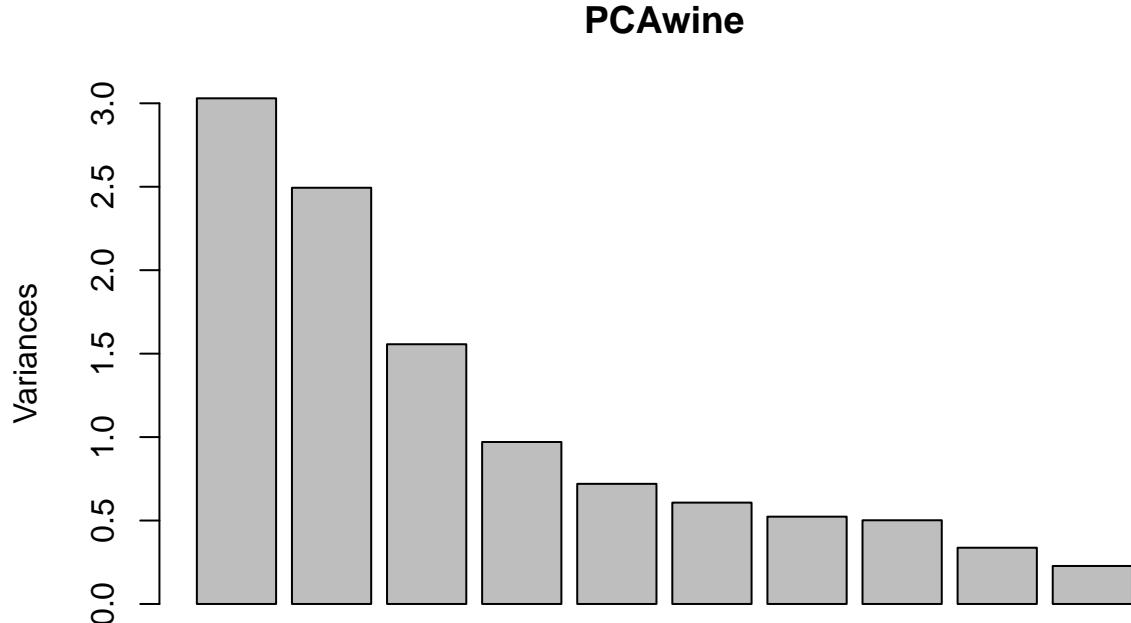
## 'data.frame': 6497 obs. of 13 variables:
## $ fixed.acidity      : num  7.4 7.8 7.8 11.2 7.4 7.4 7.9 7.3 7.8 7.5 ...
## $ volatile.acidity    : num  0.7 0.88 0.76 0.28 0.7 0.66 0.6 0.65 0.58 0.5 ...
## $ citric.acid         : num  0 0 0.04 0.56 0 0 0.06 0 0.02 0.36 ...
## $ residual.sugar      : num  1.9 2.6 2.3 1.9 1.9 1.8 1.6 1.2 2 6.1 ...
## $ chlorides           : num  0.076 0.098 0.092 0.075 0.076 0.075 0.069 0.065 0.073 0.071 ...
## $ free.sulfur.dioxide: num  11 25 15 17 11 13 15 15 9 17 ...
## $ total.sulfur.dioxide: num  34 67 54 60 34 40 59 21 18 102 ...
## $ density              : num  0.998 0.997 0.997 0.998 0.998 ...
```

```

## $ pH : num 3.51 3.2 3.26 3.16 3.51 3.51 3.3 3.39 3.36 3.35 ...
## $ sulphates : num 0.56 0.68 0.65 0.58 0.56 0.56 0.46 0.47 0.57 0.8 ...
## $ alcohol : num 9.4 9.8 9.8 9.8 9.4 9.4 9.4 10 9.5 10.5 ...
## $ quality : int 5 5 5 6 5 5 5 7 7 5 ...
## $ color : Factor w/ 2 levels "red","white": 1 1 1 1 1 1 1 1 1 1 ...
wine_variables = wine %>% select(-quality, -color)

PCAwine = prcomp(wine_variables, scale=TRUE)
plot(PCAwine)

```



```

summary(PCAwine)

## Importance of components:
##          PC1    PC2    PC3    PC4    PC5    PC6    PC7
## Standard deviation   1.7407 1.5792 1.2475 0.98517 0.84845 0.77930 0.72330
## Proportion of Variance 0.2754 0.2267 0.1415 0.08823 0.06544 0.05521 0.04756
## Cumulative Proportion 0.2754 0.5021 0.6436 0.73187 0.79732 0.85253 0.90009
##          PC8    PC9    PC10   PC11
## Standard deviation   0.70817 0.58054 0.4772 0.18119
## Proportion of Variance 0.04559 0.03064 0.0207 0.00298
## Cumulative Proportion 0.94568 0.97632 0.9970 1.00000

round(PCAwine$rotation[,1:3],2)

##          PC1    PC2    PC3
## fixed.acidity -0.24  0.34 -0.43
## volatile.acidity -0.38  0.12  0.31

```

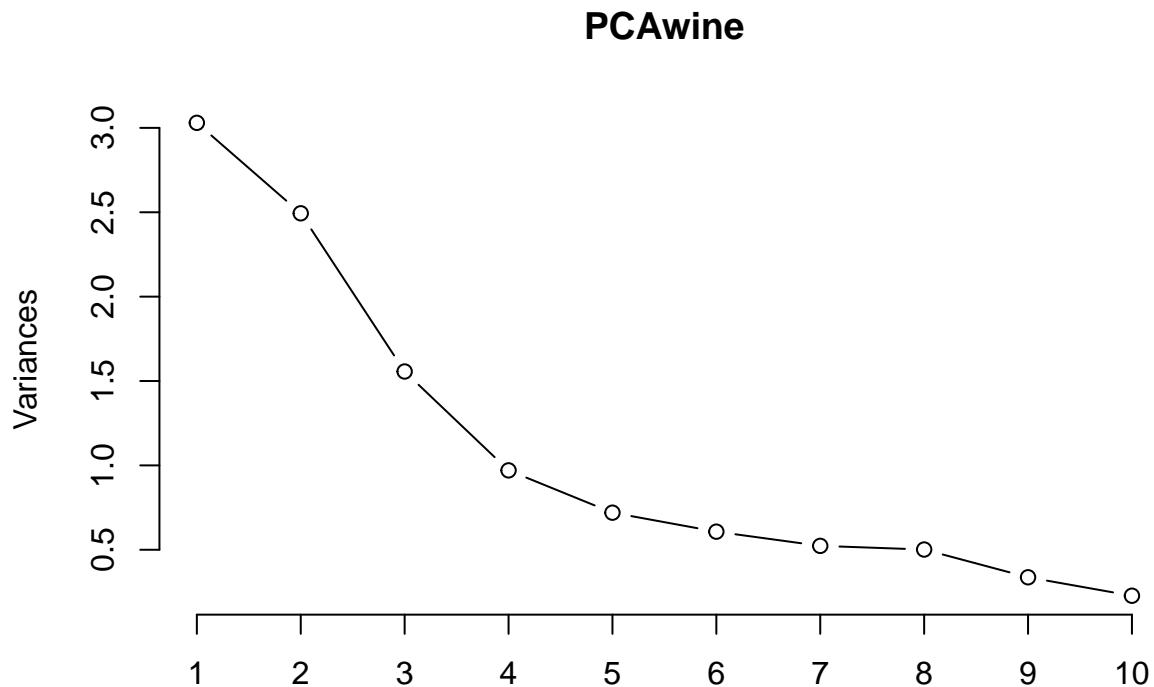
```

## citric.acid      0.15  0.18 -0.59
## residual.sugar  0.35  0.33  0.16
## chlorides       -0.29  0.32  0.02
## free.sulfur.dioxide  0.43  0.07  0.13
## total.sulfur.dioxide  0.49  0.09  0.11
## density        -0.04  0.58  0.18
## pH              -0.22 -0.16  0.46
## sulphates      -0.29  0.19 -0.07
## alcohol         -0.11 -0.47 -0.26

```

After running PCA, we can see that from PC1 and PC2, the cumulative proportion of variance accounts for already roughly 50.21% of the data.

```
# scree plot
screeplot(PCAwine, n pcs = min(10, length(PCAwine$sdev)), type="lines")
```



The scree plot entails that our proportion variance does decrease over time as the number of principal components increases.

```
wine_scores = PCAwine$x
summary(wine_scores)
```

```

##          PC1            PC2            PC3            PC4
##  Min. :-7.3997    Min. :-4.59929   Min. :-6.32835   Min. :-9.42558
##  1st Qu.:-1.0419   1st Qu.:-1.16961  1st Qu.:-0.78278  1st Qu.:-0.56465
##  Median : 0.3069   Median :-0.02271   Median : 0.04275   Median : 0.04002
##  Mean   : 0.0000   Mean   : 0.00000   Mean   : 0.00000   Mean   : 0.00000
##  3rd Qu.: 1.2342   3rd Qu.: 1.09054   3rd Qu.: 0.82321   3rd Qu.: 0.66230

```

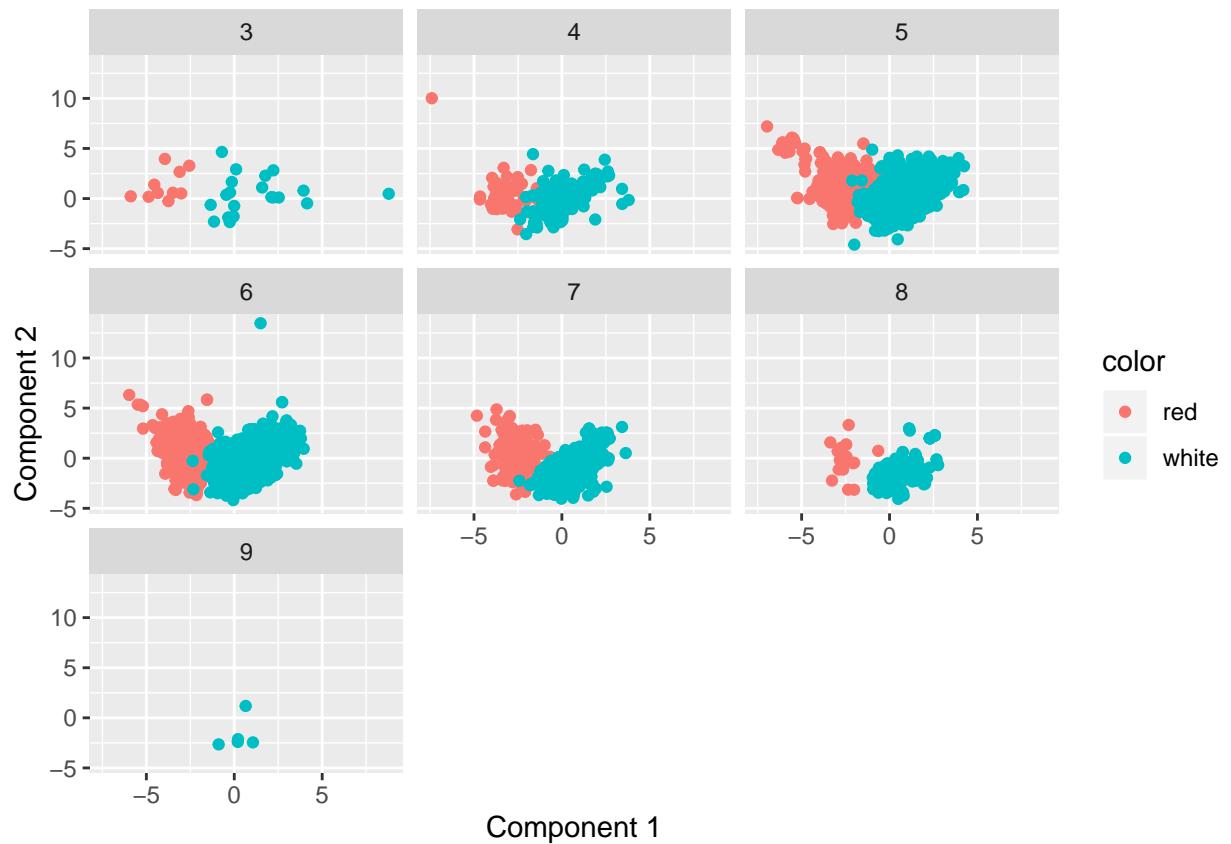
```

##  Max.    : 8.7854   Max.    :13.46906  Max.    : 4.57902  Max.    : 3.12036
##  PC5          PC6          PC7          PC8
##  Min.    :-9.46732  Min.    :-5.25162  Min.    :-4.81874  Min.    :-7.96928
##  1st Qu.:-0.48395  1st Qu.:-0.48713  1st Qu.:-0.45287  1st Qu.:-0.39471
##  Median :-0.04172  Median : 0.01252  Median :-0.01849  Median : 0.05515
##  Mean    : 0.00000  Mean    : 0.00000  Mean    : 0.00000  Mean    : 0.00000
##  3rd Qu.: 0.40509  3rd Qu.: 0.50448  3rd Qu.: 0.44982  3rd Qu.: 0.45527
##  Max.    : 8.79522  Max.    : 2.77108  Max.    : 5.22828  Max.    : 4.07838
##  PC9          PC10         PC11
##  Min.    :-3.30490  Min.    :-3.00229  Min.    :-1.675654
##  1st Qu.:-0.32719  1st Qu.:-0.27608  1st Qu.:-0.109061
##  Median : 0.03997  Median : 0.03022  Median :-0.002271
##  Mean    : 0.00000  Mean    : 0.00000  Mean    : 0.000000
##  3rd Qu.: 0.36710  3rd Qu.: 0.30301  3rd Qu.: 0.103198
##  Max.    : 3.46756  Max.    : 2.84793  Max.    : 4.448971

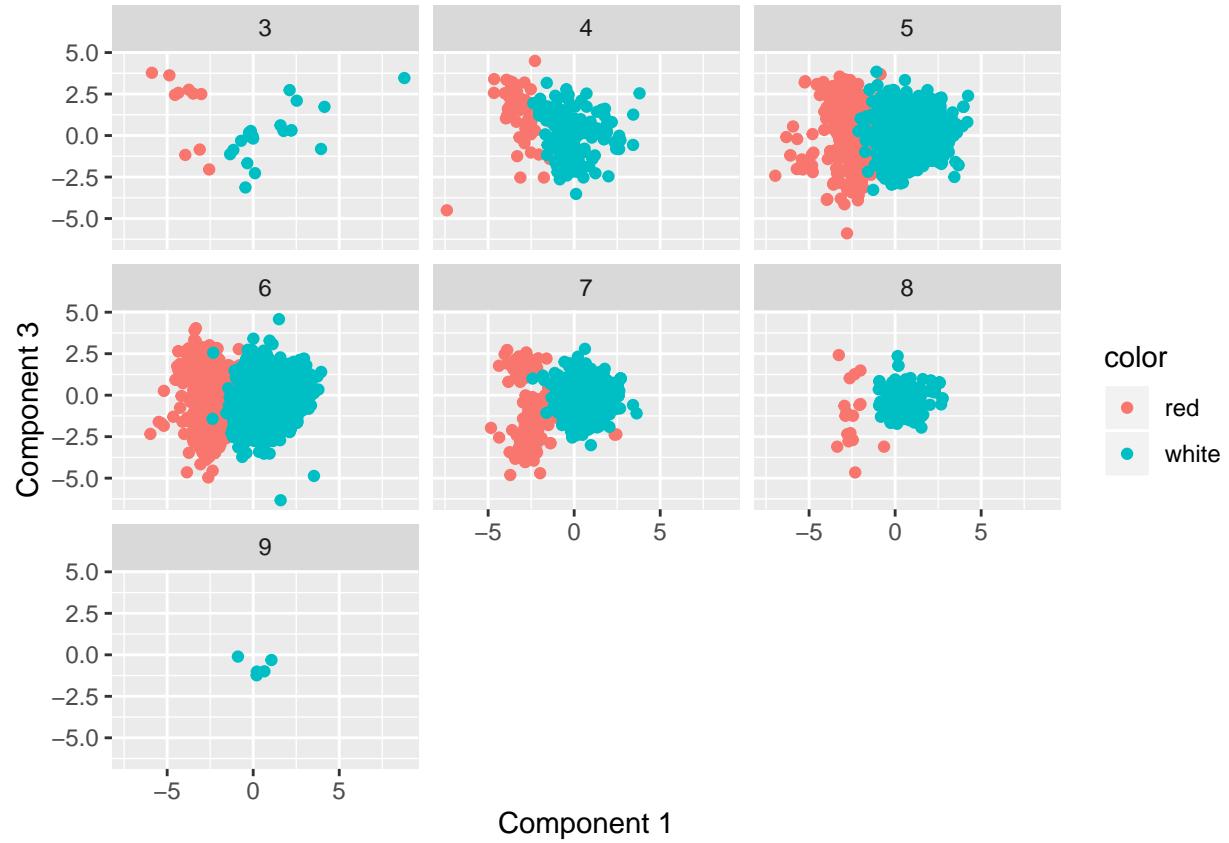
```

# PCA1 gives the best separation

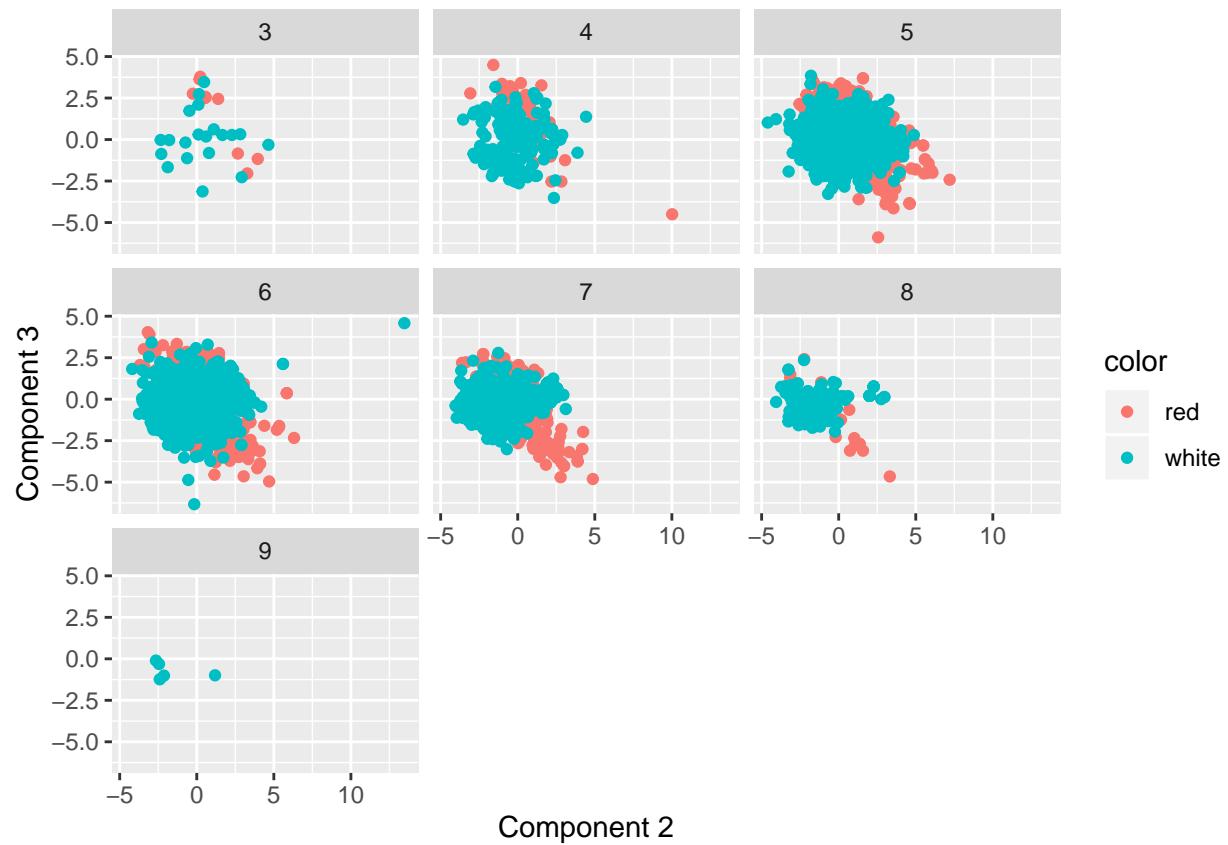
```
qplot(wine_scores[,1], wine_scores[,2], color=color, facets=~quality, xlab='Component 1', ylab='Component 2')
```



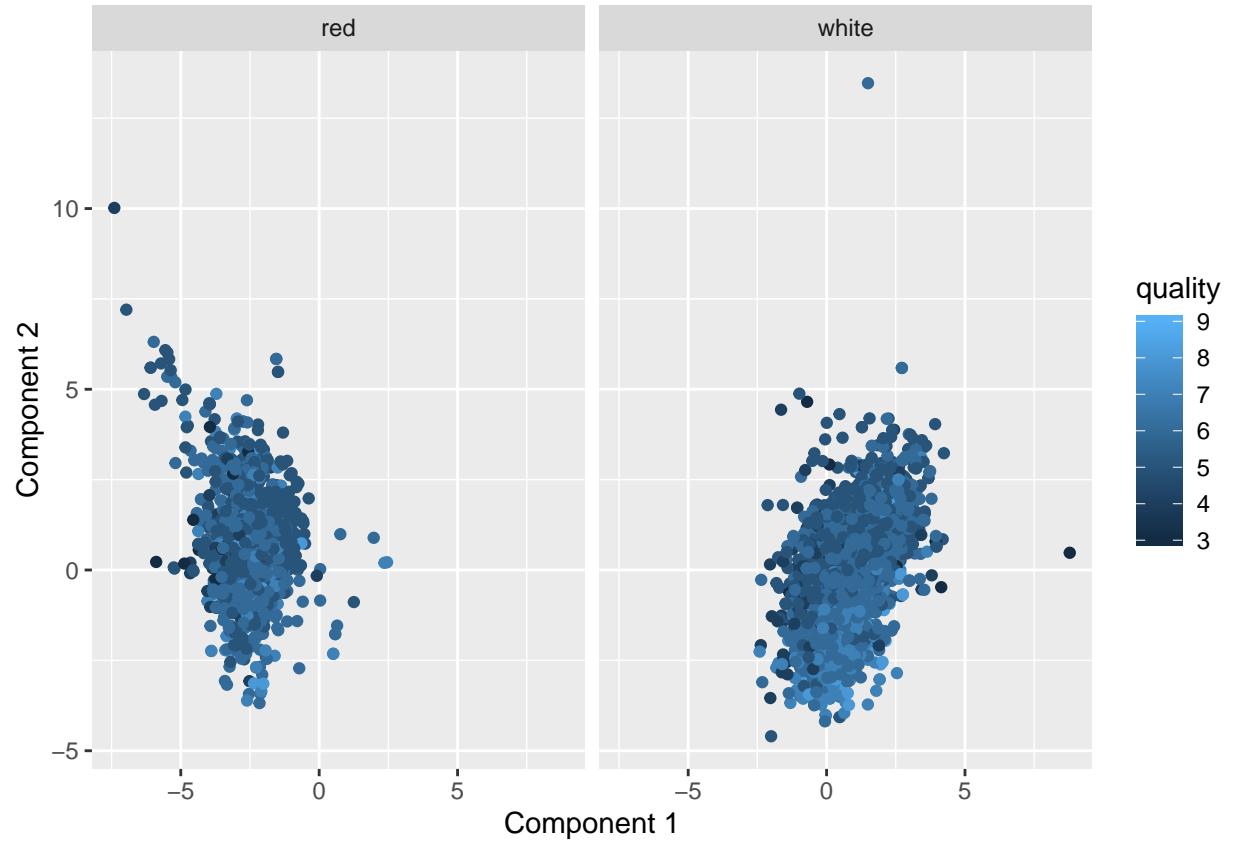
```
qplot(wine_scores[,1], wine_scores[,3], color=color, facets=~quality, xlab='Component 1', ylab='Component 2')
```



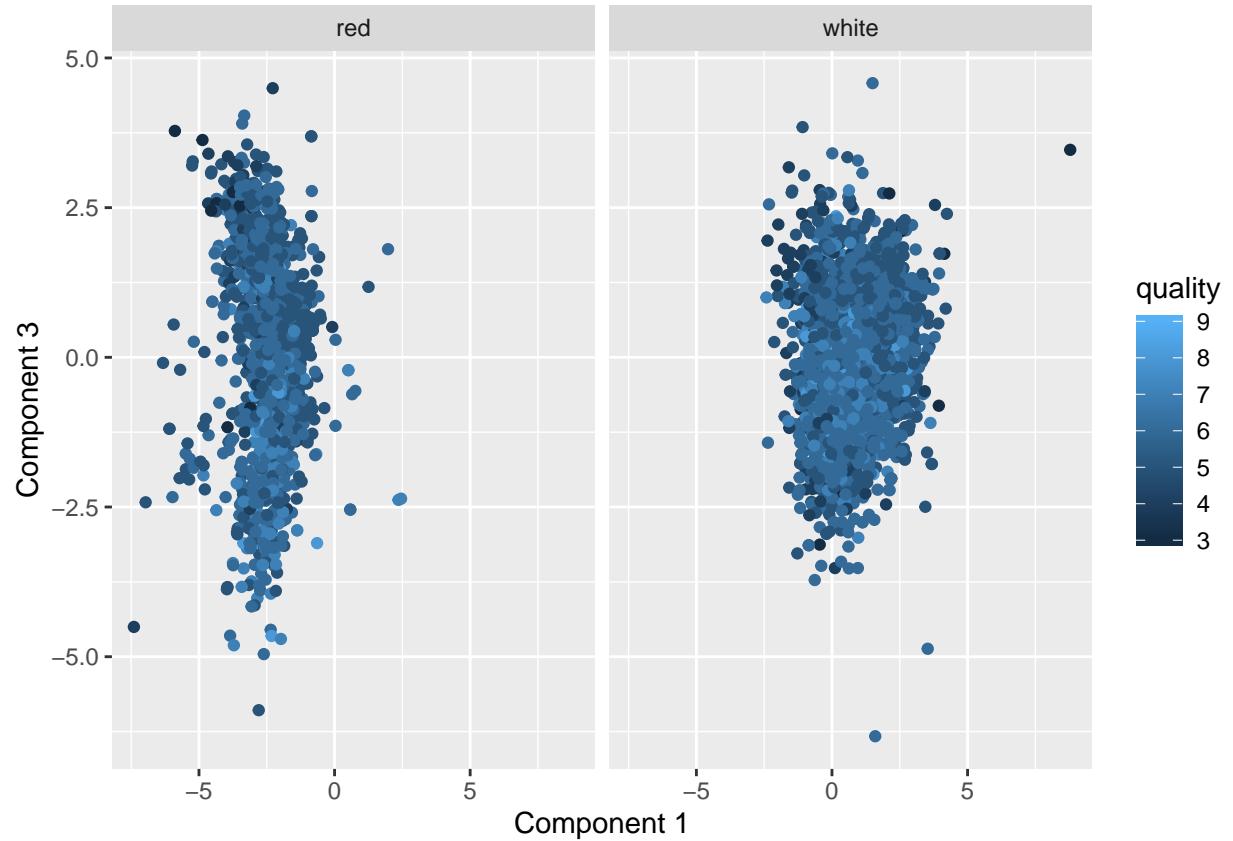
```
# too much overlap, wrong direction
qplot(wine_scores[,2], wine_scores[,3], color=color, facets=~quality, xlab='Component 2', ylab='Component 1')
```



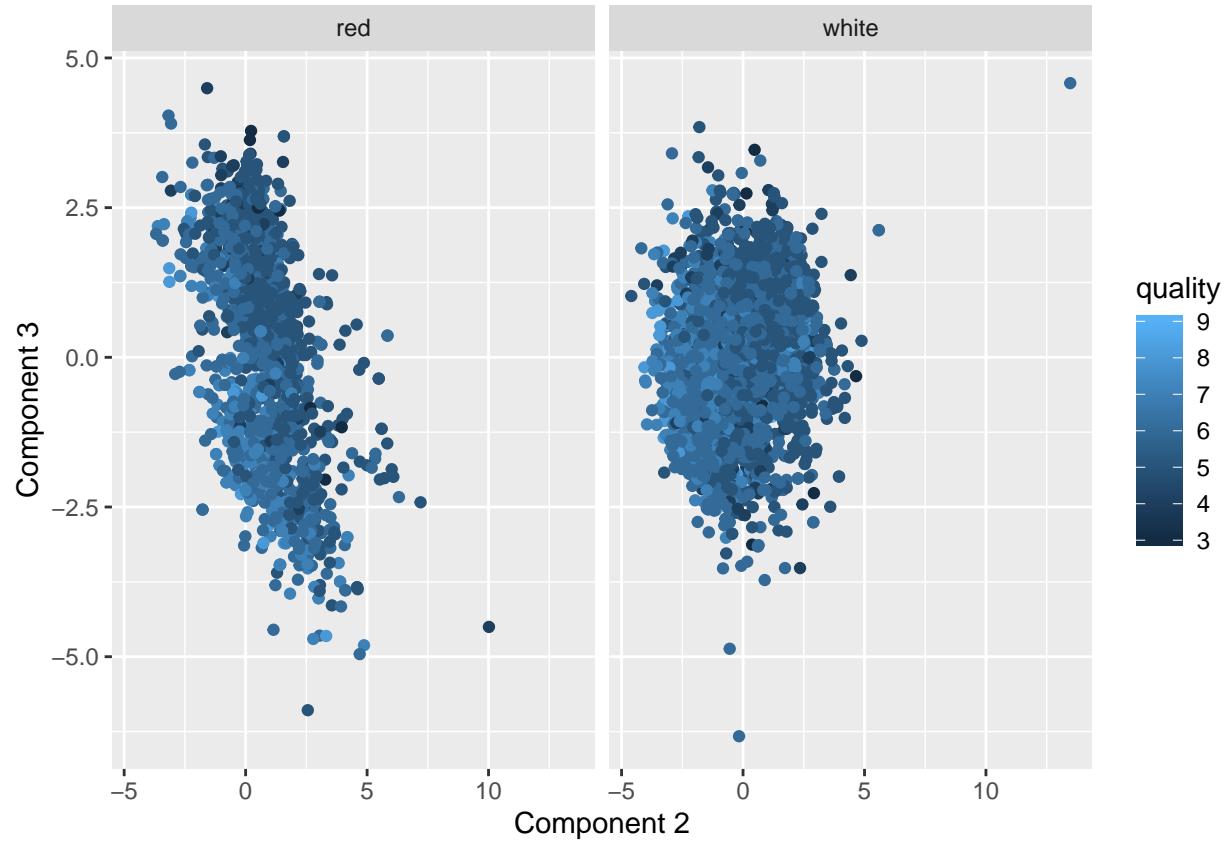
```
# facetting by color and look at quality trends, doesn't work well
qplot(wine_scores[,1], wine_scores[,2], color=quality, facets=~color, xlab='Component 1', ylab='Component 2')
```



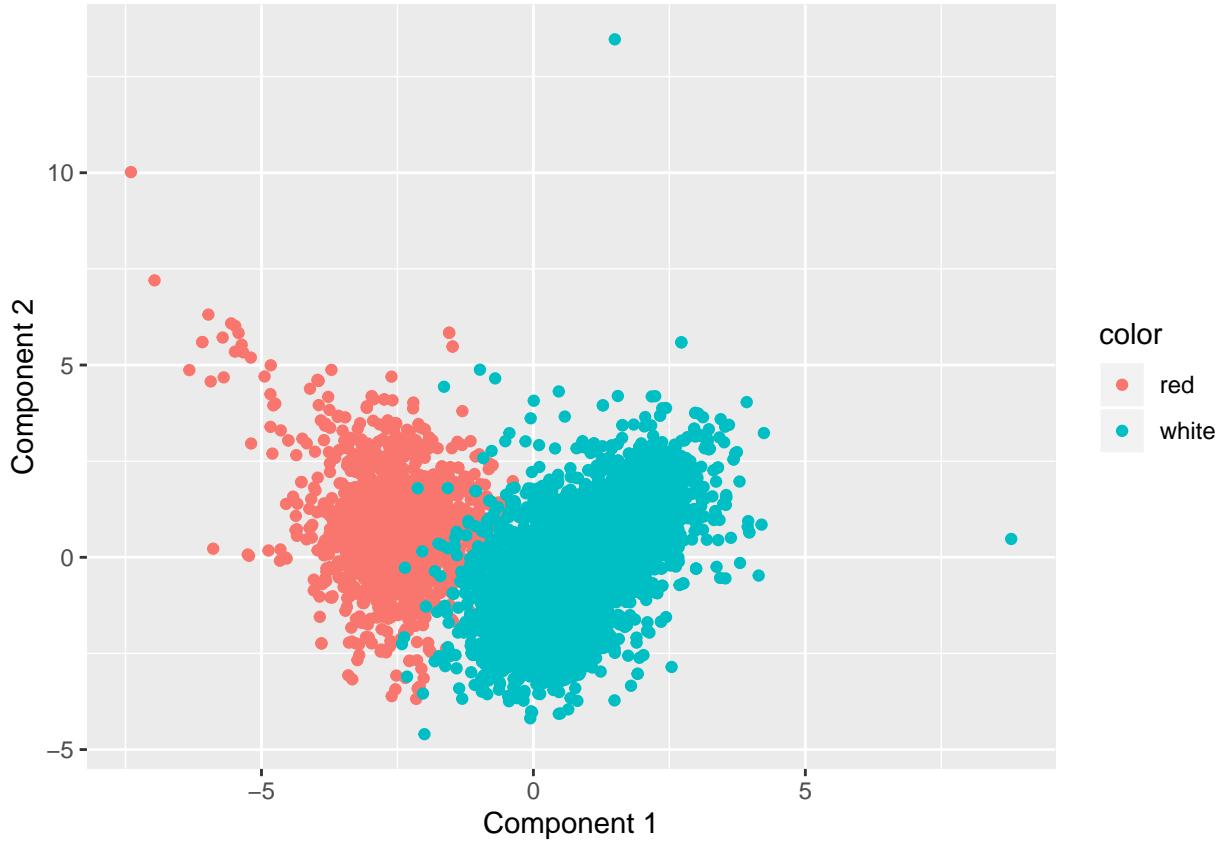
```
qplot(wine_scores[,1], wine_scores[,3], color=quality, facets=~color, xlab='Component 1', ylab='Component
```



```
qplot(wine_scores[,2], wine_scores[,3], color=quality, facets=~color, xlab='Component 2', ylab='Component
```



```
qplot(wine_scores[,1], wine_scores[,2], color=color, xlab='Component 1', ylab='Component 2', data=wine)
```



When comparing PC1, PC2, PC3 in the context of our red and white colors, we see that the noise in the data is not distinguishable as the number of principal components increase; they tend to “weave together”.

Since the PCA with all 3-8 quality points didn't work well (we could not see clear separation) we try to split quality into 2 groups, 5 and below and above 5 and run another PCA.

```

high_qual <- wine %>% filter(grep("6|7|8", quality)) %>% mutate(qual="6-8") %>% select(-quality)
low_qual <- wine %>% filter(grep("3|4|5", quality)) %>% mutate(qual="3-5") %>% select(-quality)
wine_by_qual <- full_join(high_qual, low_qual)

## Joining, by = c("fixed.acidity", "volatile.acidity", "citric.acid", "residual.sugar", "chlorides", "color")
dim(high_qual)

## [1] 4108   13
dim(low_qual)

## [1] 2384   13
wine_var_only = wine_by_qual %>% select(-color, -qual)

PCAwine_qual = prcomp(wine_var_only, scale=TRUE)
summary(PCAwine_qual)

## Importance of components:
##                               PC1      PC2      PC3      PC4      PC5      PC6      PC7
## Standard deviation     1.7408  1.5789  1.2476  0.98524  0.84846  0.77944  0.72299
## Proportion of Variance 0.2755  0.2266  0.1415  0.08824  0.06544  0.05523  0.04752
## Cumulative Proportion   0.2755  0.5021  0.6436  0.73187  0.79731  0.85254  0.90006

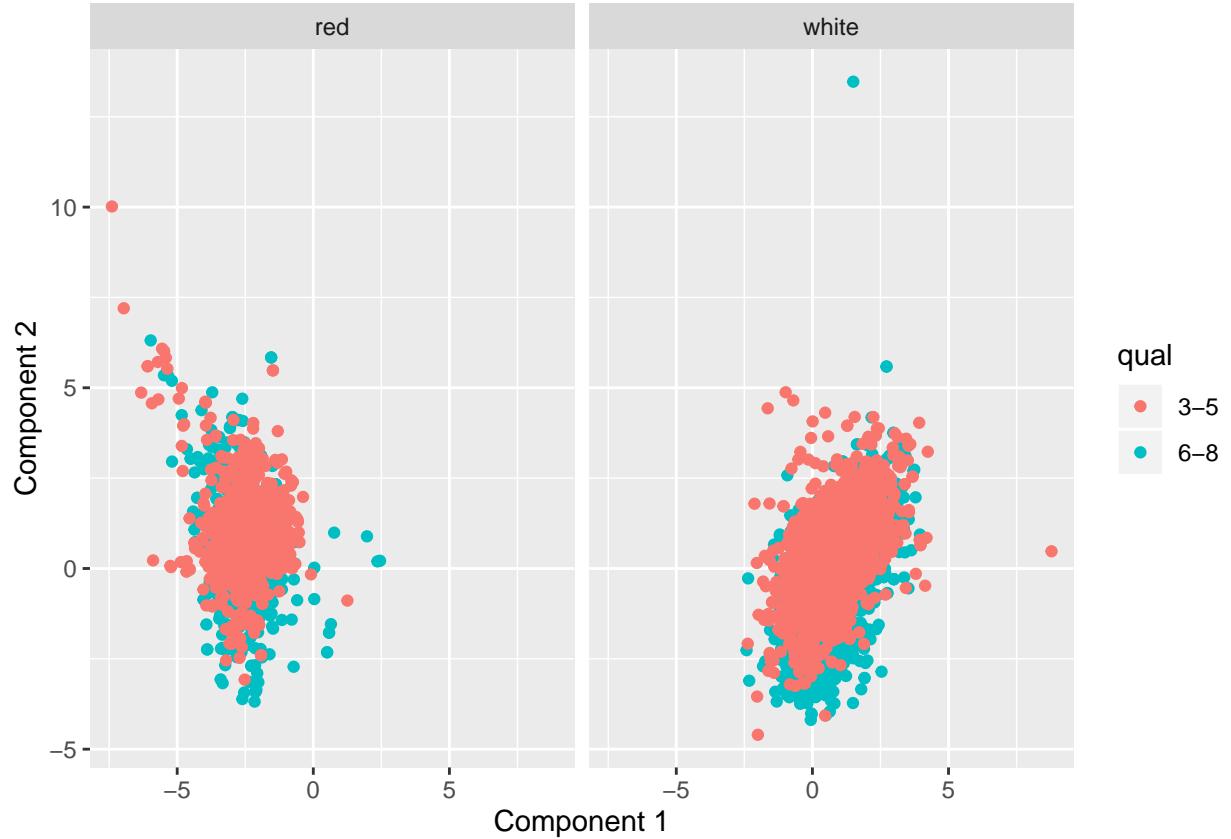
```

```

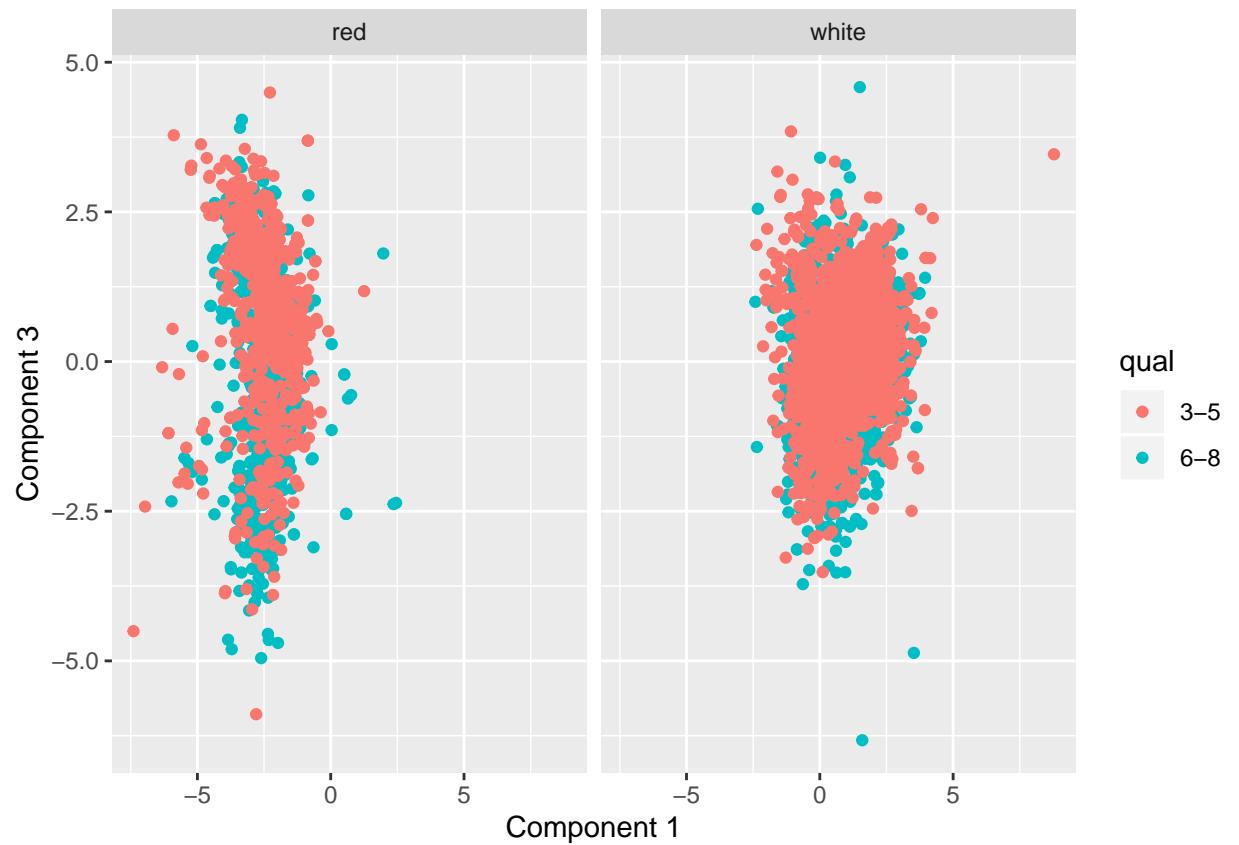
##          PC8      PC9      PC10     PC11
## Standard deviation 0.70837 0.58059 0.47709 0.18113
## Proportion of Variance 0.04562 0.03064 0.02069 0.00298
## Cumulative Proportion 0.94568 0.97632 0.99702 1.00000
qual_scores = PCAwine_qual$x

```

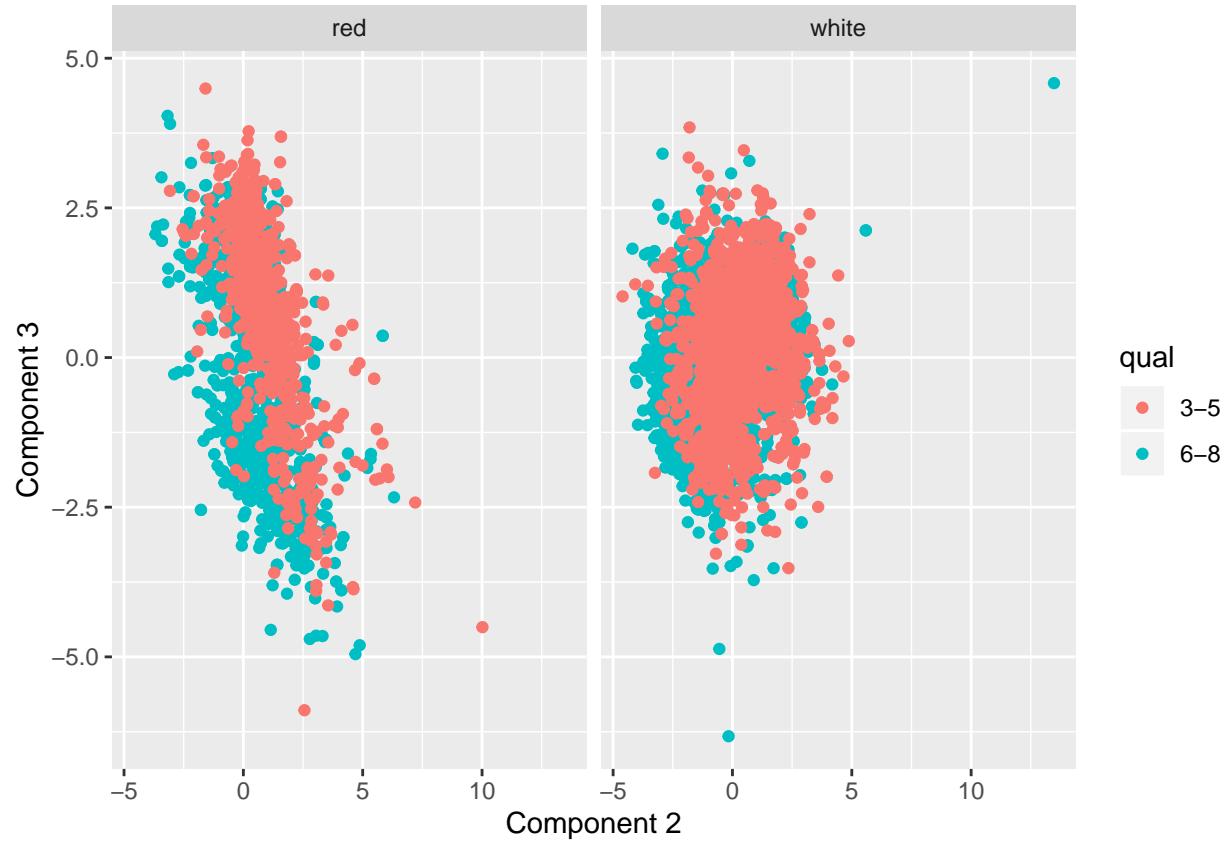
```
qpplot(qual_scores[,1], qual_scores[,2], color=qual, facets=~color,xlab='Component 1', ylab='Component 2'
```



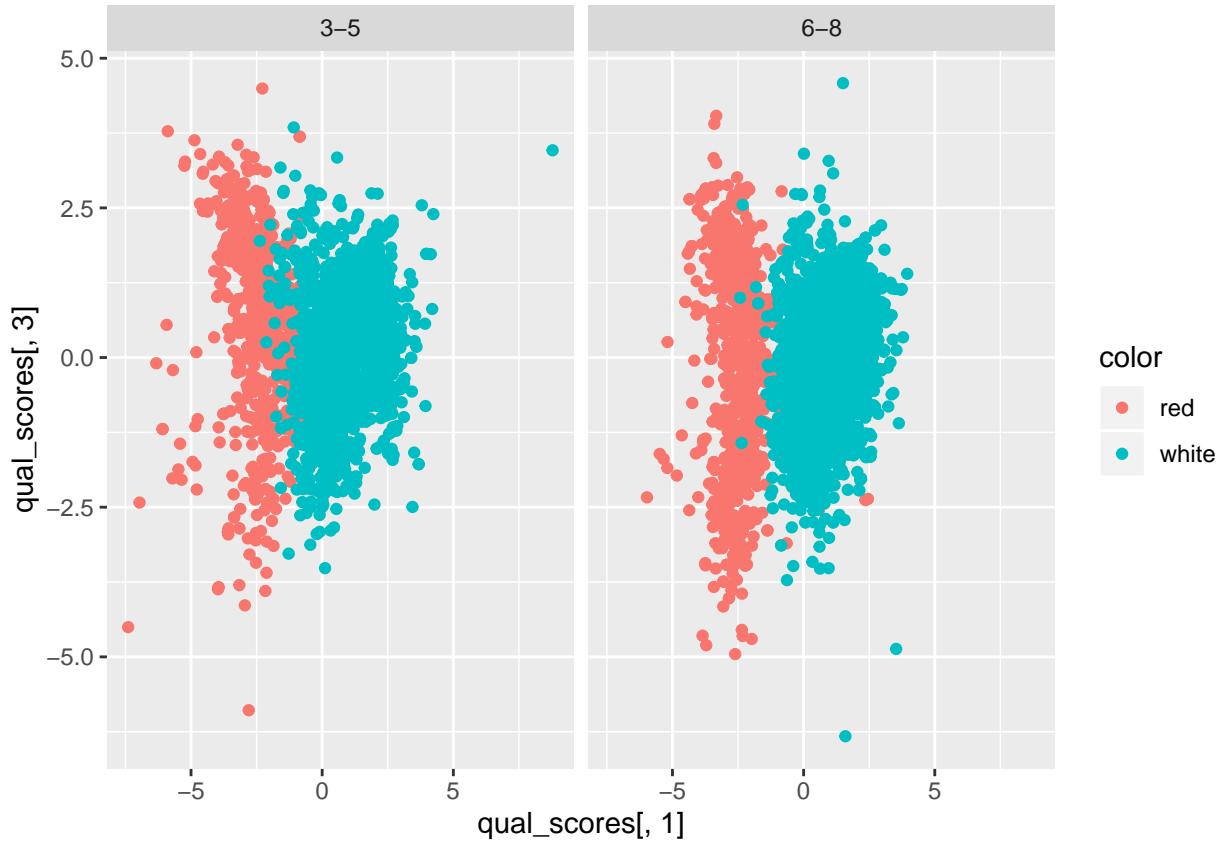
```
qpplot(qual_scores[,1], qual_scores[,3], color=qual, facets=~color,xlab='Component 1', ylab='Component 3'
```



```
qplot(qual_scores[,2], qual_scores[,3], color=qual, facets=~color, xlab='Component 2', ylab='Component 3')
```



```
# despite the color separation still seem very clear
qplot(qual_scores[,1], qual_scores[,3], color=color, facets=~qual, data=wine_by_qual)
```



```
# see that there seems to be separation between higher and lower quality, but not clear with our PCA
```

We can now see a separation between higher and lower quality, but it is still not clear using the PCA method.

```
loadings = PCAwine$rotation
summary(loadings)
```

	PC1	PC2	PC3	PC4
## Min.	-0.38076	-0.4651	-0.5905697	-0.6405
## 1st Qu.	-0.26446	0.0796	-0.1655715	-0.3108
## Median	-0.10644	0.1833	0.1074623	-0.2084
## Mean	-0.01429	0.1451	0.0004754	-0.1472
## 3rd Qu.	0.24915	0.3226	0.1701470	0.1185
## Max.	0.48742	0.5840	0.4553241	0.2128
	PC5	PC6	PC7	PC8
## Min.	-0.45338	-0.51838	-0.41891	-0.52487
## 1st Qu.	-0.24773	-0.31849	-0.34032	-0.34632
## Median	-0.14748	-0.20455	-0.13891	-0.02864
## Mean	-0.05401	-0.13935	-0.12403	-0.08153
## 3rd Qu.	0.15479	0.08986	-0.04641	0.14722
## Max.	0.61439	0.29658	0.52534	0.40124
	PC9	PC10	PC11	
## Min.	-0.49693	-0.713664	-0.449765	
## 1st Qu.	-0.26423	-0.200842	-0.145566	
## Median	0.11288	-0.003908	-0.043438	
## Mean	0.01634	-0.057306	-0.007395	
## 3rd Qu.	0.27417	0.099068	0.031897	

```

##  Max.    : 0.36666   Max.    : 0.480243   Max.    : 0.715162
# 3 most negatively associated variables
loadings[,1] %>% sort %>% head(3)

## volatile.acidity      sulphates      chlorides
## -0.3807575      -0.2941352      -0.2901126

lm1 = lm(quality ~ residual.sugar+free.sulfur.dioxide+total.sulfur.dioxide, data=wine)
summary(lm1)

##
## Call:
## lm(formula = quality ~ residual.sugar + free.sulfur.dioxide +
##     total.sulfur.dioxide, data = wine)
##
## Residuals:
##    Min     1Q   Median     3Q    Max
## -4.3532 -0.7655  0.1214  0.3071  3.2549
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 5.8576627  0.0248018 236.179 <2e-16 ***
## residual.sugar -0.0060152  0.0026065 -2.308  0.021 *
## free.sulfur.dioxide  0.0088917  0.0008758 10.153 <2e-16 ***
## total.sulfur.dioxide -0.0024015  0.0002898 -8.288 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8657 on 6493 degrees of freedom
## Multiple R-squared:  0.01767,   Adjusted R-squared:  0.01721
## F-statistic: 38.93 on 3 and 6493 DF,  p-value: < 2.2e-16

# 3 most positively associated variables
loadings[,1] %>% sort %>% tail(3)

##      residual.sugar  free.sulfur.dioxide total.sulfur.dioxide
##            0.3459199          0.4309140          0.4874181

lm2 = lm(quality ~ volatile.acidity+sulphates+chlorides, data=wine)
summary(lm2)

##
## Call:
## lm(formula = quality ~ volatile.acidity + sulphates + chlorides,
##     data = wine)
##
## Residuals:
##    Min     1Q   Median     3Q    Max
## -3.0898 -0.6770  0.0204  0.4259  3.2028
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 5.98593   0.04033 148.42 <2e-16 ***
## volatile.acidity -1.25159   0.06755 -18.53 <2e-16 ***
## sulphates     0.94610   0.07536 12.55 <2e-16 ***
## chlorides     -4.37342   0.33669 -12.99 <2e-16 ***

```

```

## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.8267 on 6493 degrees of freedom
## Multiple R-squared: 0.1041, Adjusted R-squared: 0.1037
## F-statistic: 251.5 on 3 and 6493 DF, p-value: < 2.2e-16

```

Now, we perform clustering to see if it provides us better information about the quality and color variables of the data. This is done using k-means clustering.

```

library(ggplot2)
library(LICORS) # for kmeans++
library(foreach)
library(mosaic)

wine = read.csv("https://raw.githubusercontent.com/jgscott/SDS323/master/data/wine.csv", header=TRUE)

summary(wine)

## fixed.acidity    volatile.acidity    citric.acid    residual.sugar
## Min. : 3.800    Min. :0.0800    Min. :0.0000    Min. : 0.600
## 1st Qu.: 6.400   1st Qu.:0.2300   1st Qu.:0.2500   1st Qu.: 1.800
## Median : 7.000   Median :0.2900   Median :0.3100   Median : 3.000
## Mean   : 7.215   Mean   :0.3397   Mean   :0.3186   Mean   : 5.443
## 3rd Qu.: 7.700   3rd Qu.:0.4000   3rd Qu.:0.3900   3rd Qu.: 8.100
## Max.   :15.900   Max.   :1.5800   Max.   :1.6600   Max.   :65.800
## chlorides      free.sulfur.dioxide total.sulfur.dioxide density
## Min. :0.00900   Min. : 1.00     Min. : 6.0       Min. :0.9871
## 1st Qu.:0.03800  1st Qu.:17.00    1st Qu.:77.0     1st Qu.:0.9923
## Median :0.04700  Median :29.00    Median :118.0    Median :0.9949
## Mean   :0.05603  Mean   :30.53    Mean   :115.7    Mean   :0.9947
## 3rd Qu.:0.06500  3rd Qu.:41.00    3rd Qu.:156.0    3rd Qu.:0.9970
## Max.   :0.61100  Max.   :289.00   Max.   :440.0    Max.   :1.0390
## pH            sulphates      alcohol      quality      color
## Min. : 2.720   Min. :0.2200   Min. : 8.00   Min. :3.000   red :1599
## 1st Qu.:3.110   1st Qu.:0.4300   1st Qu.: 9.50  1st Qu.:5.000   white:4898
## Median :3.210   Median :0.5100   Median :10.30  Median :6.000
## Mean   :3.219   Mean   :0.5313   Mean   :10.49  Mean   :5.818
## 3rd Qu.:3.320   3rd Qu.:0.6000   3rd Qu.:11.30  3rd Qu.:6.000
## Max.   :4.010   Max.   :2.0000   Max.   :14.90  Max.   :9.000

# Convert variable "color" to a numeric column vector.
# "1" for red and "2" for white
wine$color <- as.numeric(wine$color)

# Center and scale the data
X = wine[,-(12:13)]
X = scale(X, center=TRUE, scale=TRUE)

# Extract the centers and scales from the rescaled data (which are named attributes)
mu = attr(X,"scaled:center")
sigma = attr(X,"scaled:scale")

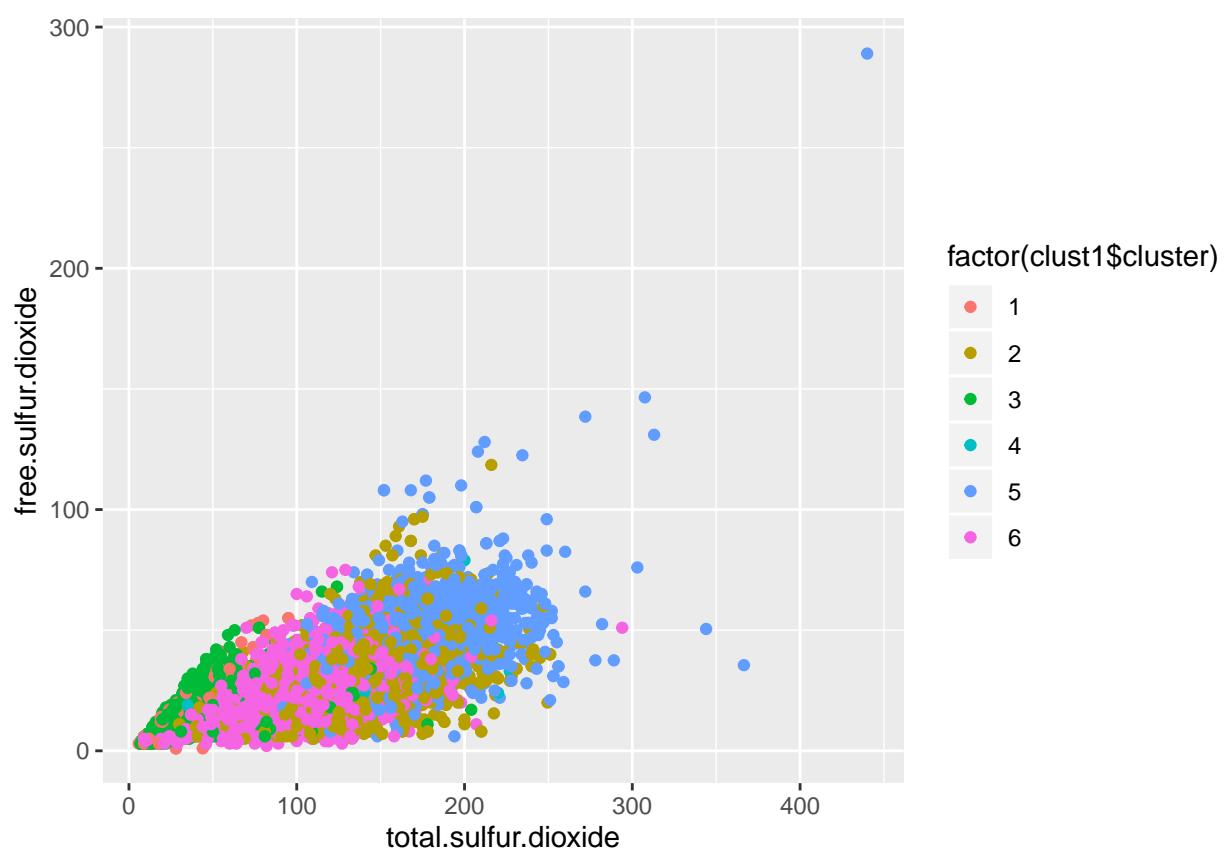
# Run k-means with 6 clusters and 50 starts
clust1 = kmeans(X, 6, nstart=50)

```

```
## Warning: did not converge in 10 iterations
cluster_number = 1
```

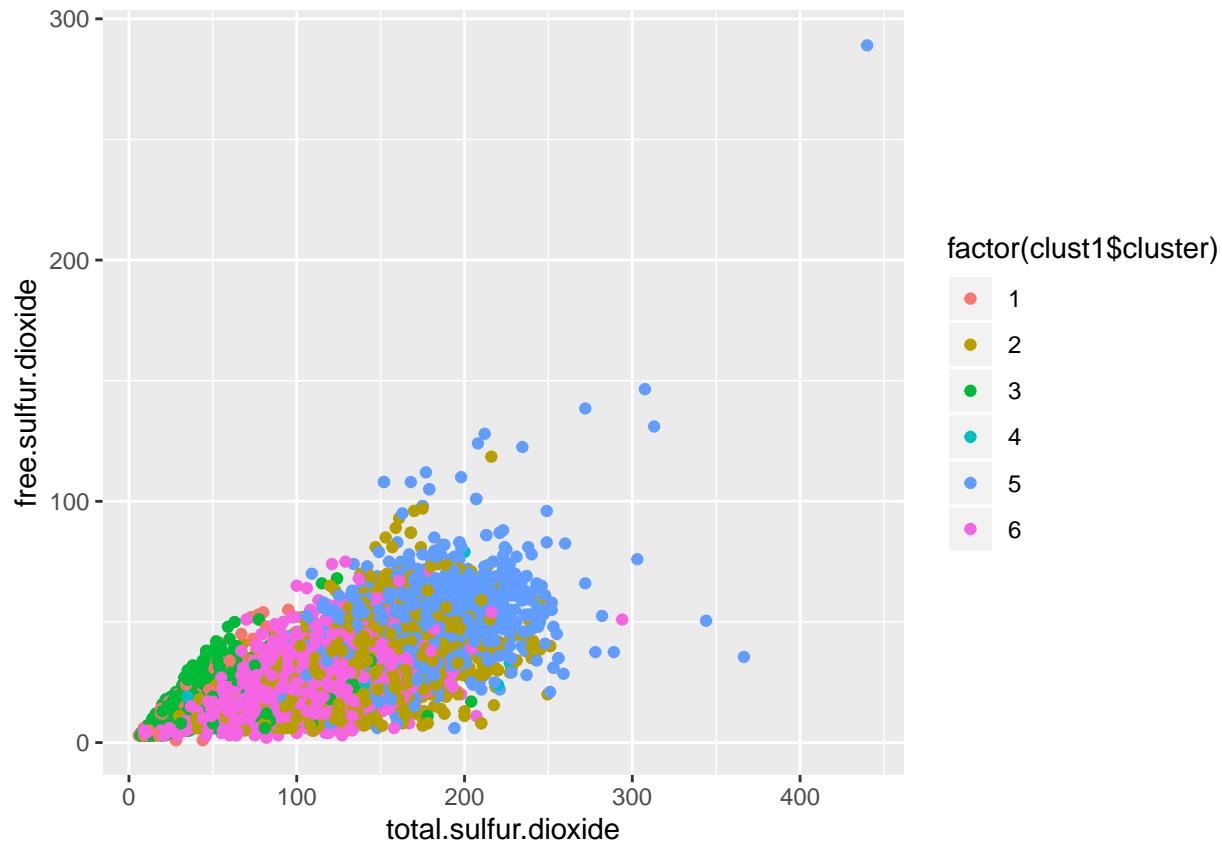
```
##      fixed.acidity    volatile.acidity    citric.acid
##      9.78717949     0.41036199     0.45242836
##      residual.sugar   chlorides    free.sulfur.dioxide
##      2.78061840     0.08449925    14.68024133
##      total.sulfur.dioxide density      pH
##      45.30316742    0.99750226    3.21259427
##      sulphates       alcohol
##      0.71941176    10.61126194
```

```
## [1] "total.sulfur.dioxide"
## [1] "free.sulfur.dioxide"
```



```
##      fixed.acidity    volatile.acidity    citric.acid
##      7.28680042     0.61903907     0.13004224
##      residual.sugar   chlorides    free.sulfur.dioxide
##      2.46737064     0.07932524    16.39651531
##      total.sulfur.dioxide density      pH
##      50.40760296    0.99610666    3.37548046
##      sulphates       alcohol
##      0.59145723    10.23440690
```

```
## [1] "total.sulfur.dioxide"
## [1] "free.sulfur.dioxide"
```



```
# What are the clusters?
clust1$center # not super helpful
```

```
##   fixed.acidity volatile.acidity citric.acid residual.sugar   chlorides
## 1    1.98380550      0.4294066  0.92070676   -0.5596315  0.8125166
## 2   -0.37768962     -0.4917289 -0.02367606   -0.3364811 -0.1639556
## 3    0.05514617      1.6969087 -1.29778246   -0.6254702  0.6648296
## 4    0.74551523      1.0520924  1.21887533   -0.4670271  8.7931881
## 5   -0.16076017     -0.3549877  0.30908083    1.4196876 -0.1608094
## 6   -0.31637697     -0.3427032  0.08034457   -0.4211476 -0.5680458
##   free.sulfur.dioxide total.sulfur.dioxide      density       pH sulphates
## 1      -0.89271064      -1.2462685  0.9356234 -0.03673537  1.2643552
## 2       0.08759749      0.3616378 -0.3048354  0.25558927 -0.1483107
## 3      -0.79601588     -1.1559594  0.4702178  0.97631911  0.4044797
## 4      -0.68689568     -0.6924492  0.7568104 -0.84066089  3.4798534
## 5       0.91917593      0.9861635  0.8803160 -0.48933952 -0.2744992
## 6      -0.13027169     -0.1166923 -1.1934858 -0.31693115 -0.4016445
##   alcohol
## 1  0.1001592
## 2 -0.2904650
## 3 -0.2158056
## 4 -0.7832782
## 5 -0.8325142
## 6  1.1763139
```

```

clust1$center[1,]*sigma + mu

##      fixed.acidity    volatile.acidity      citric.acid
##      9.78717949      0.41036199      0.45242836
##      residual.sugar      chlorides free.sulfur.dioxide
##      2.78061840      0.08449925     14.68024133
## total.sulfur.dioxide      density          pH
##      45.30316742      0.99750226     3.21259427
##      sulphates      alcohol
##      0.71941176      10.61126194

clust1$center[2,]*sigma + mu

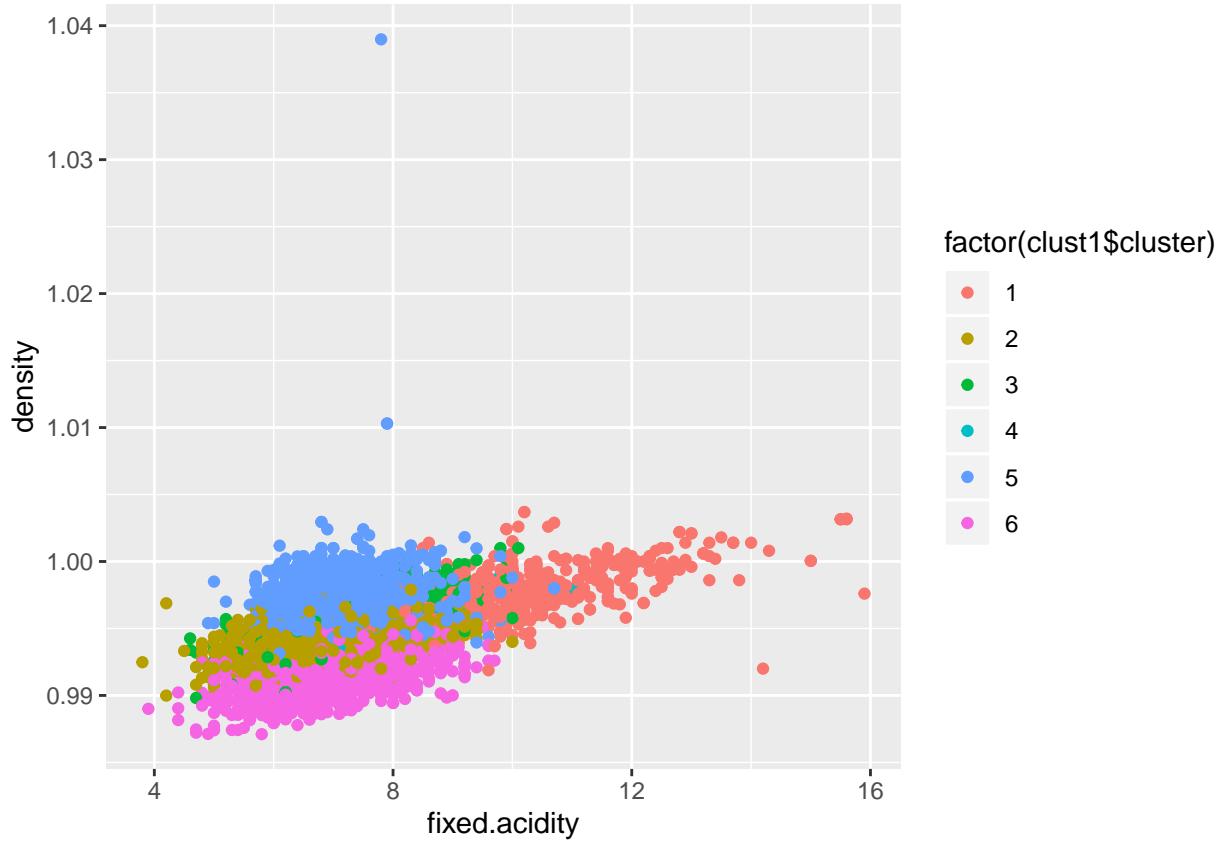
##      fixed.acidity    volatile.acidity      citric.acid
##      6.72565749      0.25870948      0.31519266
##      residual.sugar      chlorides free.sulfur.dioxide
##      3.84232416      0.05028991     32.08012232
## total.sulfur.dioxide      density          pH
##      136.18501529      0.99378253     3.25959633
##      sulphates      alcohol
##      0.50919878      10.14535984

clust1$center[4,]*sigma + mu

##      fixed.acidity    volatile.acidity      citric.acid
##      8.1818182       0.5128788      0.4957576
##      residual.sugar      chlorides free.sulfur.dioxide
##      3.2212121       0.3640909     18.3333333
## total.sulfur.dioxide      density          pH
##      76.6060606      0.9969661     3.0833333
##      sulphates      alcohol
##      1.0490909       9.5575758

# A few plots with cluster membership shown
# qplot is in the ggplot2 library
#qplot(total.sulfur.dioxide, citric.acid, data=wine, #color=factor(clust1$cluster))
qplot(fixed.acidity, density, data=wine, color=factor(clust1$cluster))

```



After performing k-means clustering, we see that the data can be distinguished better in terms of clusters.

### Conclusion:

From both PCA and k-means clustering, we can see that k-means clustering offers a better visualization of distinguishing the quality and reds/whites in the data since k-means allows us to look at it in clusters.

PCA, on the other hand, does not: it is unable to distinguish the red and white wines and the higher and lower quality wines.

However in term of what “makes sense”, PCA has better interpretability compared to k-means clustering, because it summarize the information in components while k-means summarizes the data in chunks.

## Problem 4: Market Segmentation

### Problem Overview:

NutrientH2O has collected twitter data on 7882 of their followers. For each of their followers, the data collectors have counted the frequencies the follower will tweet about 36 given topics, such as sports or politics, over a week. Proper analysis of this data could lead to better targeting and marketing for products. This report will aim to identify different groups or types of customers, or segments of the market, that are part of the customer base for Nutrient H2O.

### Data and Analysis Process

Following the principle of satisfice and some trial and error, 5 clusters will be identified using the KMEANS clustering algorithm. Much more than 5 will probably be unfeasible for NutrientH2O to carry out targeted ad campaigns for all clustered groups.

Each cluster will represent a segment of the market and a “kind” of customer to target for NutrientH20. What each cluster represents will be determined by the Euclidean coordinates of the centroid (the average point in the cluster). If the coordinates for any one topic are particularly high, it means that the frequency that the average person in that cluster tweets about that topic is also high. Among the topics, the topic of “chatter” has been discarded from the dataset as it is very general and not that useful for marketing purposes.

## Results

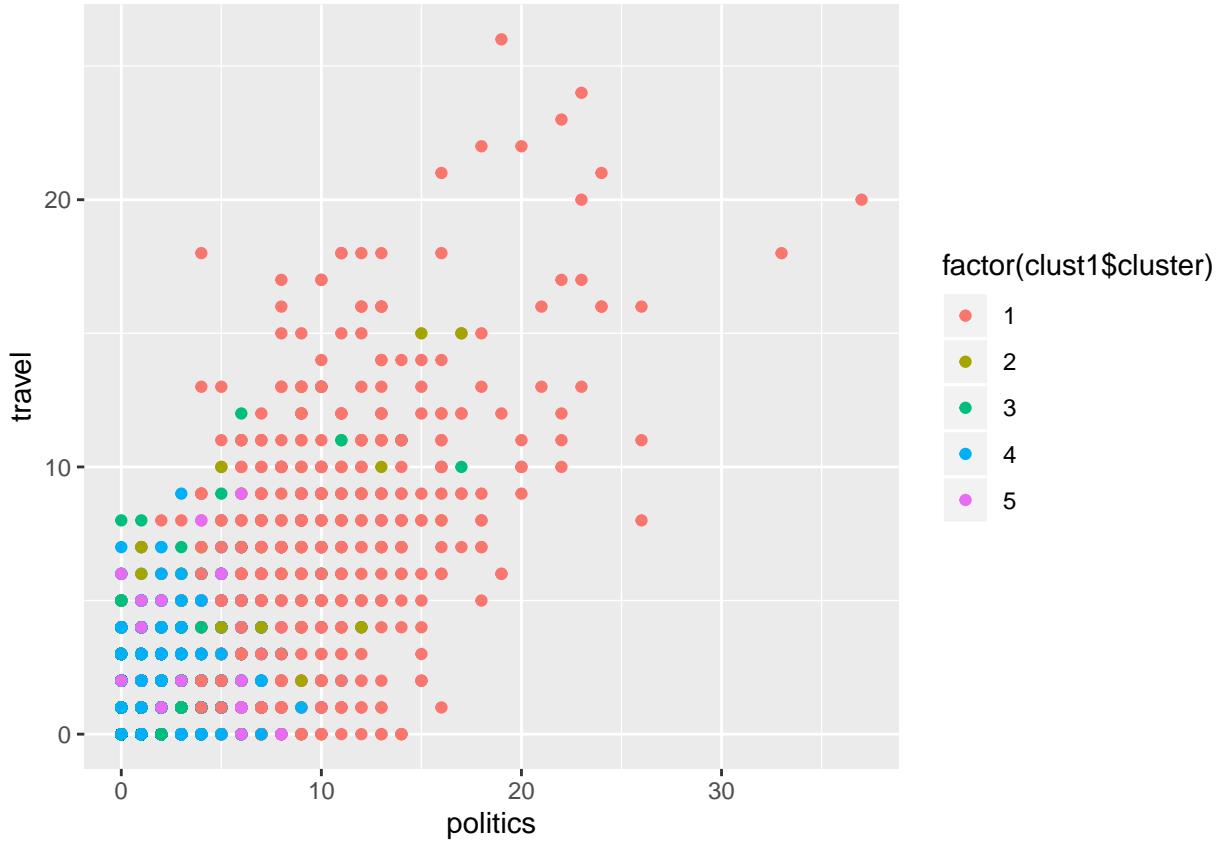
The centroid coordinates for the first cluster are:

```
##   current_events          travel photo_sharing uncategorized
##   1.646153846      5.555244755    2.516083916     0.787412587
##   tv_film      sports_fandom    politics           food
##   1.219580420      1.997202797    8.826573427    1.441958042
##   family   home_and_garden      music            news
##   0.918881119      0.612587413    0.634965035    5.208391608
##   online_gaming        shopping health_nutrition college_uni
##   1.186013986      1.374825175    1.658741259    1.682517483
##   sports_playing       cooking         eco computers
##   0.703496503      1.267132867    0.597202797    2.464335664
##   business        outdoors        crafts automotive
##   0.671328671      0.921678322    0.648951049    2.318881119
##   art             religion        beauty parenting
##   0.752447552      1.009790210    0.467132867    0.924475524
##   dating          school personal_fitness fashion
##   1.060139860      0.714685315    0.998601399    0.682517483
##   small_business        spam        adult
##   0.476923077      0.008391608    0.269930070
```

The topics that stand out most in this cluster are:

```
## [1] "politics"
## [1] "travel"
```

The following is a plot of the users based on the frequency with which they tweet about politics and travel. The users are color coded by cluster.



As can be seen from the plot above, the most frequent tweeters on topics politics and travel (top right of the graph) are mostly all in cluster 1.

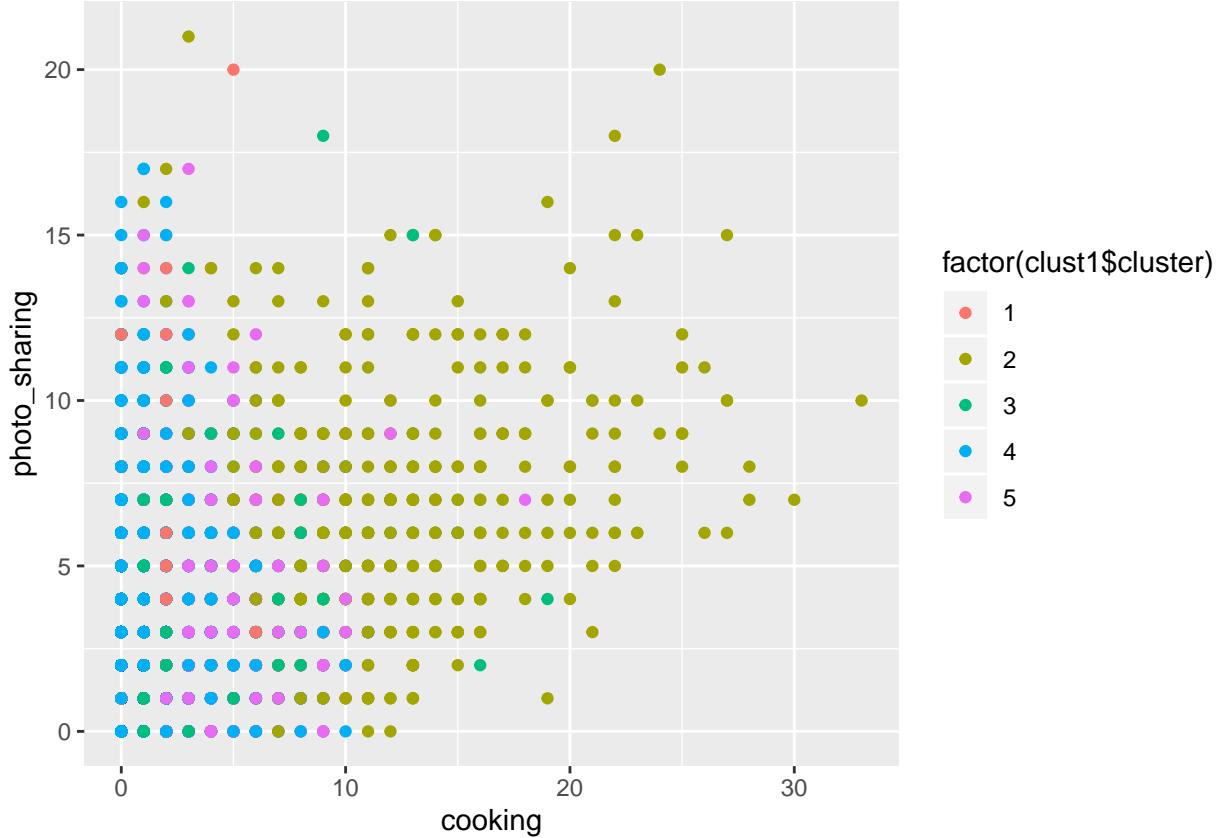
The centroid coordinates for the second cluster are:

	current_events	travel	photo_sharing	uncategorized
##	1.766556291	1.514900662	6.082781457	1.293046358
##	tv_film	sports_fandom	politics	food
##	1.180463576	1.165562914	1.403973510	1.102649007
##	family	home_and_garden	music	news
##	0.912251656	0.642384106	1.284768212	1.044701987
##	online_gaming	shopping	health_nutrition	college_uni
##	1.478476821	2.104304636	2.271523179	2.038079470
##	sports_playing	cooking	eco	computers
##	0.932119205	10.607615894	0.579470199	0.746688742
##	business	outdoors	crafts	automotive
##	0.617549669	0.841059603	0.647350993	0.902317881
##	art	religion	beauty	parenting
##	0.995033113	0.870860927	3.804635762	0.807947020
##	dating	school	personal_fitness	fashion
##	0.996688742	0.985099338	1.352649007	5.480132450
##	small_business	spam	adult	
##	0.529801325	0.003311258	0.415562914	

The topics that stand out most in this cluster are:

```
## [1] "cooking"
## [1] "photo_sharing"
```

The following is a plot of the users based on the frequency with which they tweet about cooking and photo\_sharing. The users are color coded by cluster.



As can be seen from the plot above, the most frequent tweeters on topics cooking and photo\_sharing (top right of the graph) are mostly all in cluster 2.

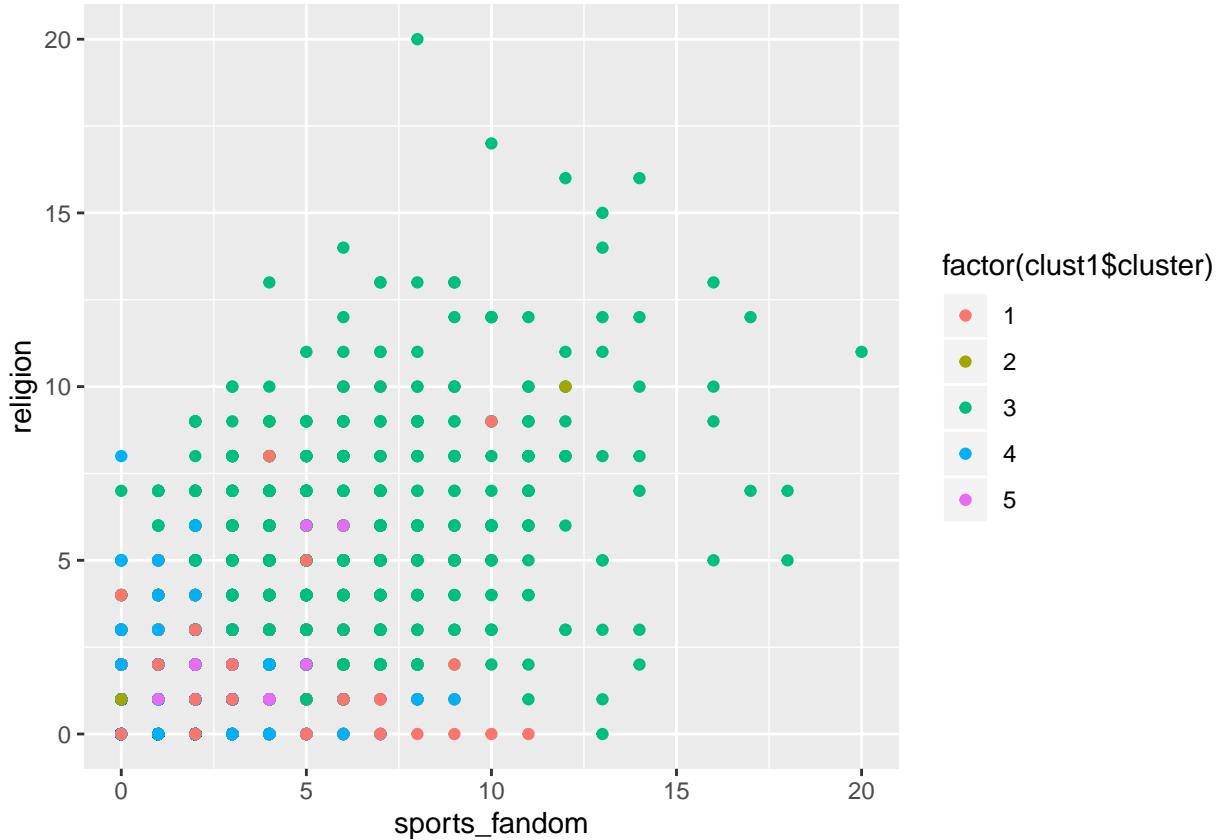
The centroid coordinates for the third cluster are:

	current_events	travel	photo_sharing	uncategorized
##	1.67598475	1.36467598	2.62134689	0.76493011
##	tv_film	sports_fandom	politics	food
##	1.10546379	5.86912325	1.17662008	4.53748412
##	family	home_and_garden	music	news
##	2.48030496	0.66200762	0.75984752	1.04828463
##	online_gaming	shopping	health_nutrition	college_uni
##	1.28589581	1.46759848	1.86277001	1.54129606
##	sports_playing	cooking	eco	computers
##	0.80432020	1.59085133	0.65311309	0.74459975
##	business	outdoors	crafts	automotive
##	0.49936468	0.70775095	1.07369759	1.04955527
##	art	religion	beauty	parenting
##	0.88818297	5.24396442	1.08005083	4.01905972
##	dating	school	personal_fitness	fashion
##	0.77255400	2.68869123	1.20203304	1.00635324
##	small_business	spam	adult	
##	0.41041931	0.00635324	0.40406607	

The topics that stand out most in this cluster are:

```
## [1] "sports_fandom"
## [1] "religion"
```

The following is a plot of the users based on the frequency with which they tweet about sports\_fandom and religion. The users are color coded by cluster.



As can be seen from the plot above, the most frequent tweeters on topics sports\_fandom and religion (top right of the graph) are mostly all in cluster 3.

The centroid coordinates for the fourth cluster are:

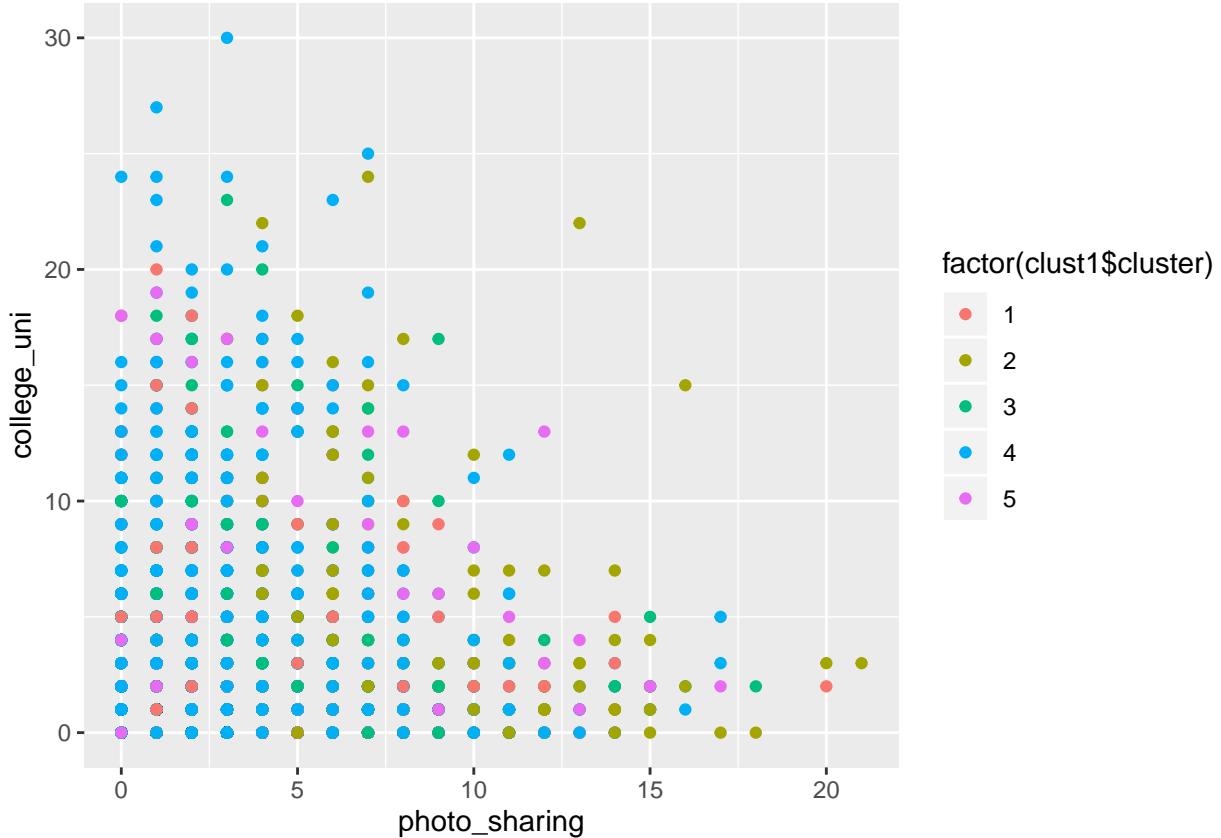
```
##   current_events      travel photo_sharing uncategorized
## 1 1.449059723 1.110766687 2.311634635 0.734242612
##   tv_film    sports_fandom     politics        food
## 1 1.035751188 0.976648068 1.001239926 0.779706551
##   family   home_and_garden       music        news
## 1 0.600537301 0.446786526 0.580078529 0.682579045
##   online_gaming      shopping health_nutrition college_uni
## 1 1.168216574 1.270510436 1.056829924 1.511675966
##   sports_playing      cooking         eco computers
## 1 0.556106634 0.854928704 0.391196528 0.372184336
##   business        outdoors       crafts automotive
## 1 0.339532961 0.401115933 0.373217607 0.594957636
##   art            religion       beauty parenting
## 1 0.655713990 0.527381690 0.348419095 0.463525522
##   dating          school personal_fitness fashion
## 1 0.552800165 0.470551767 0.639181649 0.524281876
##   small_business        spam        adult
```

```
##      0.287662740      0.006612937      0.417854929
```

The topics that stand out most in this cluster are:

```
## [1] "photo_sharing"
## [1] "college_uni"
```

The following is a plot of the users based on the frequency with which they tweet about photo\_sharing and college\_uni. The users are color coded by cluster.



As can be seen from the plot above, the most frequent tweeters on topics photo\_sharing and college\_uni (top right of the graph) are mostly all in cluster 4.

The centroid coordinates for the fifth cluster are:

```
##   current_events          travel    photo_sharing     uncategorized
##   1.552828175  1.234791889  2.704375667  0.970117396
##   tv_film        sports_fandom   politics       food
##   1.034151547  1.160085379  1.246531483  2.106723586
##   family         home_and_garden music        news
##   0.792956243  0.635005336  0.767342583  1.087513340
##   online_gaming   shopping     health_nutrition college_uni
##   1.197438634  1.487726788  11.843116329 1.335112060
##   sports_playing  cooking      eco        computers
##   0.691568837  3.252934899  0.911419424  0.550693703
##   business       outdoors     crafts      automotive
##   0.477054429  2.672358591  0.597652081  0.675560299
##   art            religion     beauty      parenting
##   0.749199573  0.754535752  0.416221985  0.754535752
```

```

##          dating           school      personal_fitness      fashion
## 1.024546425 0.589114194 6.354322305 0.776947705
## small_business          spam        adult
## 0.293489861 0.006403415 0.421558164

```

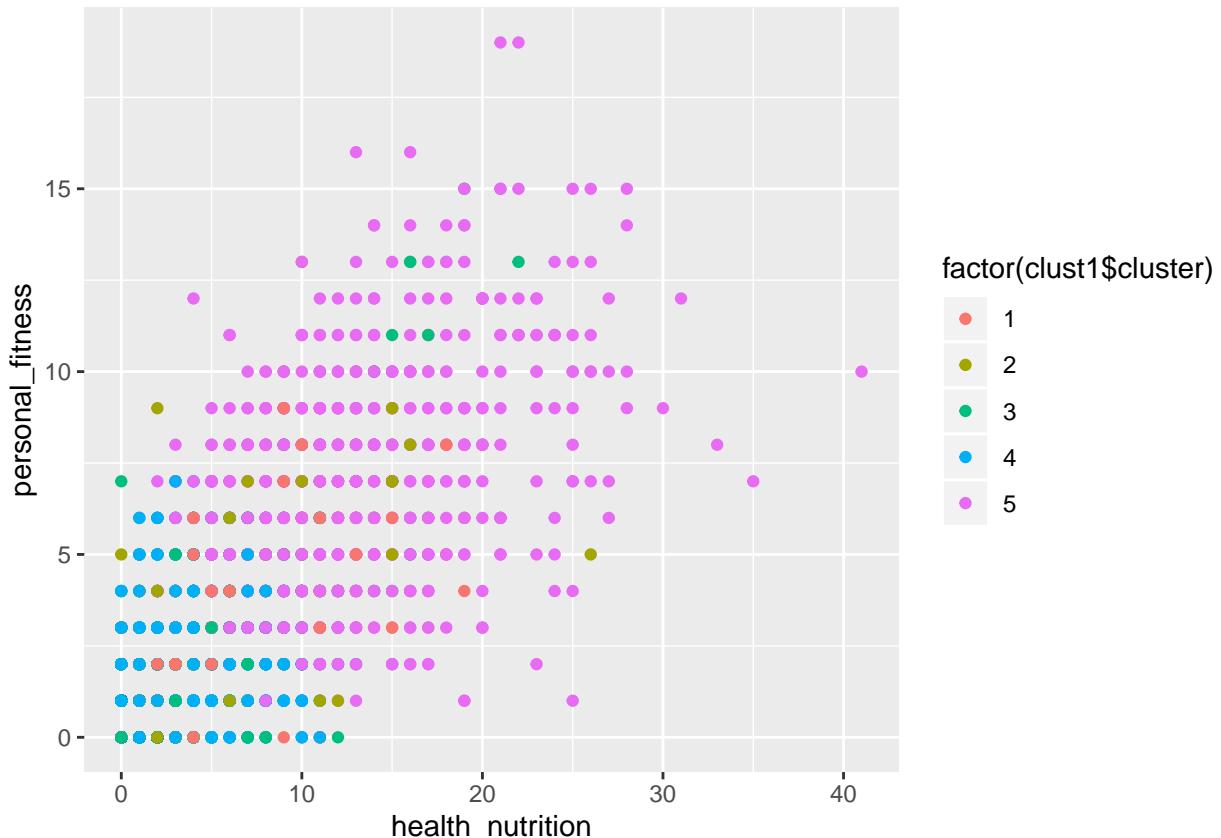
The topics that stand out most in this cluster are:

```

## [1] "health_nutrition"
## [1] "personal_fitness"

```

The following is a plot of the users based on the frequency with which they tweet about health\_nutrition and personal\_fitness. The users are color coded by cluster.



As can be seen from the plot above, the most frequent tweeters on topics health\_nutrition and personal\_fitness (top right of the graph) are mostly all in cluster 5.

## Conclusions

The four market segments that best seem to be indicated from the Twitter data are:

- Users that cook and share photos often (cooking, photo\_sharing)
- Consumers that are into health and nutrition and personal fitness (health\_nutrition, personal\_fitness)
- Politically informed/opinionated travellers (politics, travel)
- And religious sports fanatics (religion, sports\_fandom)

A fifth cluster for photo sharing college/university students was tested but the data didn't cluster well.

NutrientH2O can make good use of these clusterings by marketing differently to customers in those different groups.