# STA 380 Homework 1

Henry Chang, Tiffany Sung, Joseph Chin

8/8/2018

## Problem 1 Probability Practice

### Part A

Given the following: Prob(RC) = 0.3 Prob(Yes|RC) = Prob(No|RC) = 0.5 Prob(TC) = 0.7 Prob(Yes) = 0.65 Prob(No) = 0.35

We want to know Prob(Yes|TC) = ?

Solution: Prob(Yes) = Prob(Yes|RC) * Prob(RC) + Prob(Yes|TC) * Prob(TC) 0.65 = 0.5 * 0.3 + Prob(Yes|TC) * 0.7

```
(0.65-0.3*0.5) / 0.7
## [1] 0.7142857
```

### Part B

The probability of having disease given test positive: D: having disease, ND: not having disease P: getting postivie in the test, N: getting negative in the test

P(D|P) = P(P|D) * P(D) / ( P(P|D) * P(D) + P(P|N) * P(N) )

```
(0.993) * 0.000025 / ((0.993) * 0.000025 + 0.0001 * 0.999975)
## [1] 0.1988824
```

Problem of the test lies in having too many "False Positive". Which is giving too many positive on tests for those who don't have a desease. This kind of implementing a universal testing policy for the disease will lead to panic and chaos.
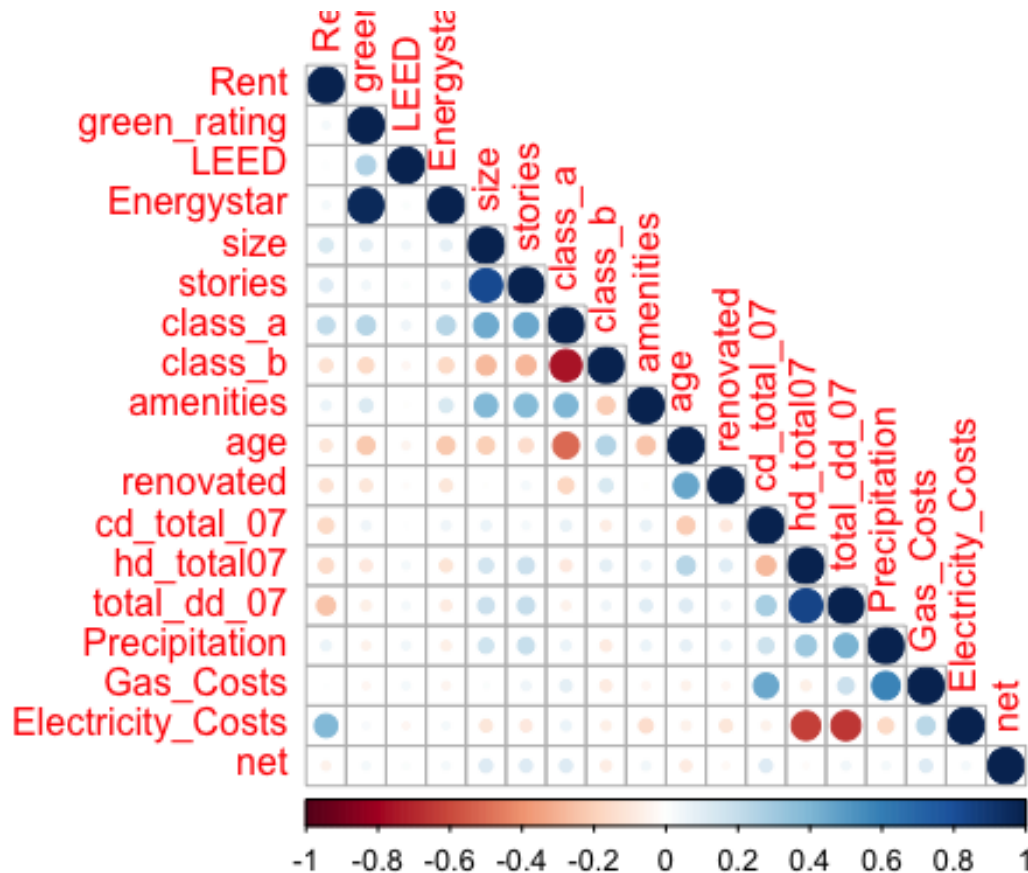
## Problem 2 Green Buildings

```
library(corrplot)
## corrplot 0.84 loaded
library(ggplot2)
plot_data =
read.csv(url("https://raw.githubusercontent.com/jgscott/STA380/master/data/greenbuildings.csv"))
#colnames(plot_data)
re_order_var = c("Rent", "green_rating", "LEED","Energystar", "size",
"stories", "class_a",
 "class_b", "amenities",  "age", "renovated", "cd_total_07",
```

```
"hd_total07","total_dd_07", "Precipitation","Gas_Costs" ,
"Electricity_Costs", "net")
```

After reading our data into R. I reordered the variables by their categories to make correlation plots easier to read. Our new variable list starts with rent, green building certificate, characteristics of a building, weather indicators, ends with energy and other cost.
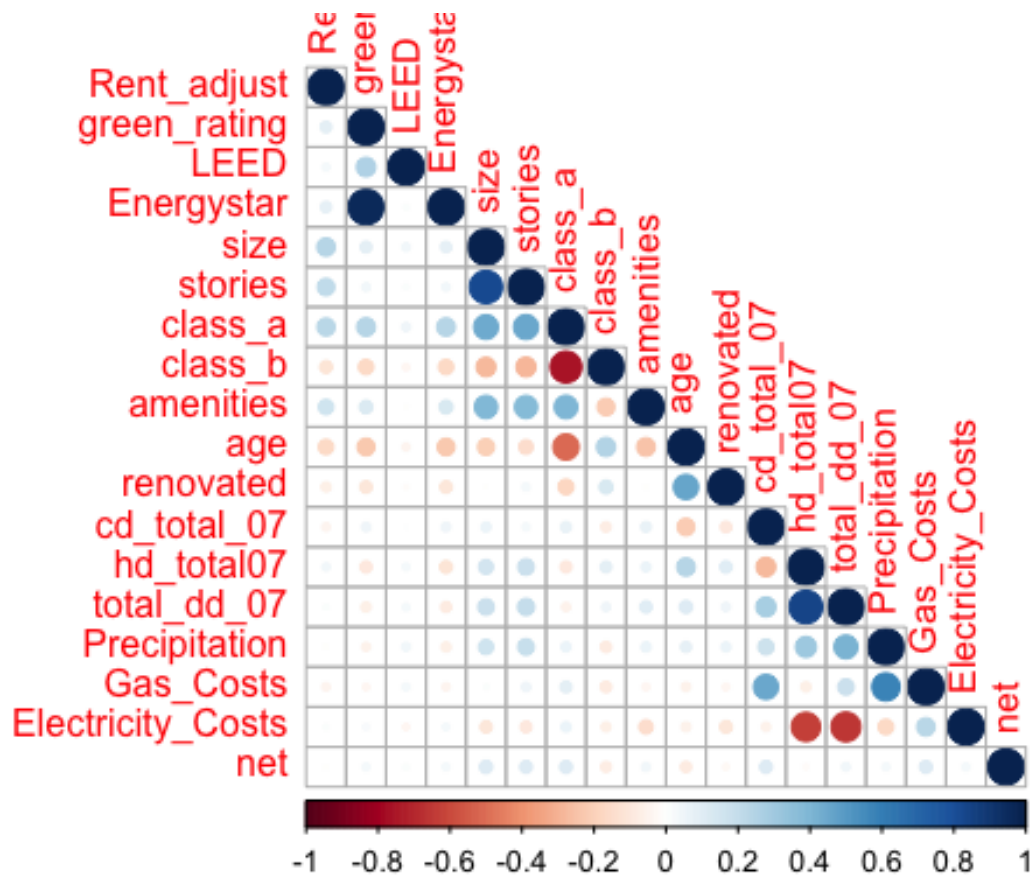
**Variable correlation plot**
```
corrplot(cor(plot_data[re_order_var]), type = "lower")
```



The correlations seems reasonable: rent postiively correlated to good characteristics of a building and negatively correlated to worse building conditions and less favorable weather conditions. However, "cluster_rent" provides aggregated information by cluster. Which cannot be included in our correlation plot but provides valuable information.

Therefore, we created varaiable "Rent_adjust" by dividing rent using its cluster's rent. Which tells us the house's rent proportional to the cluster's rent.

```
plot_data["Rent_adjust"] = plot_data["Rent"] / plot_data["cluster_rent"]
corrplot(cor(plot_data[c("Rent_adjust", re_order_var[-1])]), type = "lower")
```

After replacing Rent using "Rent_adjust", we found that correlations between "Rent_adjust" and weather related variables has dissappeared. Which seems to tell us that the effect of weather on rent are reflected on the cluster's rent.
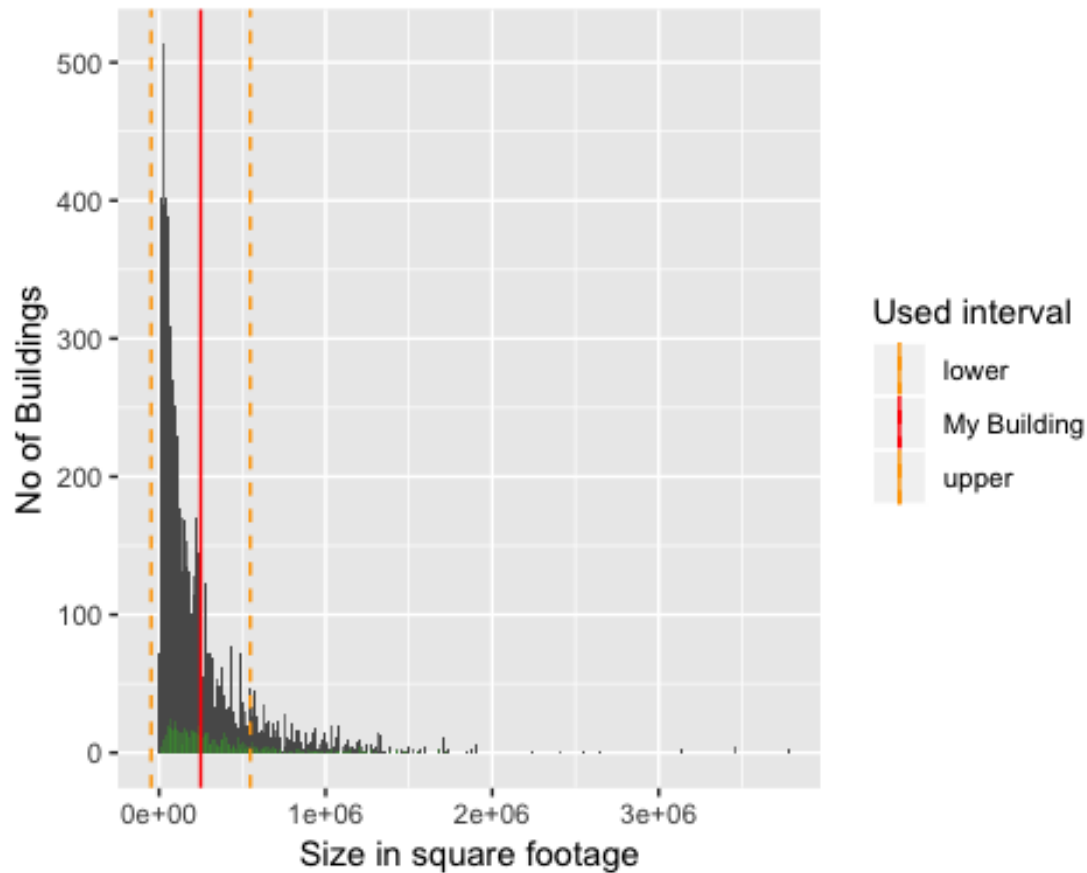
**Examine rent difference between green and non-green buildings having similar characteristic with our new building.**

After we have studied our variables, we look into information we have on our new building : size, age, and stories to help us decide whether we should build a green building.
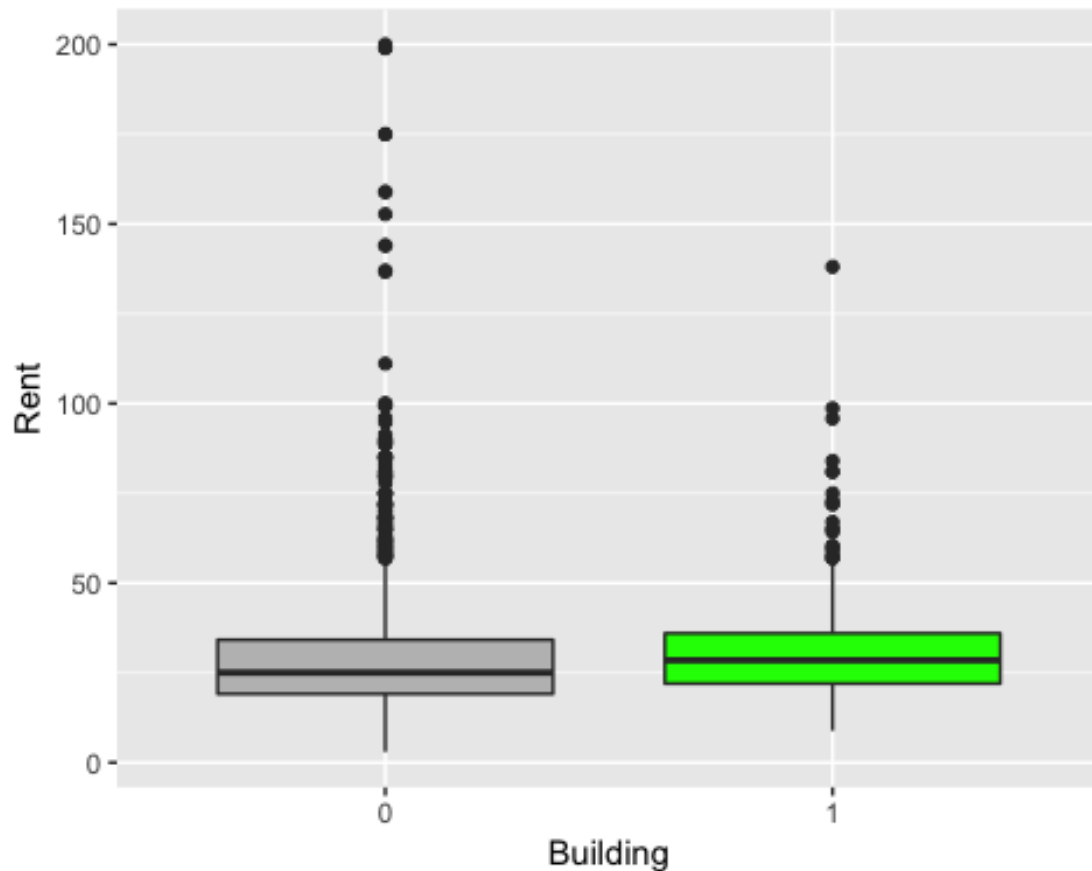
**Size**

```
size_sd = sd(plot_data$size)
ggplot(plot_data, aes(x = size)) +
  geom_histogram(binwidth = 10000) +
  geom_histogram(data=subset(plot_data, green_rating == "1"),binwidth =
10000,fill = "green", alpha = 0.2) +
  labs( x = "Size in square footage", y ="No of Buildings" ) +
  geom_vline(aes(xintercept = 250000, color = "My Building")) +
  geom_vline(aes(xintercept = 250000 + size_sd, color = "upper") , linetype =
"dashed") +
  geom_vline(aes(xintercept = 250000 - size_sd, color = "lower"), linetype =
"dashed") +
```

```
    scale_color_manual(name="Used interval", values=c(`My Building` = "red",
upper = "orange", lower = "orange"))
```
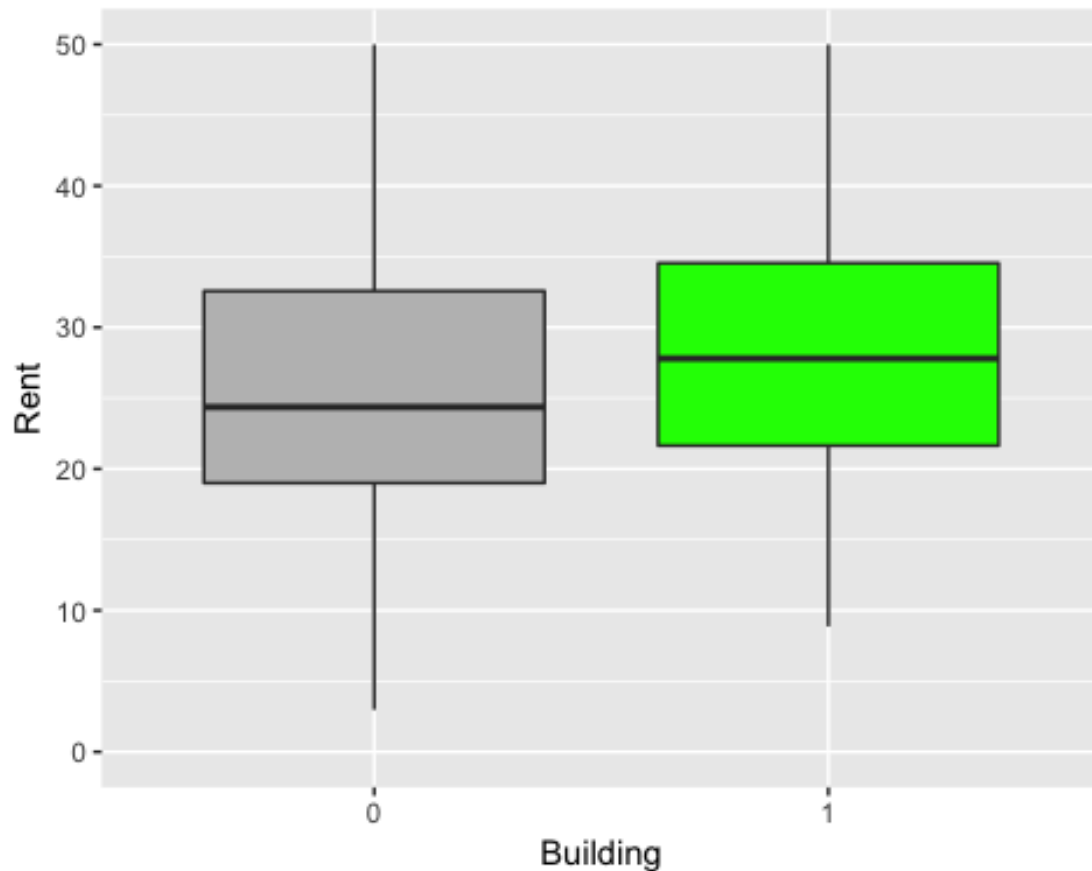


We plot the "size" variable's distribution of our entire data along with the green building's size distribution colored in green. The vertical red line shows the our building's size : 250000. And the orange dash line shows interval of one standard deviation, in which are the data points we are carrying to our boxplot.

```
size_explore = plot_data[plot_data$size <=  250000 + size_sd, ]
size_explore = size_explore[size_explore$size >= 250000 - size_sd, ]
size_explore$green_rating = as.factor(size_explore$green_rating)
ggplot(size_explore, aes(x = green_rating, y = Rent) ) +
  geom_boxplot(fill = c("grey", "green")) +
  #scale_y_continuous(limit = c(0, 50)) +
  labs(x = "Building", y = "Rent")
```

Using the data points having less than one standard deviation on "size" with our new building, we plotted the rent difference between green and non-green buildings.

```
size_explore = plot_data[plot_data$size <=  250000 + size_sd, ]
size_explore = size_explore[size_explore$size >= 250000 - size_sd, ]
size_explore$green_rating = as.factor(size_explore$green_rating)
ggplot(size_explore, aes(x = green_rating, y = Rent) ) +
  geom_boxplot(fill = c("grey", "green")) +
  scale_y_continuous(limit = c(0, 50)) +  # 363 among 6998 have rent > 50
  labs(x = "Building", y = "Rent")
```
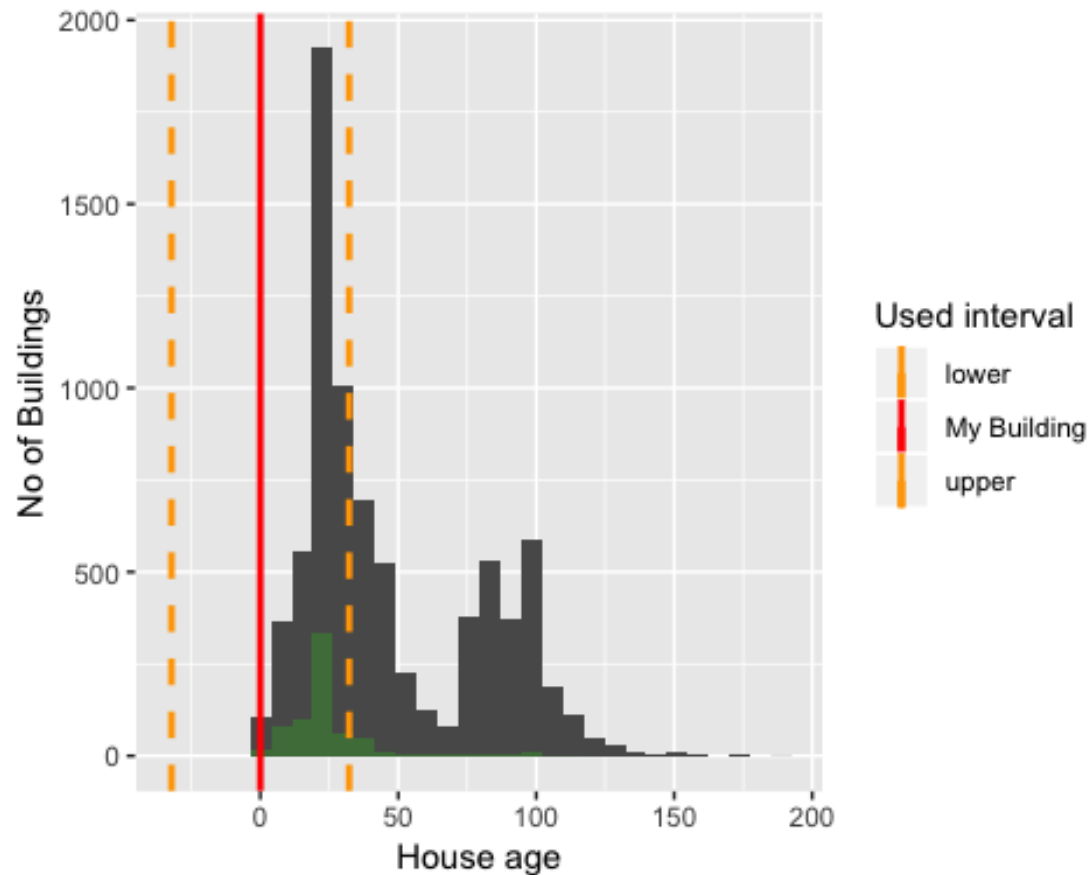
We zoom in our plot by resetting Y limit. We found that 25, 50, 75 rent quantile of green buildings are all higher than non-green buildings.
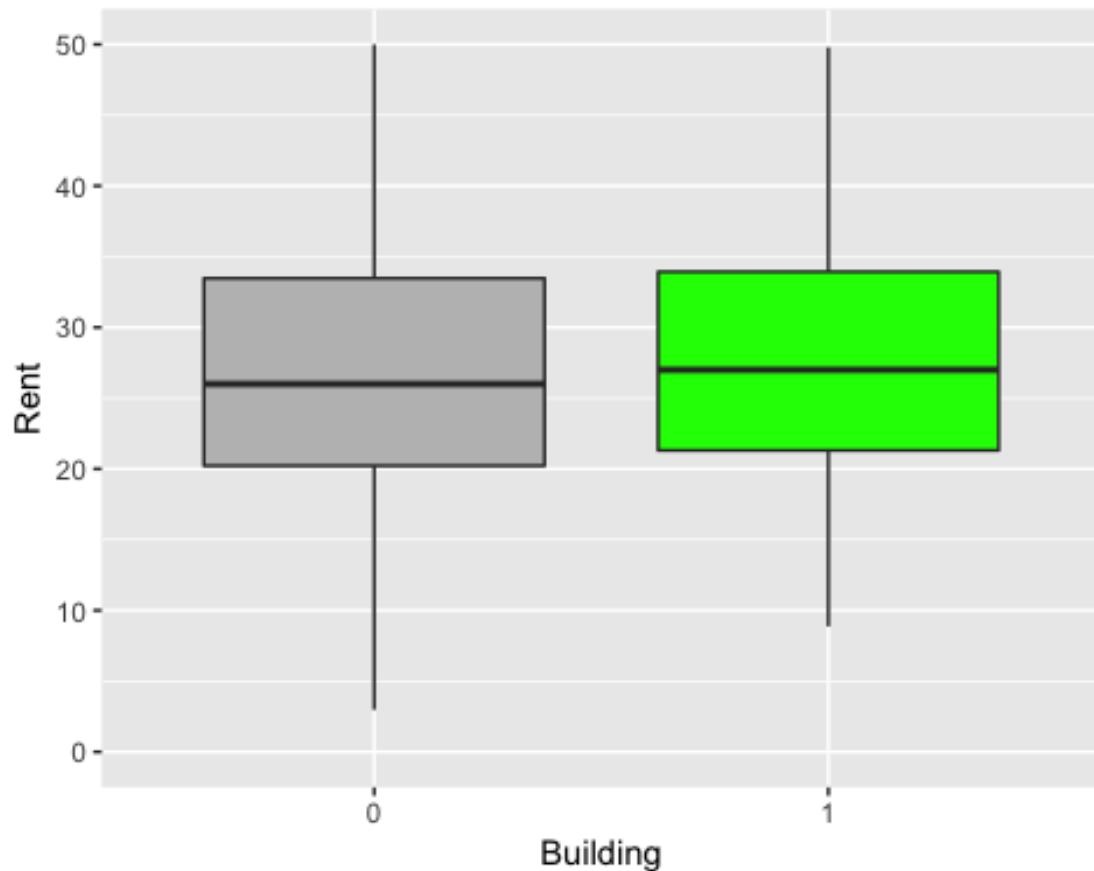
### Age

```
age_sd = sd(plot_data$age)
ggplot(plot_data, aes(x = age)) +
  geom_histogram() +
  geom_histogram(data=subset(plot_data, green_rating == "1"),fill = "green",
alpha = 0.2) +
  labs( x = "House age", y ="No of Buildings" ) +
  geom_vline(aes(xintercept = 0, color = "My Building"), size = 1) +
  geom_vline(aes(xintercept = 0 + age_sd, color = "upper") , linetype =
"dashed", size = 1) +
  geom_vline(aes(xintercept = 0 - age_sd, color = "lower"), linetype =
"dashed", size = 1) +
  scale_color_manual(name="Used interval", values=c(`My Building` = "red",
                                                    upper = "orange",
                                                    lower = "orange"))
```

We plot the "age" variable's distribution along with "age"'s distribution among green buildings. We found new buildings (age = 0) only consist of a small portion of our data.

```
age_explore = plot_data[plot_data$age <=  age_sd, ]
age_explore$green_rating = as.factor(age_explore$green_rating)
ggplot(age_explore, aes(x = green_rating, y = Rent)) +
    geom_boxplot(fill = c("grey", "green")) +
    scale_y_continuous(limit = c(0, 50)) + #173 among 3740 have rent > 50
    labs(x = "Building", y = "Rent")
```
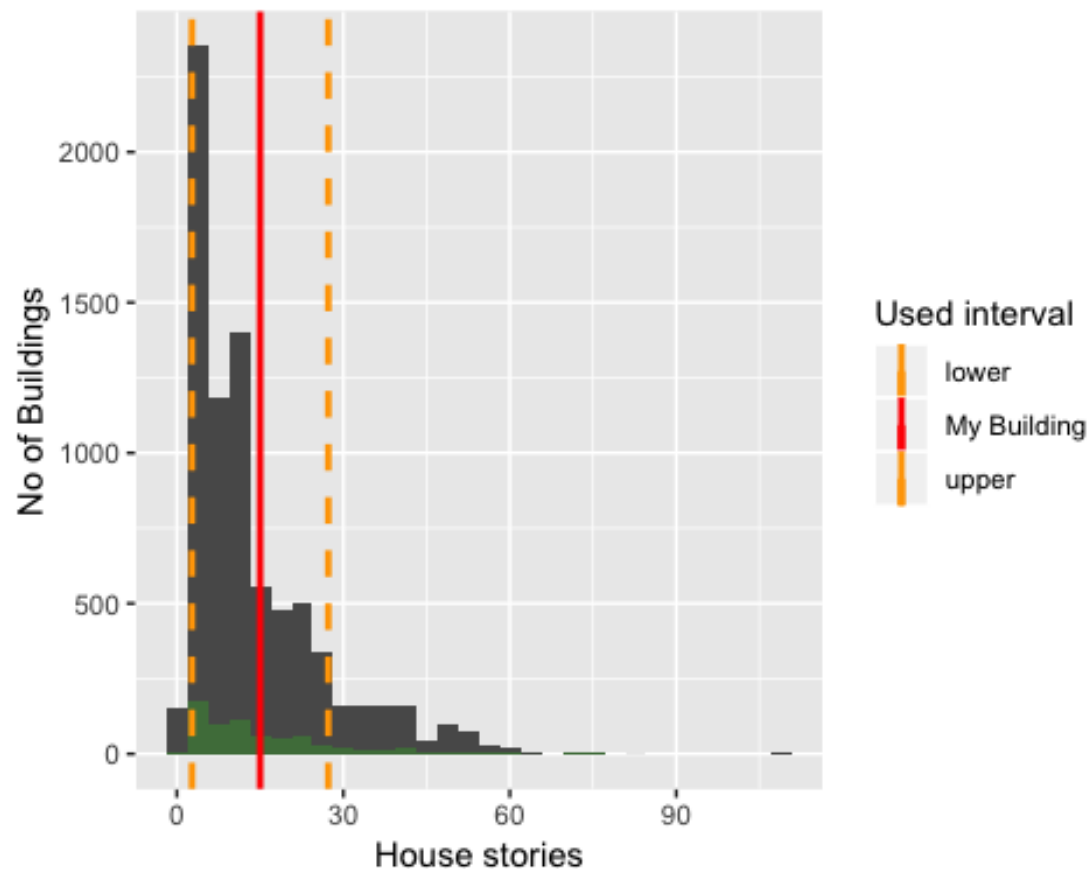
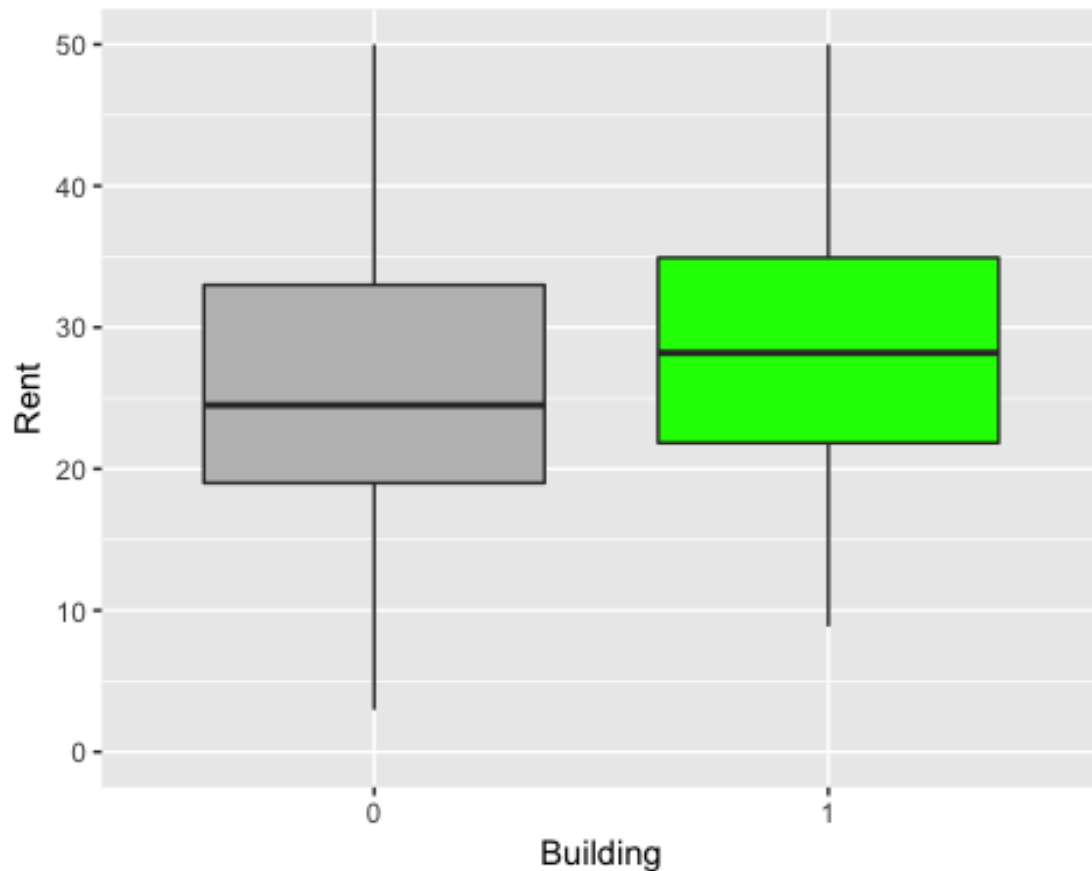Green buildings have higher rents in relatively new buildings.

### Stories

```
stories_sd = sd(plot_data$stories)
ggplot(plot_data, aes(x = stories)) +
  geom_histogram() +
  geom_histogram(data=subset(plot_data, green_rating == "1"),fill = "green",
alpha = 0.2) +
  labs( x = "House stories", y ="No of Buildings" ) +
  geom_vline(aes(xintercept = 15, color = "My Building"), size = 1) +
  geom_vline(aes(xintercept = 15 + stories_sd, color = "upper") , linetype =
"dashed", size = 1) +
  geom_vline(aes(xintercept = 15 - stories_sd, color = "lower"), linetype =
"dashed", size = 1) +
  scale_color_manual(name="Used interval", values=c(`My Building` = "red",
                                          upper = "orange",
                                          lower = "orange"))
```

Under the distribution of "stories", we can see that our new buildings story is around 70 quantile among our data set.

```
stories_explore = plot_data[plot_data$stories <= 15 + stories_sd, ]
stories_explore = stories_explore[stories_explore$stories >= 15 - stories_sd,
]
stories_explore$green_rating = as.factor(stories_explore$green_rating)
ggplot(stories_explore, aes(x = green_rating, y = Rent)) +
  geom_boxplot(fill = c("grey", "green")) +
  scale_y_continuous(limit = c(0, 50)) +   #173 among 3740 have rent > 50
  labs(x = "Building", y = "Rent")
```

Green buildings can bring higher rent in our stories interval.

## Validation using regression

On plots above we obsereved higher rents from green buildings under different condition controled. Before coming up with our final advice, it's always good to check the result of linear regression, which is good at controling variables we consider might affect our rent.

We used Rent_adjust as dependent variable to control the effect from clusters. Then added class_a to control the raised rent from high-quality buildings .

```
summary(lm( Rent_adjust ~ log(size) + age + stories+ class_a +green_rating ,
            data = plot_data))

##
## Call:
## lm(formula = Rent_adjust ~ log(size) + age + stories + class_a +
##     green_rating, data = plot_data)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -0.9002 -0.1413 -0.0251  0.0899  3.4442
##
## Coefficients:
```

```
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.6550154  0.0454151  14.423  < 2e-16 ***
## log(size)      0.0307891  0.0042013   7.328 2.56e-13 ***
## age           -0.0005894  0.0001151  -5.121 3.11e-07 ***
## stories        0.0014500  0.0004035   3.594 0.000328 ***
## class_a        0.0528742  0.0088073   6.003 2.02e-09 ***
## green_rating   0.0333510  0.0117117   2.848 0.004416 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2797 on 7888 degrees of freedom
## Multiple R-squared:  0.08143,    Adjusted R-squared:  0.08085
## F-statistic: 139.9 on 5 and 7888 DF,  p-value: < 2.2e-16
```

Green building provides 3.3% increase on local cluster rent per square foot. (Our building 250,000 square foot, extra cost for green: 5,000,000)

## Problem 3 Bootstrapping

```
library(mosaic)
library(quantmod)
library(foreach)
library(GGally)

# Import the five ETFs
mystocks = c("SPY", "TLT", "LQD", "EEM", "VNQ")
myprices = getSymbols(mystocks, from = "2007-01-01")
getSymbols(mystocks)
## [1] "SPY" "TLT" "LQD" "EEM" "VNQ"

#Get the adjusted closing price of each ETF
for(ticker in mystocks) {

  expr = paste0(ticker, "a = adjustOHLC(", ticker, ")")
  eval(parse(text=expr))
}

# Combine all the returns in a matrix
all_returns = cbind(ClCl(LQDa),ClCl(TLTa),ClCl(SPYa), ClCl(VNQa), ClCl(EEMa))
all_returns = as.matrix(na.omit(all_returns))

head(all_returns)
##                ClCl.LQDa     ClCl.TLTa     ClCl.SPYa    ClCl.VNQa
## 2007-01-04   0.0075152938  0.006063328  0.0021221123  0.001296655
## 2007-01-05  -0.0006526807 -0.004352668 -0.0079763183 -0.018518518
## 2007-01-08  -0.0002798843  0.001793566  0.0046250821  0.001451392
## 2007-01-09   0.0001866169  0.000000000 -0.0008498831  0.012648208
## 2007-01-10  -0.0013063264 -0.004475797  0.0033315799  0.012880523
## 2007-01-11  -0.0030832664 -0.005844712  0.0043804651  0.012973706
```
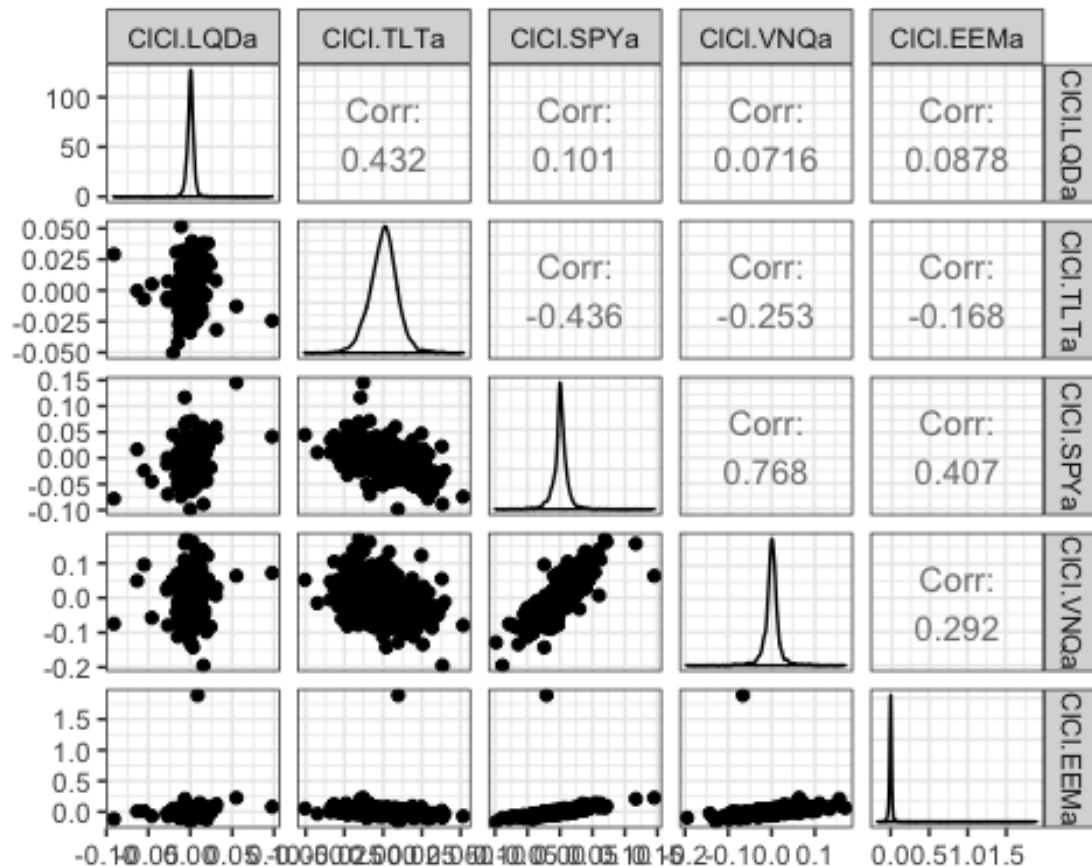
```
##               ClCl.EEMa
## 2007-01-04 -0.013809353
## 2007-01-05 -0.029238205
## 2007-01-08  0.007257535
## 2007-01-09 -0.022336235
## 2007-01-10 -0.002303160
## 2007-01-11  0.012650112
```

```
# The sample correlation matrix
cor(all_returns)
```

```
##              ClCl.LQDa  ClCl.TLTa  ClCl.SPYa   ClCl.VNQa   ClCl.EEMa
## ClCl.LQDa 1.00000000  0.4323732  0.1013347  0.07156075  0.08784764
## ClCl.TLTa 0.43237319  1.0000000 -0.4362148 -0.25332123 -0.16758098
## ClCl.SPYa 0.10133468 -0.4362148  1.0000000  0.76813129  0.40676425
## ClCl.VNQa 0.07156075 -0.2533212  0.7681313  1.00000000  0.29228612
## ClCl.EEMa 0.08784764 -0.1675810  0.4067643  0.29228612  1.00000000
```

From the correlation plot, we can see that SPY is highly correlated with VNQ and EEM which indicates that emerging-market and real estate relatively more dominant by the market condition. Also, we noticed that TLT is negatively correlated with SPY. This can be explained as when the stock market goes up, people tend to sell treasury bonds to buy stocks, and when it goes down, people buy treasury bonds because they are safe.

```
ggpairs(as.data.frame(all_returns))+theme_bw()
```

From the mean and variance of each ETFs, we can see that EEM, VNQ have relatively higher return, but higher risk consequently. Thus, we can then assume EEM and VNQ is a risker investment compare to the other ETFs. This is also very intuitive because investing in bonds is less risky than investing in emerging-market or real estate which are both highly volatile.

```
#The mean and variance
sort(apply(all_returns,2,mean),decreasing = TRUE)

##     ClCl.EEMa     ClCl.VNQa     ClCl.SPYa     ClCl.TLTa     ClCl.LQDa
## 0.0009813682 0.0004157265 0.0003981241 0.0002787785 0.0002095494

cat("\n")

sort(apply(all_returns,2,var),decreasing = TRUE)

##     ClCl.EEMa     ClCl.VNQa     ClCl.SPYa     ClCl.TLTa     ClCl.LQDa
## 1.616072e-03 4.467609e-04 1.546740e-04 8.369447e-05 2.718333e-05

par(mfrow=c(2,3))
hist(all_returns[SPYa], 30, xlim=c(-0.04, 0.04))
hist(all_returns[TLTa], 30, xlim=c(-0.04, 0.04))
hist(all_returns[LQDa], 30, xlim=c(-0.04, 0.04))
```
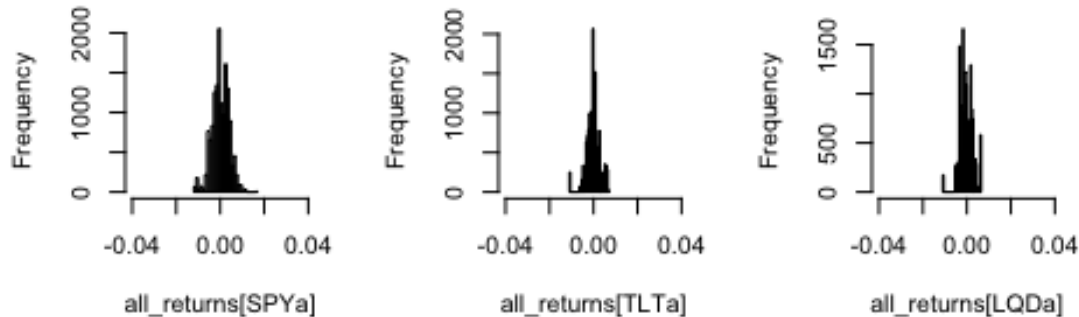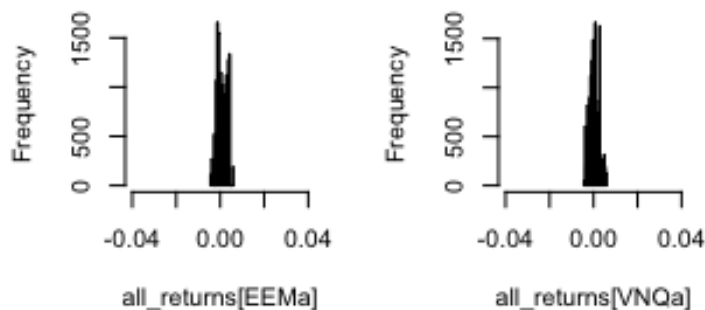
```
hist(all_returns[EEMa], 30, xlim=c(-0.04, 0.04))
hist(all_returns[VNQa], 30, xlim=c(-0.04, 0.04))
```

**Histogram of all_returns[S|Histogram of all_returns[T|Histogram of all_returns[LO**



**Histogram of all_returns[E|Histogram of all_returns[V|**



**Evenly Weighted Portfolio** First, we construct a portfolio which the initial wealth is evenly distributed into each ETFs.

```
set.seed(1)
initial_wealth = 100000
sim_even_split = foreach(i=1:5000, .combine='rbind') %do% {
  total_wealth = initial_wealth
  weights = c(0.2, 0.2, 0.2, 0.2, 0.2) #seperate evenly
  holdings = weights * total_wealth
  n_days = 20
  wealthtracker = rep(0, n_days)
  for(today in 1:n_days) {
    return.today = resample(all_returns, 1, orig.ids=FALSE)
    holdings = holdings + holdings*return.today #p+r(stock)
    total_wealth = sum(holdings)#p+r(portfolio)
    wealthtracker[today] = total_wealth
    holdings = total_wealth * weights #daily rebalance
  }
  wealthtracker
}
```

With an even weighted portfolio, we have a 0.9% return and 20-day 5% VaR at -6300.89

```
even_mean <- mean(sim_even_split[,n_days])
even_return <- (mean(sim_even_split[,n_days])-initial_wealth)/initial_wealth
even_5_perc <- quantile(sim_even_split[,n_days], 0.05) - initial_wealth

# Print out results
cat("For an even split investment of $100,000,\nExpected holdings:",
even_mean, "\nReturn Rate:", even_return, "\n5% VaR:", even_5_perc)

## For an even split investment of $100,000,
## Expected holdings: 101009.1
## Return Rate: 0.01009115
## 5% VaR: -5941.083
```

**Safer Portfolio** Then, according to the mean and variance analysis of each ETFs, we want to construct a more risk-adverse portfolio by weighted toward the less-volatile ETFs for instance, SPY,TLT and LQD. We decided not to put any weight on EEM and VNQ because accoding to the correlation plot, these two ETF has relatively higher correlation with SPY which is the market index. If we put them in the portfolio, it would not diverse and reduce the risk.

```
set.seed(1)
initial_wealth = 100000
sim_safe_split = foreach(i=1:5000, .combine='rbind') %do% {
  total_wealth = initial_wealth
  weights = c(0.5, 0.3, 0.2, 0, 0) # 0.2 SPY, 0.3 TLT, 0.5 LQD
  holdings = weights * total_wealth
  n_days = 20
  wealthtracker = rep(0, n_days)
  for(today in 1:n_days) {
    return.today = resample(all_returns, 1, orig.ids=FALSE)
    holdings = holdings + holdings*return.today
    total_wealth = sum(holdings)
    wealthtracker[today] = total_wealth
    holdings = total_wealth * weights
  }
  wealthtracker
}
```

With a safer portfolio, we have a lower 20-day 5% VaR at -3002.014, however, this also means a lower return at 0.5%.

```
safe_mean <- mean(sim_safe_split[,n_days])
safe_return <- (mean(sim_safe_split[,n_days])-initial_wealth)/initial_wealth
safe_5_perc <- quantile(sim_safe_split[,n_days], 0.05) - initial_wealth

# Print out results
cat("For an safer split investment of $100,000,\nExpected
```

```
holdings:",safe_mean, "\nReturn Rate:", safe_return, "\n5% VaR:",
safe_5_perc)

## For an safer split investment of $100,000,
## Expected holdings: 100544.5
## Return Rate: 0.005445058
## 5% VaR: -2784.43
```

**Aggressive Portfolio** For the aggressive portfolio, we merely place our initial wealth on the risker investment EEM and VNQ. We did not place any holdings on the other ETFs because if we added them in the portfolio, it would diversify the risk and make the portfolio less aggressive.

```
set.seed(1)
initial_wealth = 100000
sim_risky_split = foreach(i=1:5000, .combine='rbind') %do% {
  total_wealth = initial_wealth
  weights = c(0, 0, 0, 0.3, 0.7) #0.7 EEM, 0.3 VNQ
  holdings = weights * total_wealth
  n_days = 20
  wealthtracker = rep(0, n_days)
  for(today in 1:n_days) {
    return.today = resample(all_returns, 1, orig.ids=FALSE)
    holdings = holdings + holdings*return.today
    total_wealth = sum(holdings)
    wealthtracker[today] = total_wealth
    holdings = total_wealth * weights
  }
  wealthtracker
}
```

With an aggressive portfolio, we got a 1.65% return but follow by a really high 20-days 5% VaR at -12596.79

```
risky_mean <- mean(sim_risky_split[,n_days])
risky_return <- (mean(sim_risky_split[,n_days])-
initial_wealth)/initial_wealth
risky_5_perc <- quantile(sim_risky_split[,n_days], 0.05) - initial_wealth

# Print out results
cat("For a risky split investment of $100,000,\nExpected holdings:",
risky_mean, "\nReturn Rate:", risky_return, "\n5% VaR:", risky_5_perc)

## For a risky split investment of $100,000,
## Expected holdings: 101864.7
## Return Rate: 0.01864712
## 5% VaR: -12422.34
```

With the histogram below, again, we prove that portfolio that place higher proportions on ETFs with higher variance could be more volatile. The 'risk' portfolio has more chances of
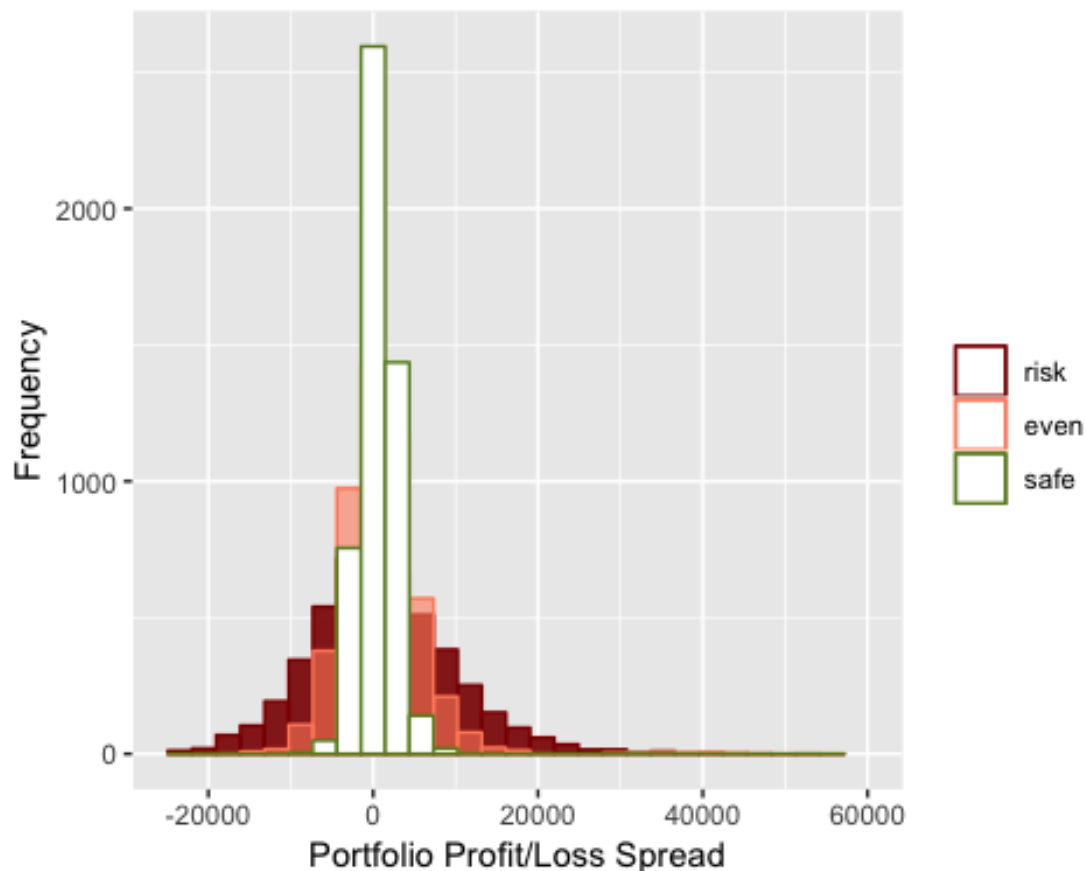
getting higher returns and also more loss. The 'safe' is the lower-volatility portfolio and as expected, it has lower chances of getting big profit or big loss.

**Plot Histogram of Portfolios**

```
even = sim_even_split[,n_days]- initial_wealth
safe = sim_safe_split[,n_days]- initial_wealth
risk = sim_risky_split[,n_days]- initial_wealth

plot_df = data.frame(risk,even, safe)

ggplot(data = plot_df) +
    geom_histogram(aes(risk,color = 'risk'), fill="red4",alpha =
0.9,position="identity")+
    geom_histogram(aes(even,color = 'even'),fill = 'salmon1',alpha =
0.6,position="identity") +
    geom_histogram(aes(safe,color = 'safe'),fill = 'white',alpha =
1.0,position="identity")+
    xlim(c(-25000,60000))+
    xlab('Portfolio Profit/Loss Spread')+
    ylab('Frequency')+
    scale_colour_manual("",
                        breaks = c("risk", "even", "safe"),
                        values = c("salmon1", "red4","olivedrab4"))
```

**Evaluate Portfolio Performance**

```
get_return = function(initial_wealth,weights,all_returns){

    sim_split = foreach(i = 1:5000, .combine = 'rbind') %do% {
        total_wealth = initial_wealth
        holdings = weights * total_wealth
        n_days = 20
        wealthtracker = rep(0, n_days)
        for(today in 1:n_days) {
            return.today = resample(all_returns, 1, orig.ids=FALSE)
            holdings = holdings + holdings*return.today #p+r(stock)
            total_wealth = sum(holdings)#p+r(portfolio)
            wealthtracker[today] = total_wealth
            holdings = total_wealth * weights #daily rebalance
        }
        wealthtracker
    }
    return(quantile(sim_split[,n_days], 0.05) - initial_wealth)
}

#weight_even = rep(0.2,5)
#even = get_return(100000,weight_uni,all_returns)
```

Once again, by randomly tuning the weight on those 5 ETFs, we can observe that if placing more holdings on the higher-volatility ETFs, we would expect it to lose more money of the portfolios during the worst 5% of possible 20-day periods.

```r
set.seed(1)
weight = runif(5)
weight = weight / sum(weight)
for (i in 1:10){
    weight = runif(5)
    weight = weight / sum(weight)
    initial_wealth = 100000
    VaR = get_return(initial_wealth,weight,all_returns)
    cat('The portfolio weight is: ')
    print(weight)
    cat('VaR of the portfolio is: ')
    print(VaR)
}

## The portfolio weight is: [1] 0.28120698 0.29569494 0.20683781 0.19692041
0.01933986
## VaR of the portfolio is:        5%
## -4208.322
## The portfolio weight is: [1] 0.06473905 0.45353950 0.19643722 0.00748451
0.27779972
## VaR of the portfolio is:        5%
## -4711.766
## The portfolio weight is: [1] 0.2317894 0.1058172 0.1834601 0.2476991
0.2312342
## VaR of the portfolio is:        5%
## -7038.451
## The portfolio weight is: [1] 0.37746475 0.15186376 0.19061732 0.18721150
0.09284268
## VaR of the portfolio is:        5%
## -4958.259
## The portfolio weight is: [1] 0.001673649 0.352725081 0.217737384
0.140157414 0.287706472
## VaR of the portfolio is:        5%
## -6268.753
## The portfolio weight is: [1] 0.0005496292 0.2491131647 0.3524609624
0.3284036533 0.0694725904
## VaR of the portfolio is:        5%
## -7292.21
## The portfolio weight is: [1] 0.18054499 0.08488587 0.24893533 0.28819775
0.19743606
## VaR of the portfolio is:        5%
## -8010.587
## The portfolio weight is: [1] 0.2102807 0.1492866 0.1785018 0.2136476
0.2482834
## VaR of the portfolio is:        5%
## -6928.782
```
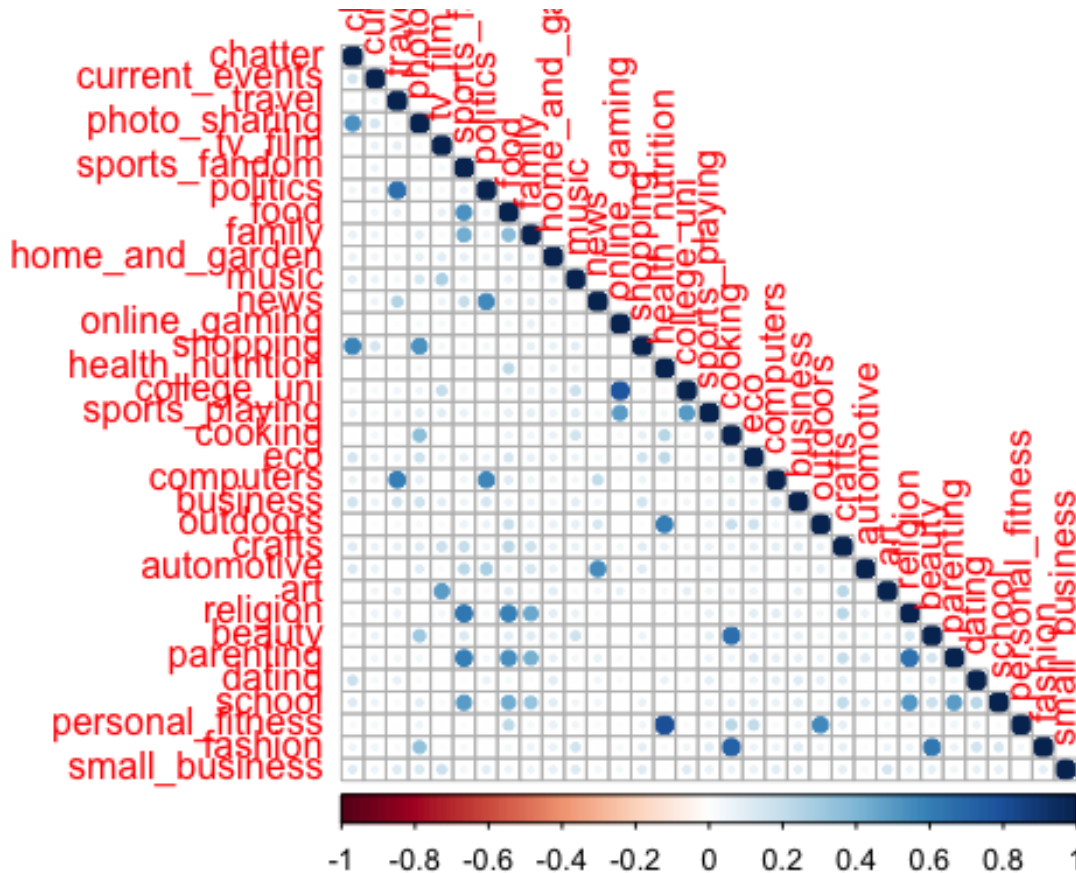
```
## The portfolio weight is: [1] 0.19093565 0.39039666 0.06379952 0.03689110
0.31797707
## VaR of the portfolio is:        5%
## -4800.12
## The portfolio weight is: [1] 0.22974677 0.40298065 0.07596301 0.11979554
0.17151403
## VaR of the portfolio is:        5%
## -4200.626
```

## Problem 4 Market segmentation

```r
library(tidyverse)  # data manipulation
library(cluster)    # clustering algorithms
library(factoextra) # clustering algorithms & visualization
library(corrplot)
```
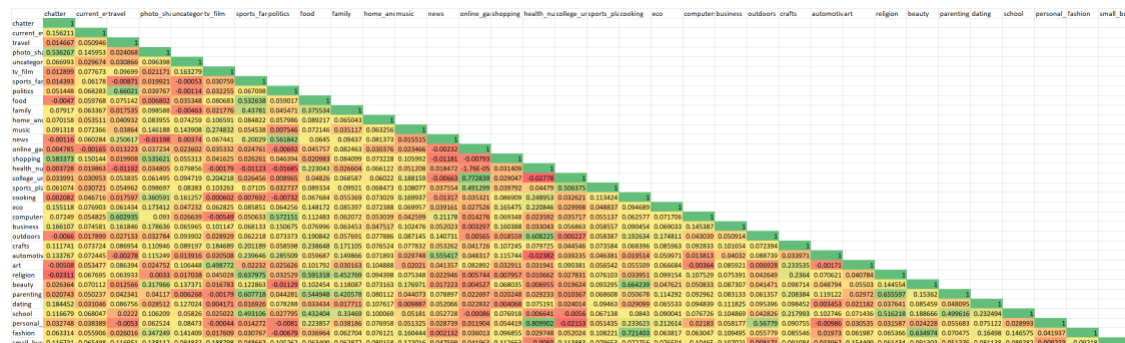
### Data Exploration

```r
# read file
social <-
read.csv(url("https://raw.githubusercontent.com/jgscott/STA380/master/data/so
cial_marketing.csv"), row.names = 1)
#drop columns spam, adult, and uncategorized from file
social = subset(social, select = -c(spam,adult,uncategorized))

corrplot(cor(social), type = "lower")
```

From the correlation plot, we can identidy several variables with high correlation. For example: personal_fitness and health_nutrition.

```r
knitr::include_graphics("https://raw.githubusercontent.com/tiffblahthegiraffe/STA380_HW/master/plot.png")
```



With excel, we looked further into the correlations between variables, and decided to manually group several highly correlated variables together into six hypothetical groups/clusters, as displayed below:

```r
knitr::include_graphics("https://raw.githubusercontent.com/tiffblahthegiraffe/STA380_HW/master/group.png")
```

| influencer | businessman | artists | familyguy | dude | fit |
|---|---|---|---|---|---|
| photo_sharing | travel | tv_film | sports_fandom | online_gaming | health_nutrition |
| shopping | politics | art | religion | college_uni | outdoors |
| current_events | computers | music | parenting | sports_playing | personal_fitness |
| dating | news | crafts | school | | eco |
| cooking | automotive | small_business | food | | |
| beauty | business | | family | | |
| fashion | | | home_and_garden | | |

Let's dig into each group a little more: Group 1 is influencer, whom can be seem as the social media savvy millennials that like sharing their lifestyle online Group 2 represents businessman, a group of corporate people that enjoy politics, news, and business Group 3 are artists that talks about art, music and films on social media Group 4 captures familyguy. These people are very family oriented, with interests including religion, family, and parenting Group 5 is "dude." Think about a college male that likes online gaming and sports. Group 6 represents "fit." These people care about their fitness, shape, nutrition consumption, and hope to stay active

Now, with six hypothetical groups, we want to use clustering to see if our correlation-based groups serve as a group proxy as NutrientH20's market segmentation

### K Means
```
#create hypothetical groups in dataframe
social["influencer"] =
social$chatter+social$photo_sharing+social$shopping+social$current_events+soc
ial$dating + social$cooking+social$beauty+social$fashion
social["businessman"] = social$travel + social$politics +social$computers
+social$news + social$automotive +social$business
social['artists'] = social$tv_film +social$art +social$music +social$crafts
+social$small_business
social['familyguy'] = social$sports_fandom + social$religion
+social$parenting +social$school +social$food +social$family
+social$home_and_garden
social['dude'] = social$online_gaming +social$college_uni
+social$sports_playing
social['fit']
=social$outdoors+social$health_nutrition+social$personal_fitness+social$eco

#select six groups and form new dataframe
social_new = social[, c(34:39)]

#normalize social_new for better comparison between rows
social_norm = social_new/rowSums(social_new)

#scale social_norm
social_scaled <- scale(social_norm, center=TRUE, scale=TRUE)
```
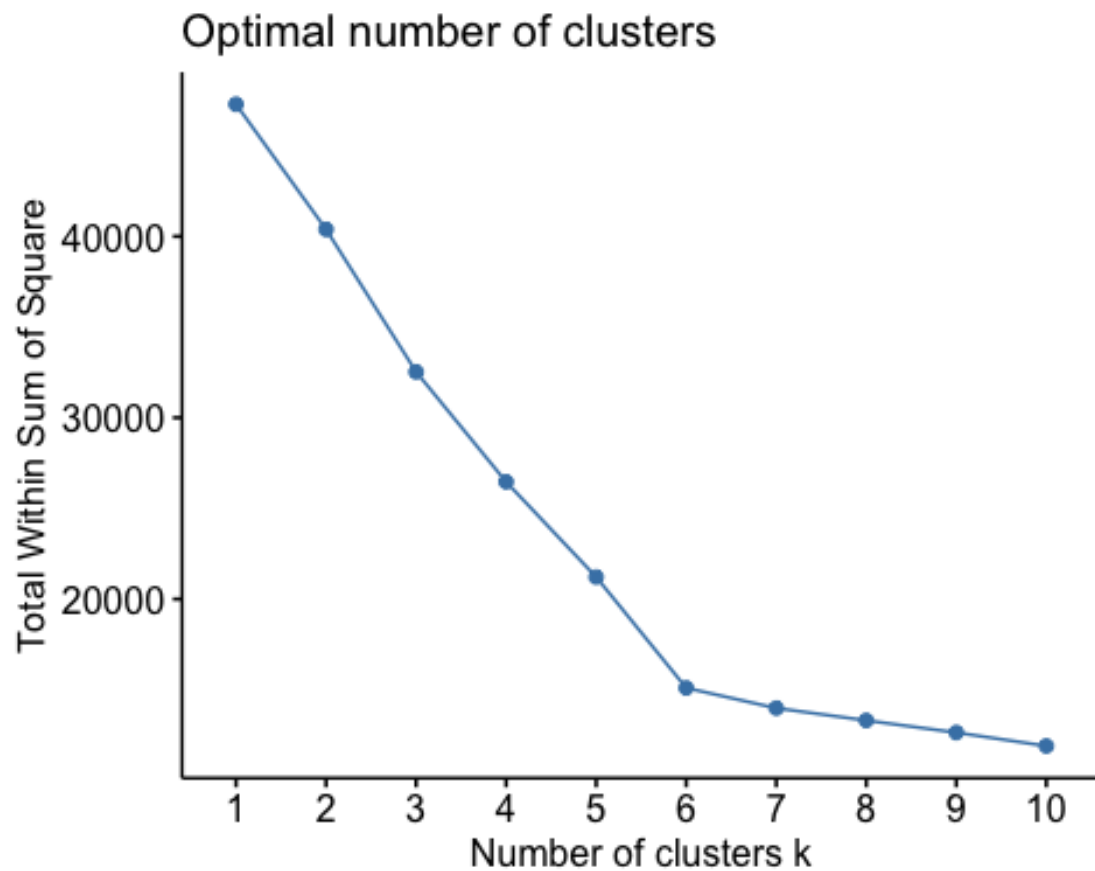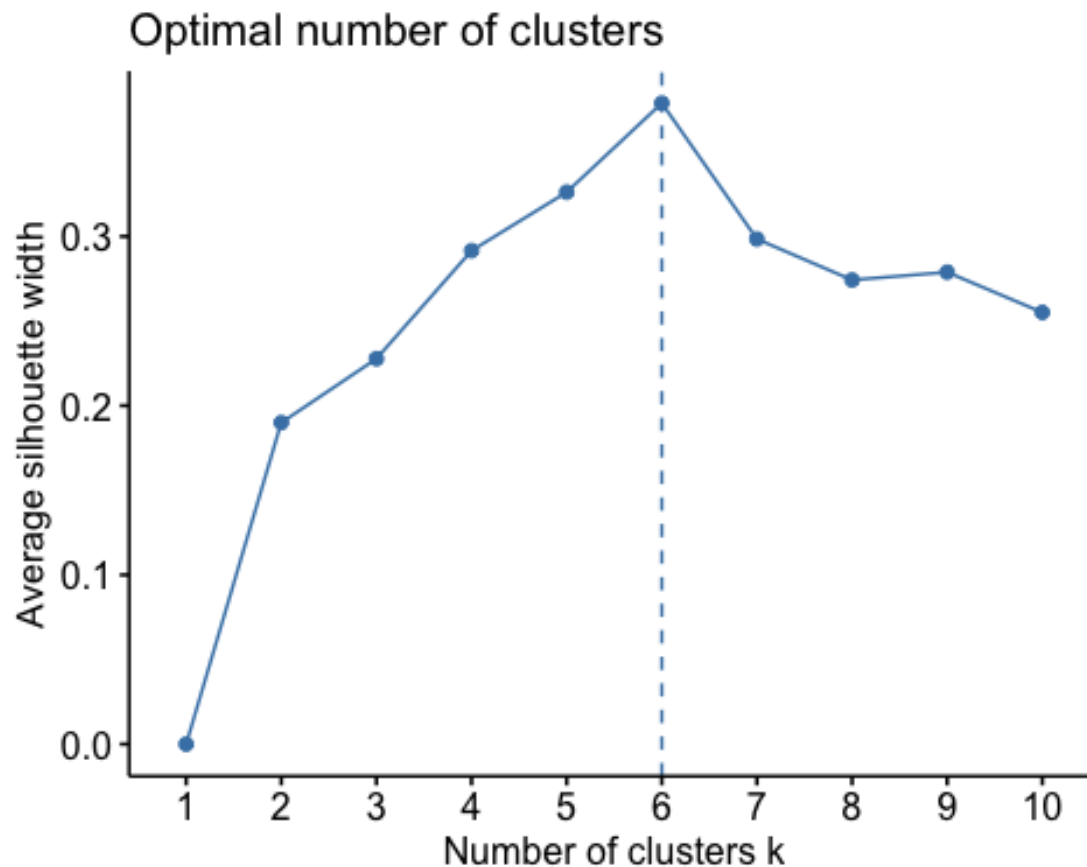
```
#select K
set.seed(123)
fviz_nbclust(social_scaled, kmeans, method = "wss")
```

## Optimal number of clusters



From elbow method, it is concluded that k=6 gives us best clustering result

```
set.seed(123)
fviz_nbclust(social_scaled, kmeans, method = "silhouette")
```

## Optimal number of clusters



Again, K=6 returns best clusting result.

To triple check, we calculate CH

```r
set.seed(123)
for (i in 2:10){
  final <- kmeans(social_scaled, centers = i , nstart = 25)
  B = final$betweenss
  final$withinss
  W = final$tot.withinss
  B/W
  n = nrow(social_scaled)
  k=i
  CH = (B/(k-1))/(W/(n-k))
  cat("k=", k, ", CH:", CH, "\n")

}

## k= 2 , CH: 1652.997
## k= 3 , CH: 1907.703
## k= 4 , CH: 2224.62
## k= 5 , CH: 2747.85
## k= 6 , CH: 3358.518
## k= 7 , CH: 3126.006
```

```
## k= 8 , CH: 2900.054
## k= 9 , CH: 2730.181
## k= 10 , CH: 2609.697
```

When K=6, CH reches its max of 3358.518 From the methods above, we get a preliminary idea that our six manually selected groups might be a good proxy for market segmentation

```
set.seed(123)
#use K=6 to run kmeans
final <- kmeans(social_scaled, centers =6 , nstart = 25)

set.seed(123)
#display centers of seven clusters to see how they are allocated
print(final$centers)

##    influencer businessman    artists  familyguy        dude        fit
## 1 -0.3958139  -0.1868781  2.2283920 -0.2565413  0.02980669 -0.3856151
## 2 -0.5158493  -0.4328407 -0.3624055 -0.3525595 -0.31408302  1.8864271
## 3 -0.6836458  -0.4425627 -0.1639833 -0.4226420  2.83906155 -0.4246151
## 4 -0.6435441  -0.4117771 -0.2058447  2.0067311 -0.30223302 -0.3872227
## 5 -0.7161993   1.9344438 -0.3028690 -0.1884180 -0.28954096 -0.4019451
## 6  1.1261119  -0.3096098 -0.2503628 -0.4000180 -0.25278227 -0.3889869
```

From the result, we can see that our hypothetical grouping method works well! Each of the six groups represents a distinct demographics.

Cluster 1 represents "artists." Variables include: tv_film, art, music, crafts, small_business

Cluster 2 is "fit." These people are the outdoor enthuiast that care about health_nutrition, outdoors, personal_fitness, and eco.

Cluster 3 centers on "dude," a proxy for people who enjoy topics like online_gaming, college_uni, sports_playing.

Cluster 4 centers on "familyguy." This category captures people who are "family oriented" and enojoy taking about topics like sports_fandom, religion, parenting, school, food, family, and home_and_garden.

Cluster 5 centers heavily on the businessman group, which includes people who like talking about topics like travel, travel, politics, computers, news, automotive, business.

Cluster 6 captures "influencers," which can be thought as the millennials that like sharing lifestyle related topics on social media. Topics include:photo_sharing, shopping, current_events, dating, cooking, beauty, fashion.
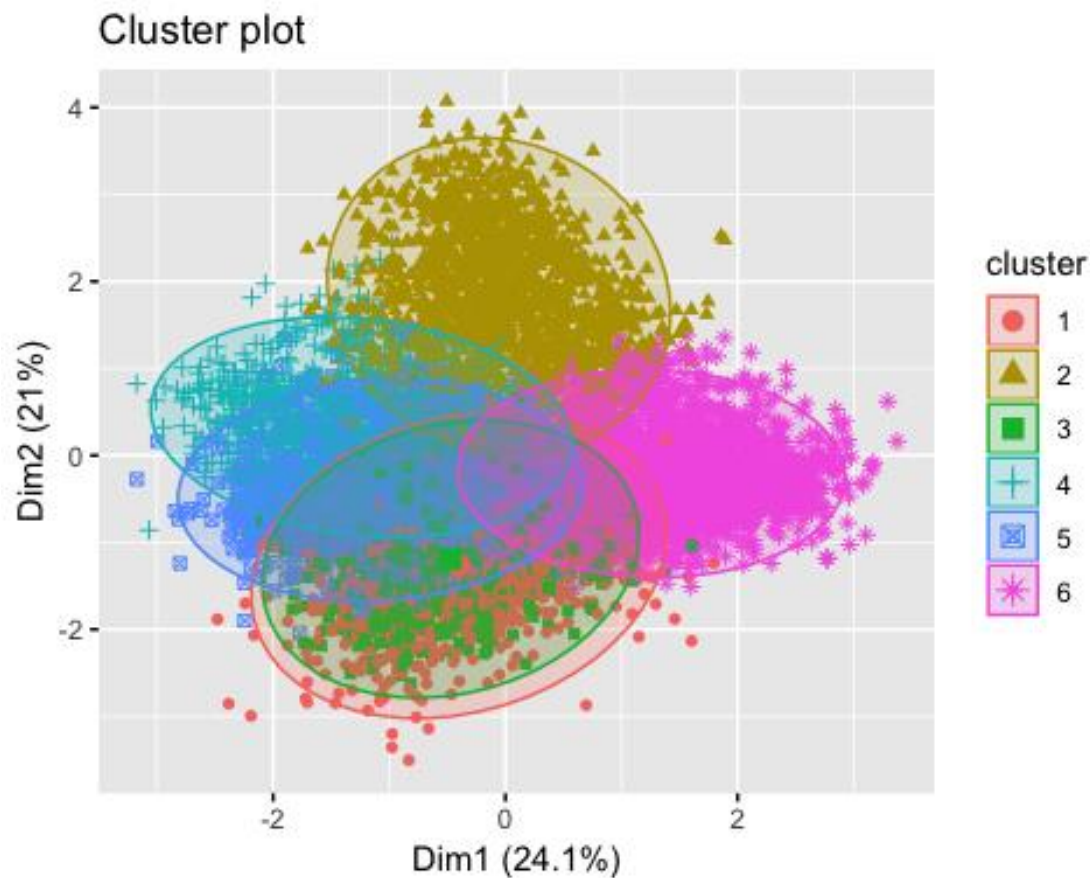
```
#lets see how many people belong in each group with this line of code:
final$size

## [1]  841 1362  627 1133 1205 2714
```

cluster 1 (artists) has 841 people cluster 2 (fit) has 1362 people cluster 3 (dude) represents 627 people cluster 4 (familyguy) includes 1133 people cluster 5 (businessman) has 1205 people cluster 6 (influencer) includes 2714 people

```
#kmeans plot to show 6 clusters
fviz_cluster(final, data = social_scaled, ellipse.type = "norm", stand =
TRUE, geom = "point")
```



The plot shows that our six groups are separated in a fairly clear fashion.

## PCA

Now, let us move on to PCA to further analyze our hypothetical groups.

```
#Run PCA with 5 ranks
pc1 = prcomp(social_norm, scale=TRUE, rank=5)
loadings = pc1$rotation
scores = pc1$x

#several biplots show the first two PCs and how these groups are segmented

q1 = qplot(scores[,1], scores[,2], color= social_norm$influencer ,
xlab='Component 1', ylab='Component 2')
q2 = qplot(scores[,1], scores[,2], color = social_norm$businessman,
```

```
xlab='Component 1', ylab='Component 2')
q3 = qplot(scores[,1], scores[,2], color = social_norm$artists,
xlab='Component 1', ylab='Component 2')
q4 = qplot(scores[,1], scores[,2], color = social_norm$familyguy,
xlab='Component 1', ylab='Component 2')
q5 = qplot(scores[,1], scores[,2], color = social_norm$dude, xlab='Component
1', ylab='Component 2')
q6 = qplot(scores[,1], scores[,2], color = social_norm$fit, xlab='Component
1', ylab='Component 2')
```
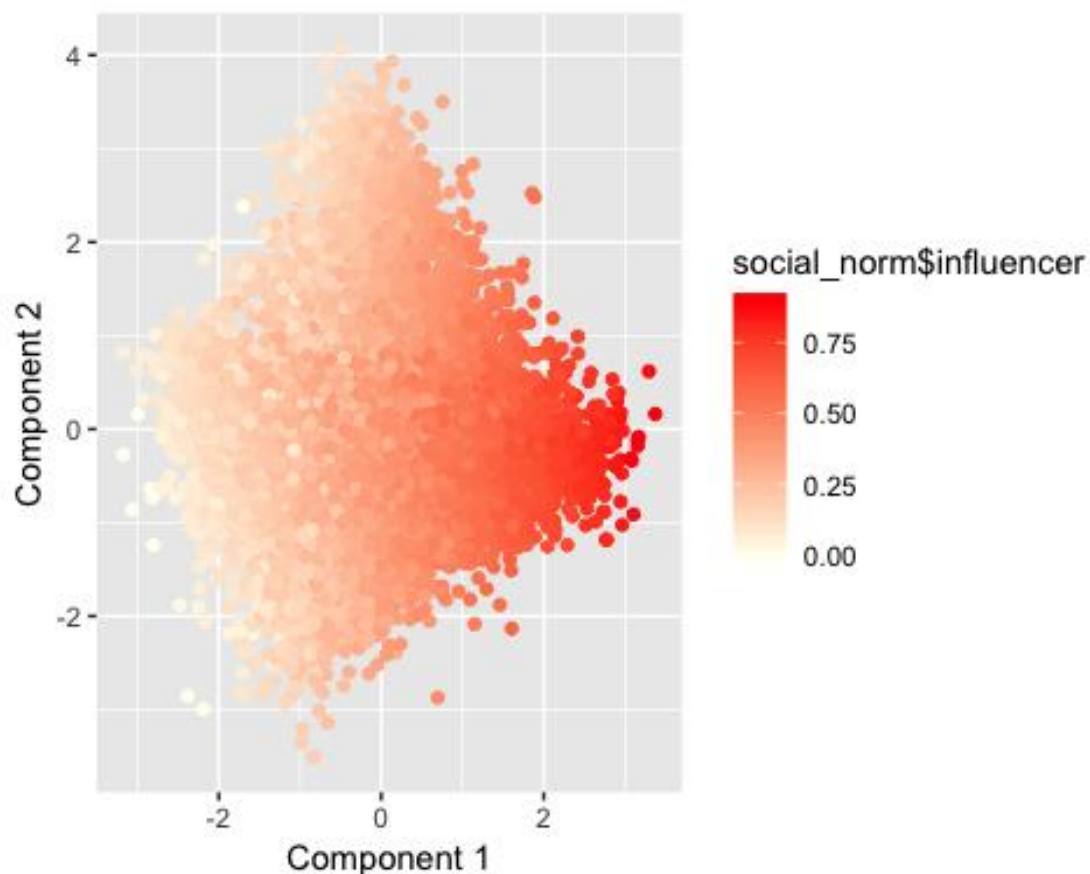
These plots showcase where each hypothetical group belong in the PCA two-domention result. A plot represent a twitter user, and the more red it is, the the more percentage of the user's posts relate to the corresponding group. In other words, the more red it is, the more the user belongs to a hypothetical group we created.

**The influencer**

```
#influencer
q1+scale_color_gradient(low="ivory", high="red")
```
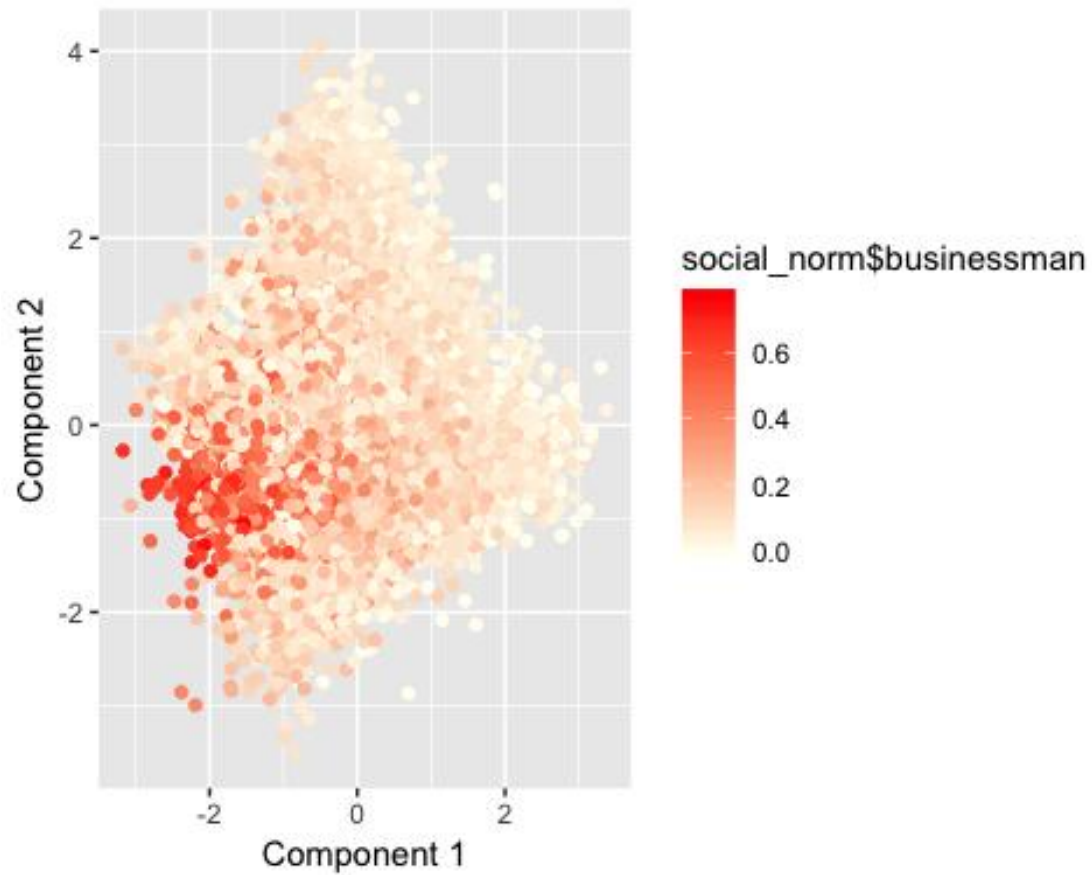


This plot shows where influencer sits in this dimension.
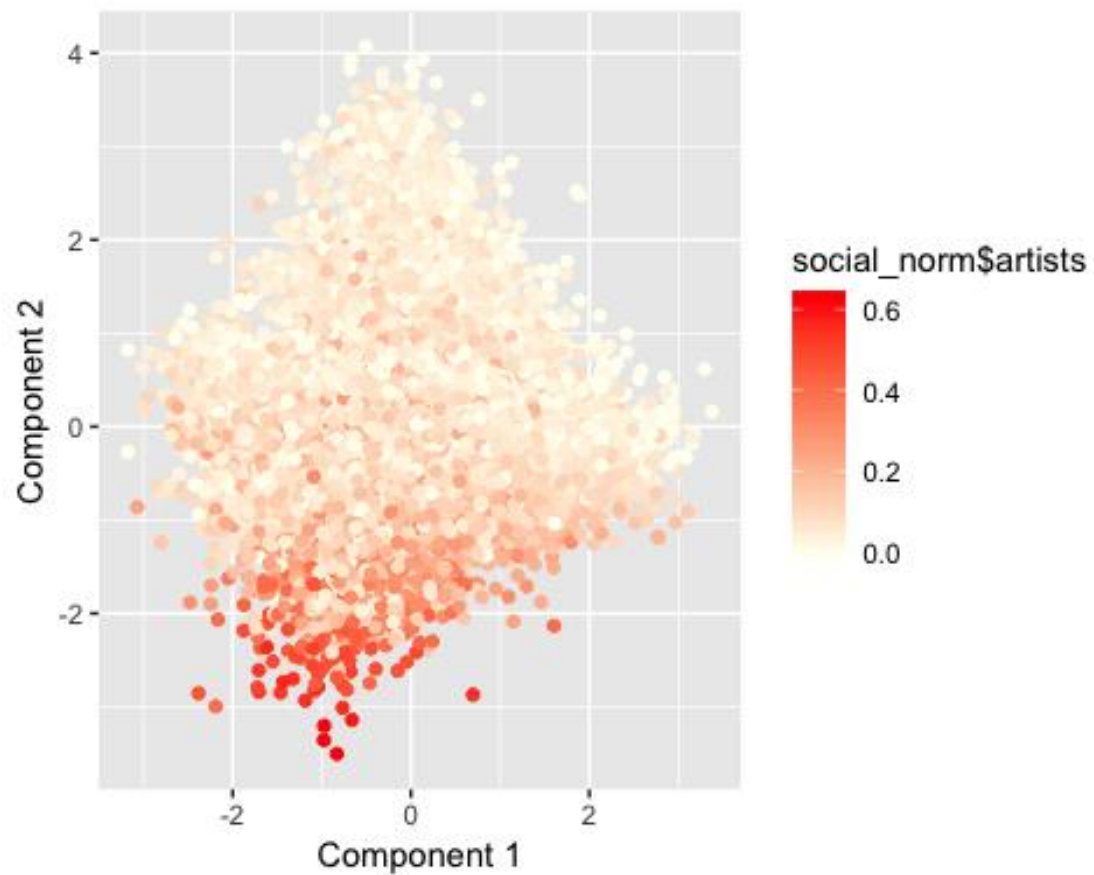
**The businessman**

```
#Businessman
q2+scale_color_gradient(low="ivory", high="red")
```



This plot shows businessman (red) versus others (white).

**The artist**
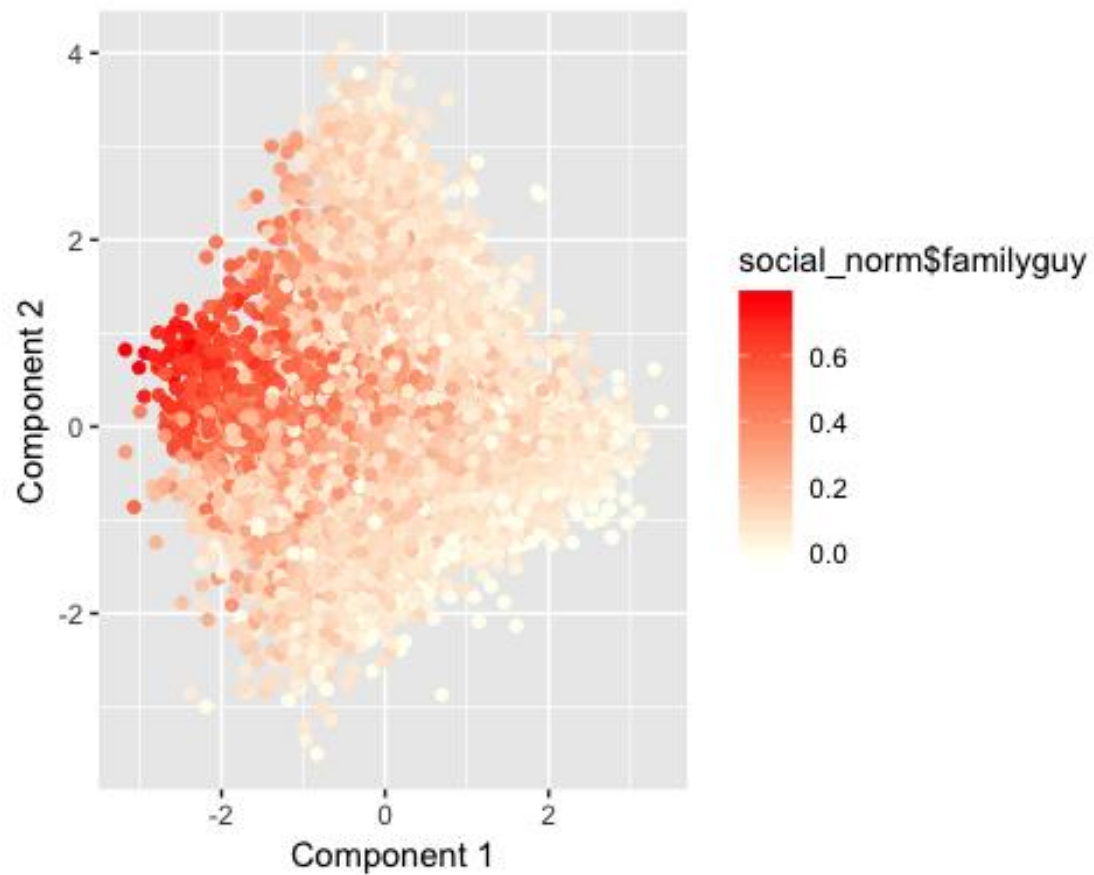
```
#Artists
q3+scale_color_gradient(low="ivory", high="red")
```

This plot points out those where are artists are.
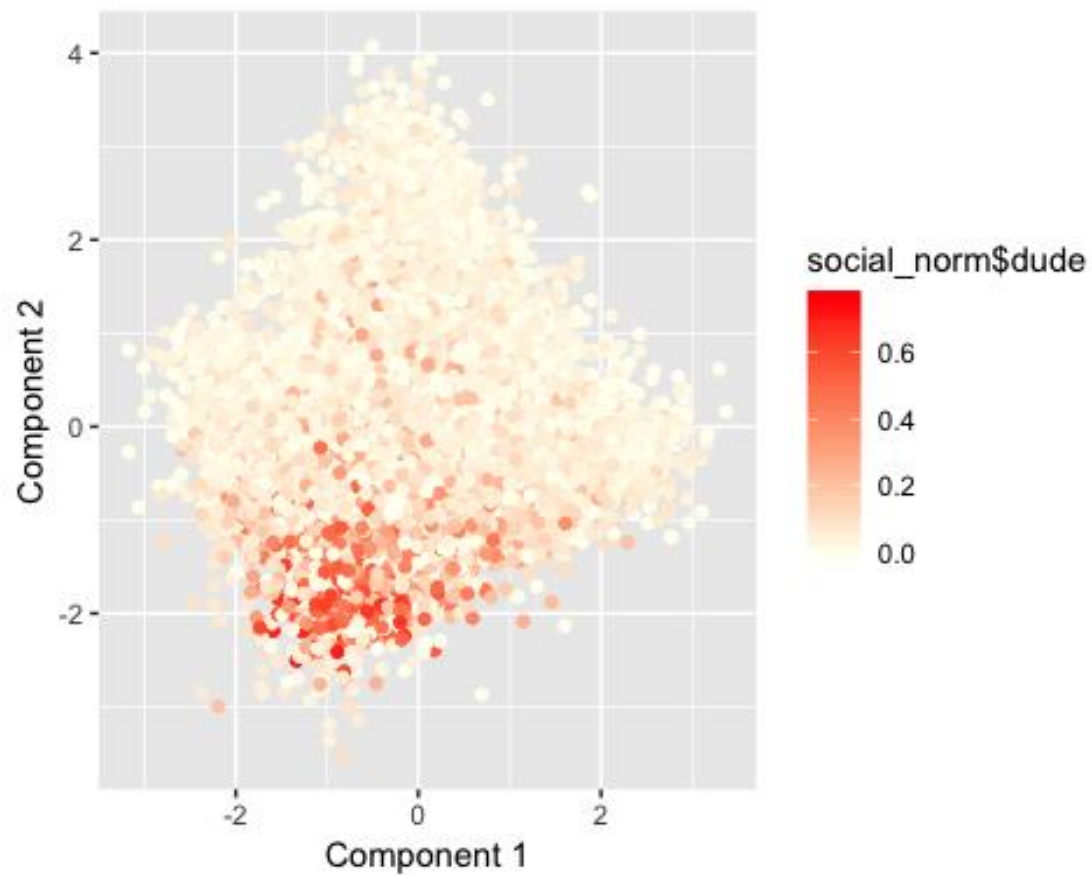
**The familyguy**

```
#Familyguy
q4+scale_color_gradient(low="ivory", high="red")
```

This plot identifies the familyguy group.

**The dude**

```
#Dude
q5+scale_color_gradient(low="ivory", high="red")
```
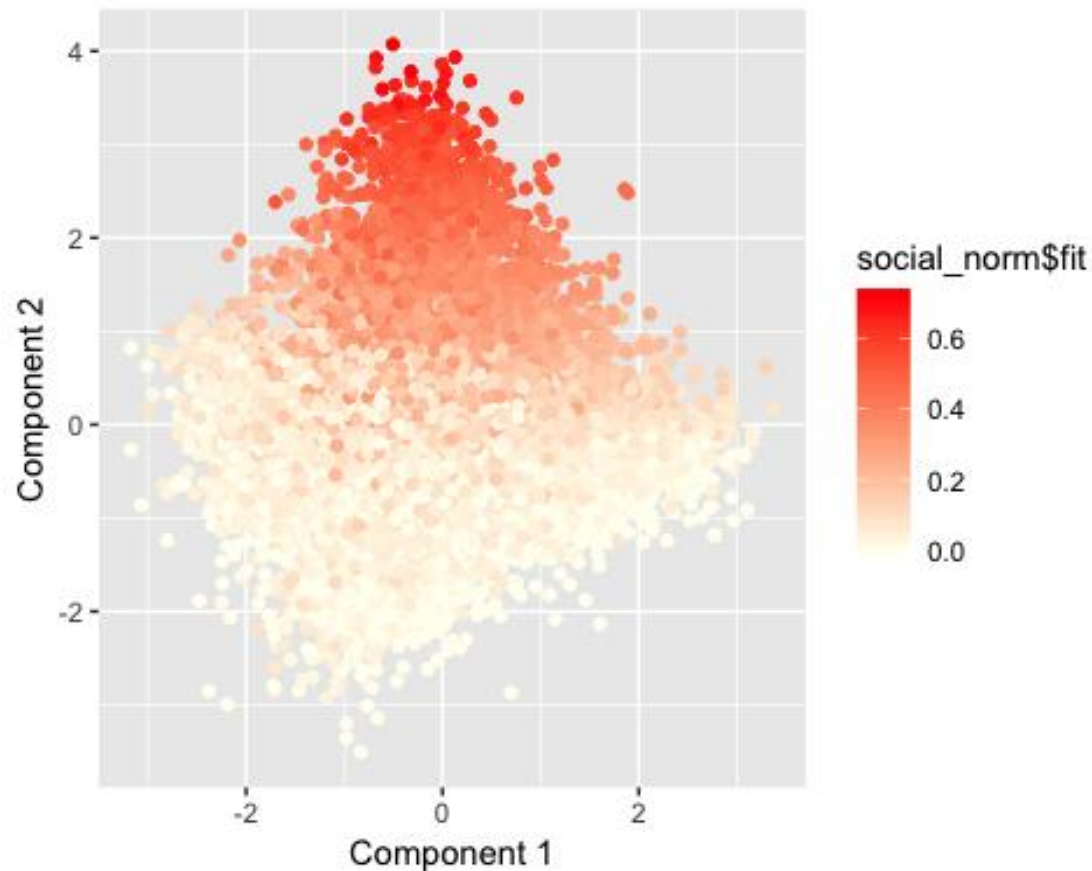
This plot shows where "dude" sit related to others

**The fit**

```
#Fit
q6+scale_color_gradient(low="ivory", high="red")
```

This plot clearly identifies the "fit" group

**Interpretation: component1 does a great job separating influencer (right) from familyguy and businessman (left of graph) Component 2 is great at separating "fit" (upper) from "artists" and "dude" (lower)

```
o1 = order(loadings[,1], decreasing=TRUE)
colnames(social_norm)[head(o1,2)]

## [1] "influencer" "fit"
```

from the formula below, it is clear that PC1, 2 has the ability to separate influencer and fit out from the rest of the data; result aligns with plots

```
loadings

##                     PC1         PC2         PC3         PC4         PC5
## influencer    0.78837083 -0.1120214  0.18897859 -0.18558798 -0.04929072
## businessman  -0.35333354 -0.2033397  0.63810200  0.50052420 -0.02318891
## artists      -0.12683750 -0.4489010 -0.30143877 -0.05915333  0.79131622
## familyguy    -0.47330634  0.1299075  0.08223908 -0.74344082 -0.15779690
## dude         -0.11307903 -0.3699373 -0.61327767  0.24497233 -0.54891098
## fit          -0.02708018  0.7686758 -0.28874700  0.31435302  0.21130124
##                     PC6
```

```
## influencer   0.5416062
## businessman  0.4190132
## artists      0.2487053
## familyguy    0.4180165
## dude         0.3360585
## fit          0.4261025
```

Looking at the vectors, we attain the same results.

## Conclusion

From kmeans clustering and PCA, we can conclude that our hypothecial grouping method works very well in identifying users with different interests on social media, or "socialgraphics." This output can help NutrientH20 better target its audience and focus their social media marketing efforts on a more defined and targeted group of people.

Once again, the groups, variables included, and number of peolpe in each group:

```
knitr::include_graphics("https://raw.githubusercontent.com/tiffblahthegiraffe
/STA380_HW/master/final.png")
```

| influencer | businessman | artists | familyguy | dude | fit |
|---|---|---|---|---|---|
| photo_sharing | travel | tv_film | sports_fandom | online_gaming | health_nutrition |
| shopping | politics | art | religion | college_uni | outdoors |
| current_events | computers | music | parenting | sports_playing | personal_fitness |
| dating | news | crafts | school | | eco |
| cooking | automotive | small_business | food | | |
| beauty | business | | family | | |
| fashion | | | home_and_garden | | |
| 2714 | 1205 | 841 | 1133 | 627 | 1362 |