

Clustering Credit Card Customers for Targeted Marketing using Unsupervised Learning

CSML1000 Project #2, by Group 8

Tamer Hanna tamerh@my.yorku.ca Pete Gray ptgray@my.yorku.ca Xiaohai Lu yu271637@my.yorku.ca
Haofeng Zhou zhf85@my.yorku.ca

```
library(dplyr);  
library(ggplot2);  
library(knitr);  
library(validate);  
library(tidyverse);  # data manipulation  
library(cluster);   # clustering algorithms  
library(clusterSim);  
library(factoextra);  
library(fpc);
```

OVERVIEW

Using data on the behaviour of credit card customers, we can use unsupervised learning to discover market segments that would be useful for targeting marketing strategies.

Our dataset has 19 columns of data on 9000 customers. Using K-means and PCA, we can determine an optimal number of market segments, discover the identifying properties of those segments, and be able to describe to the marketing department what the most significant behaviours of the people in those segments are. Marketing campaigns, then, can be targeted at, for example, impulse shoppers, big spenders, or people who have a hard time paying off their debts.

BUSINESS UNDERSTANDING

Applying specific marketing strategies to different types of customers can improve results and reduce costs. Credit card customers can be profiled by the way they use, and pay off, their card.

Business Objectives

Behavioural data can be mined to discover identifying features of clusters of customers who use their card in similar ways. The marketing department can use these insights to select target groups and optimize the marketing used to target them.

Data Mining Goals

Using shallow algorithm unsupervised learning, we can discover clusters containing customers who exhibit similar behaviours. We can examine those

DATA UNDERSTANDING

Collect Initial Data

Load the data from the local filesystem:

```
#load data
df <- read.csv("credit-card-cust-behav-data.csv", stringsAsFactors = FALSE, header = TRUE, encoding = "UTF-8")
```

Output a quick and dirty summary of the dataset, to ensure that we haven't loaded something scrambled, or the wrong thing:

```
str(df)
```

```
## 'data.frame': 8950 obs. of 18 variables:
## $ CUST_ID                  : chr  "C10001" "C10002" "C10003" "C10004" ...
## $ BALANCE                   : num  40.9 3202.5 2495.1 1666.7 817.7 ...
## $ BALANCE_FREQUENCY         : num  0.818 0.909 1 0.636 1 ...
## $ PURCHASES                 : num  95.4 0 773.2 1499 16 ...
## $ ONEOFF_PURCHASES          : num  0 0 773 1499 16 ...
## $ INSTALLMENTS_PURCHASES    : num  95.4 0 0 0 0 ...
## $ CASH_ADVANCE               : num  0 6443 0 206 0 ...
## $ PURCHASES_FREQUENCY        : num  0.1667 0 1 0.0833 0.0833 ...
## $ ONEOFF_PURCHASES_FREQUENCY: num  0 0 1 0.0833 0.0833 ...
## $ PURCHASES_INSTALLMENTS_FREQUENCY: num  0.0833 0 0 0 0 ...
## $ CASH_ADVANCE_FREQUENCY     : num  0 0.25 0 0.0833 0 ...
## $ CASH_ADVANCE_TRX           : int  0 4 0 1 0 0 0 0 0 ...
## $ PURCHASES_TRX              : int  2 0 12 1 1 8 64 12 5 3 ...
## $ CREDIT_LIMIT                : num  1000 7000 7500 7500 1200 1800 13500 2300 7000 11000 ...
## $ PAYMENTS                    : num  202 4103 622 0 678 ...
## $ MINIMUM_PAYMENTS            : num  140 1072 627 NA 245 ...
## $ PRC_FULL_PAYMENT             : num  0 0.222 0 0 0 ...
## $ TENURE                      : int  12 12 12 12 12 12 12 12 12 ...
```

Describe Data

TO BE DONE.

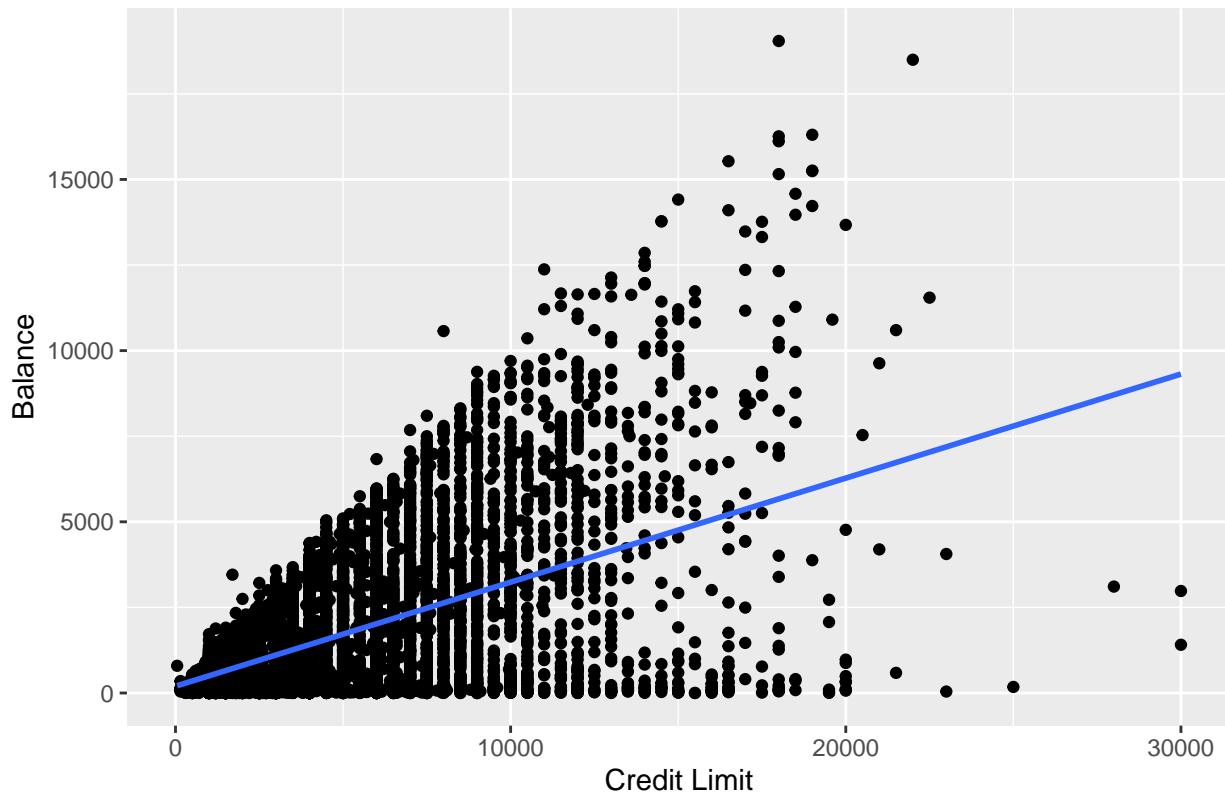
Explore Data

We'll plot some basic graphs, to ensure that the data conform to our limited domain understanding.

Balance vs. Credit Limit - we would expect to find a strong correlation here, even if only because those with small credit limits must have small balances. We certainly find it.

```
df %>% ggplot(aes(y=BALANCE, x=CREDIT_LIMIT)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) +
  labs(title="Balance vs. Credit Limit", x="Credit Limit", y="Balance")
```

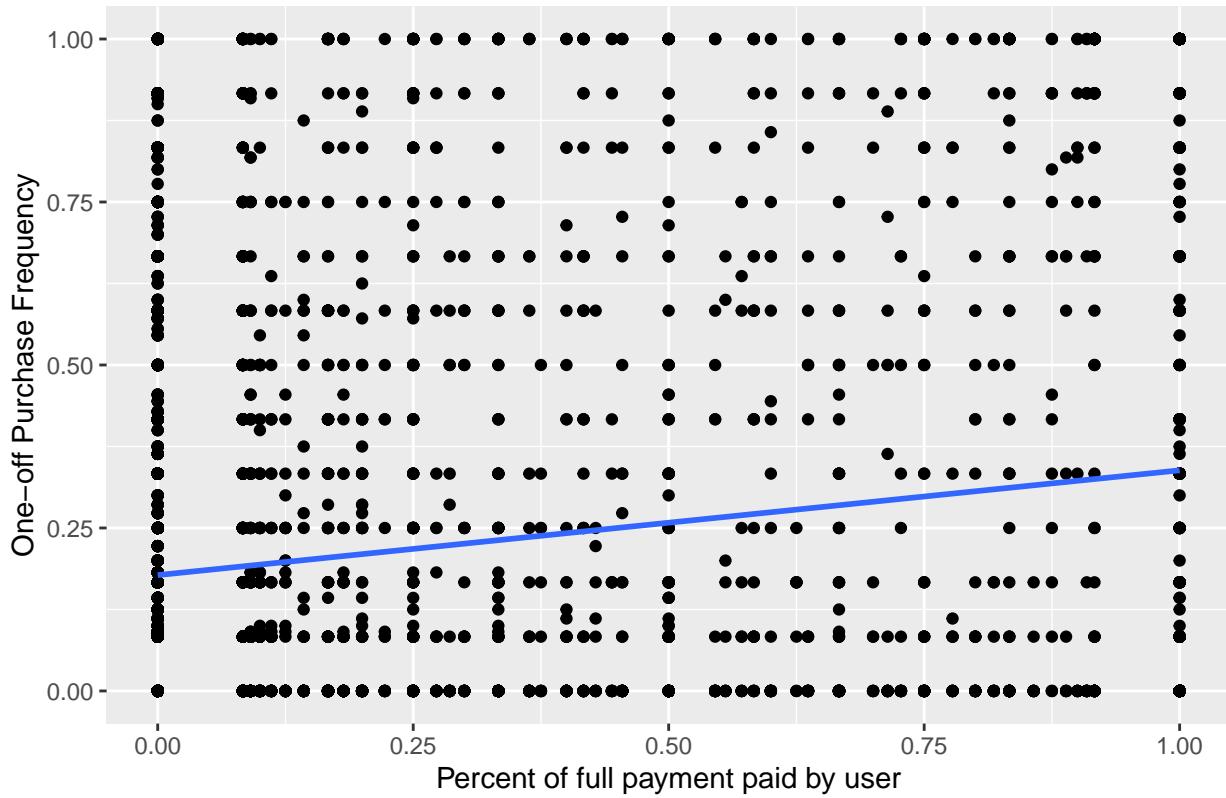
Balance vs. Credit Limit



Maybe those who pay higher portion of balance purchase more one-offs? No, this graph is pretty junky and doesn't tell us anything.

```
df %>% ggplot(aes(y=ONEOFF_PURCHASES_FREQUENCY, x=PRC_FULL_PAYMENT)) +  
  geom_point() +  
  geom_smooth(method = lm, se = FALSE) +  
  labs(title="One-off Purchase Frequency vs. Percent of full payment paid by user", x="Percent of full payment paid by user")
```

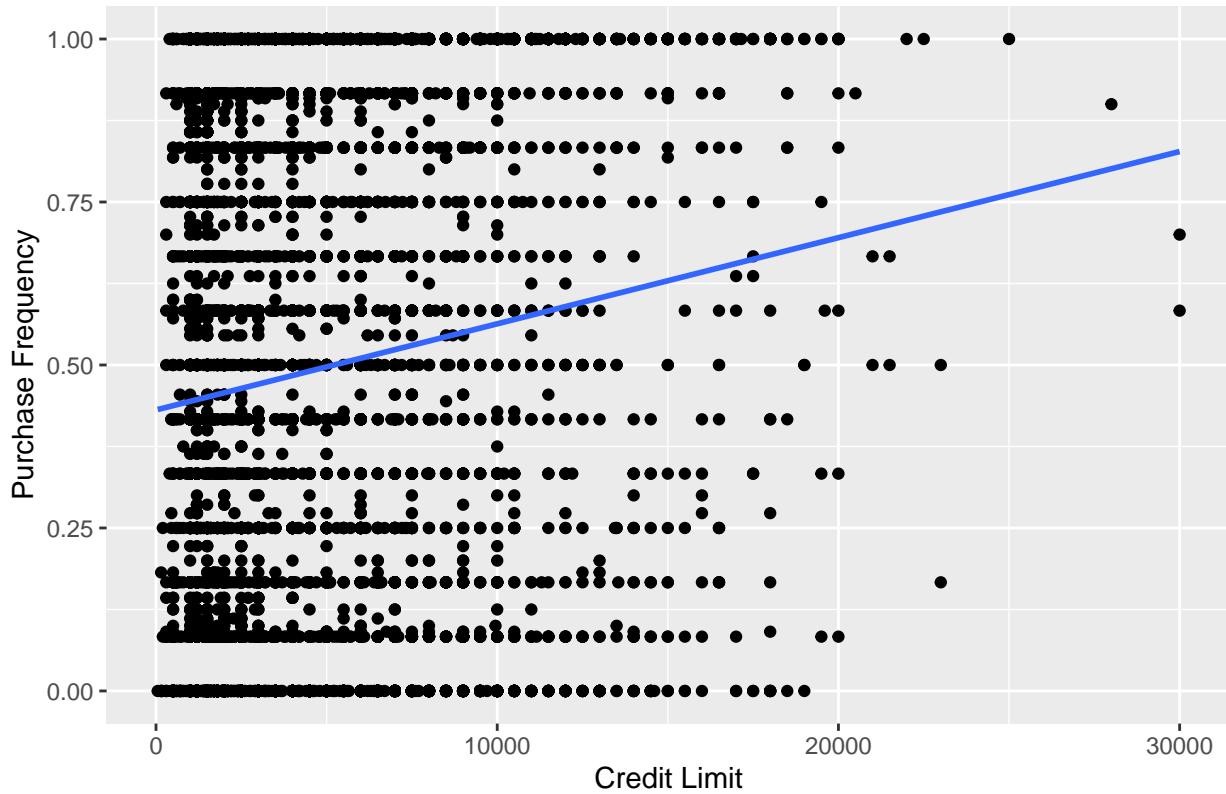
One-off Purchase Frequency vs. Percent of full payment paid by user



Here we can see that there appears to be a correlation between a customers credit limit, and the frequency of their purchases. Not surprising, given that frequent purchasing is one factor that leads to an increased credit limit.

```
df %>% ggplot(aes(y=PURCHASES_FREQUENCY, x=CREDIT_LIMIT)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) +
  labs(title="Purchase Frequency vs. Credit Limit", x="Credit Limit", y="Purchase Frequency")
```

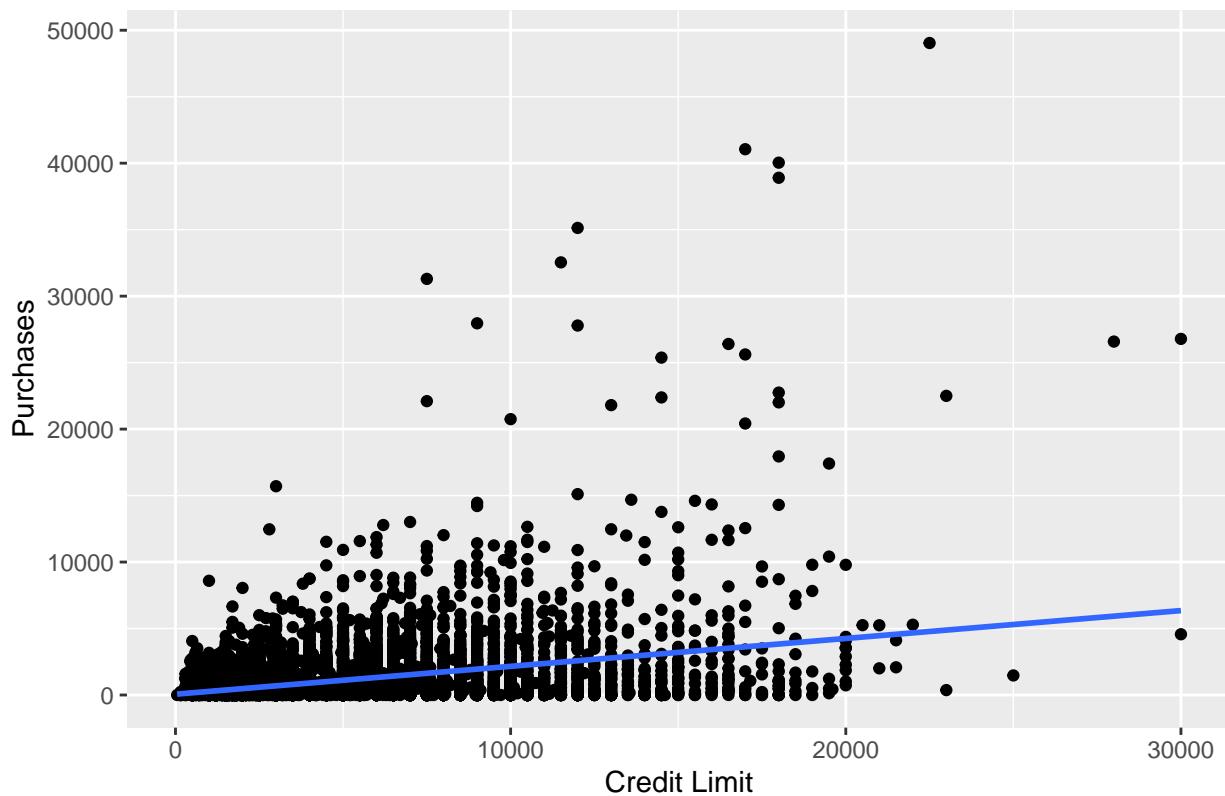
Purchase Frequency vs. Credit Limit



Higher credit limit, more purchases? Maybe. This graph is an excellent showcase of “outliers”. Only a small handful of the 9,000 records show a credit limit, or purchases, greater than \$20,000. This will be one of our guiding intuitions when we get to the Data Cleaning phase.

```
df %>% ggplot(aes(y=PURCHASES, x=CREDIT_LIMIT)) +  
  geom_point() +  
  geom_smooth(method = lm, se = FALSE) +  
  labs(title="Purchases vs. Credit Limit", x="Credit Limit", y="Purchases")
```

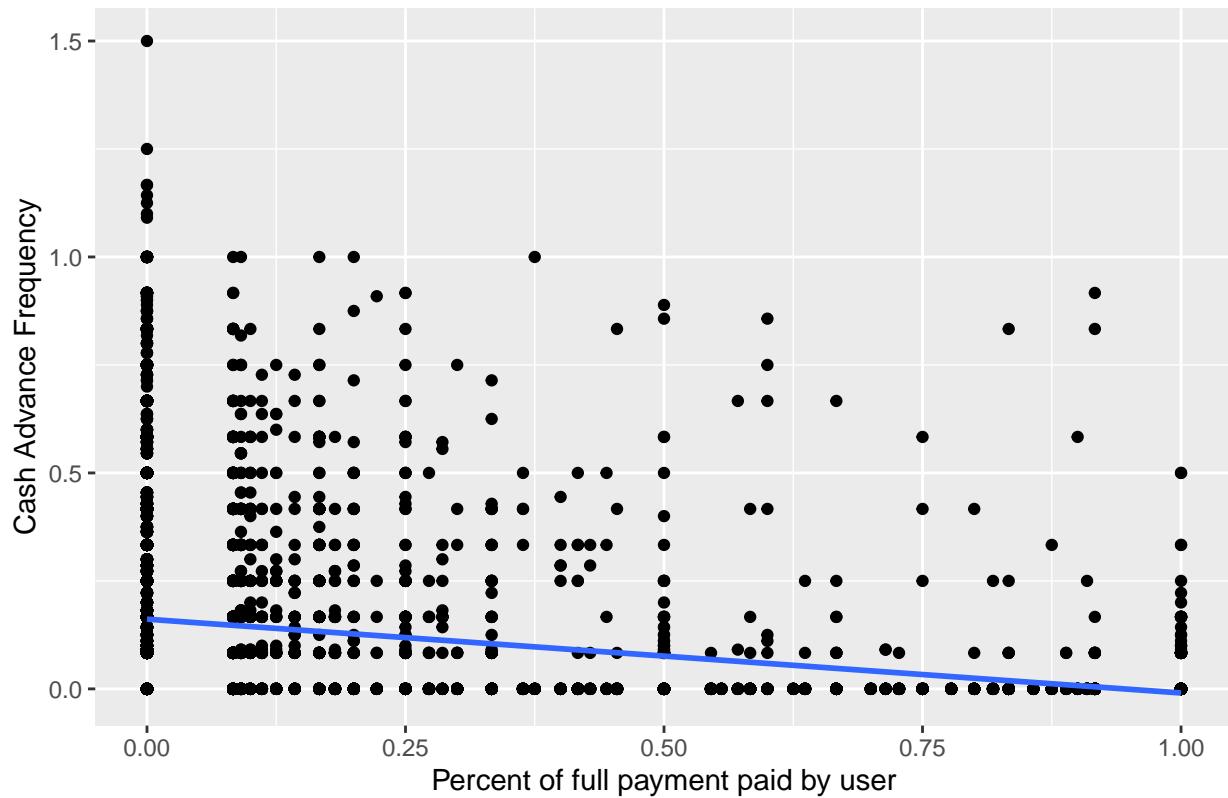
Purchases vs. Credit Limit



It would appear that those who pay off thier balance, do not take cash advances. No surprises, but a little hard to read because of the scale. One thing that appears in this graphs is a large “column” of people who are completely unlikely to pay off their entire balance, and take cash advances with very high frequency. (It is likely that the marketing department may have less interest in these people, than in others. But we'll see!)

```
df %>% ggplot(aes(y=CASH_ADVANCE_FREQUENCY, x=PRC_FULL_PAYMENT)) +  
  geom_point() +  
  geom_smooth(method = lm, se = FALSE) +  
  labs(title="Cash Advance Frequency vs. Percent of full payment paid by user", x="Percent of full payment paid by user")
```

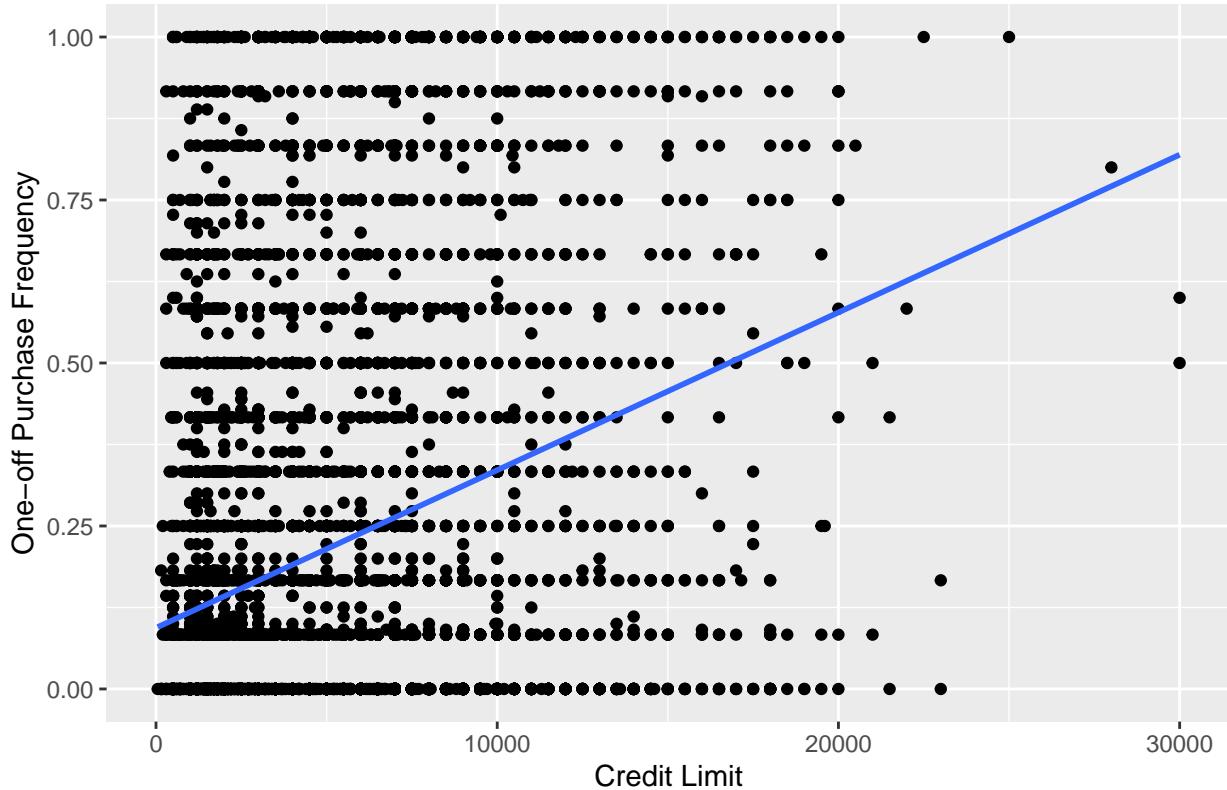
Cash Advance Frequency vs. Percent of full payment paid by user



One last graph, as part of our data exploration and sanity check. This one shows that there is a clear correlation between a customer's credit limit, and their frequency of one-off purchases. By now, we can be comfortable that our data isn't erratic, and that our limited domain knowledge isn't out to lunch.

```
df %>% ggplot(aes(y=ONEOFF_PURCHASES_FREQUENCY, x=CREDIT_LIMIT)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) +
  labs(title="One-off Purchase Frequency vs. Credit Limit", x="Credit Limit", y="One-off Purchase Frequency")
```

One-off Purchase Frequency vs. Credit Limit



Verify Data Quality

Check if data greater than zero, look for missing data

```
cf <- check_that(df, BALANCE > 0, BALANCE_FREQUENCY > 0, PURCHASES > 0, ONEOFF_PURCHASES > 0, INSTALLMENTS > 0)
summary(cf)
```

```
##      name items passes fails nNA error warning
## 1    V01  8950    8870     80     0 FALSE  FALSE
## 2    V02  8950    8870     80     0 FALSE  FALSE
## 3    V03  8950    6906   2044     0 FALSE  FALSE
## 4    V04  8950    4648   4302     0 FALSE  FALSE
## 5    V05  8950    5034   3916     0 FALSE  FALSE
## 6    V06  8950    4322   4628     0 FALSE  FALSE
## 7    V07  8950    6907   2043     0 FALSE  FALSE
## 8    V08  8950    4648   4302     0 FALSE  FALSE
## 9    V09  8950    5035   3915     0 FALSE  FALSE
## 10   V10  8950    4322   4628     0 FALSE  FALSE
## 11   V11  8950    6906   2044     0 FALSE  FALSE
## 12   V12  8950    8949     0     1 FALSE  FALSE
## 13   V13  8950    8637     0   313 FALSE  FALSE
## 14   V14  8950    3047   5903     0 FALSE  FALSE
## 15   V15  8950    8950     0     0 FALSE  FALSE
## 16   V16  8950    8722   227     1 FALSE  FALSE
##                               expression
## 1                         BALANCE > 0
```

```

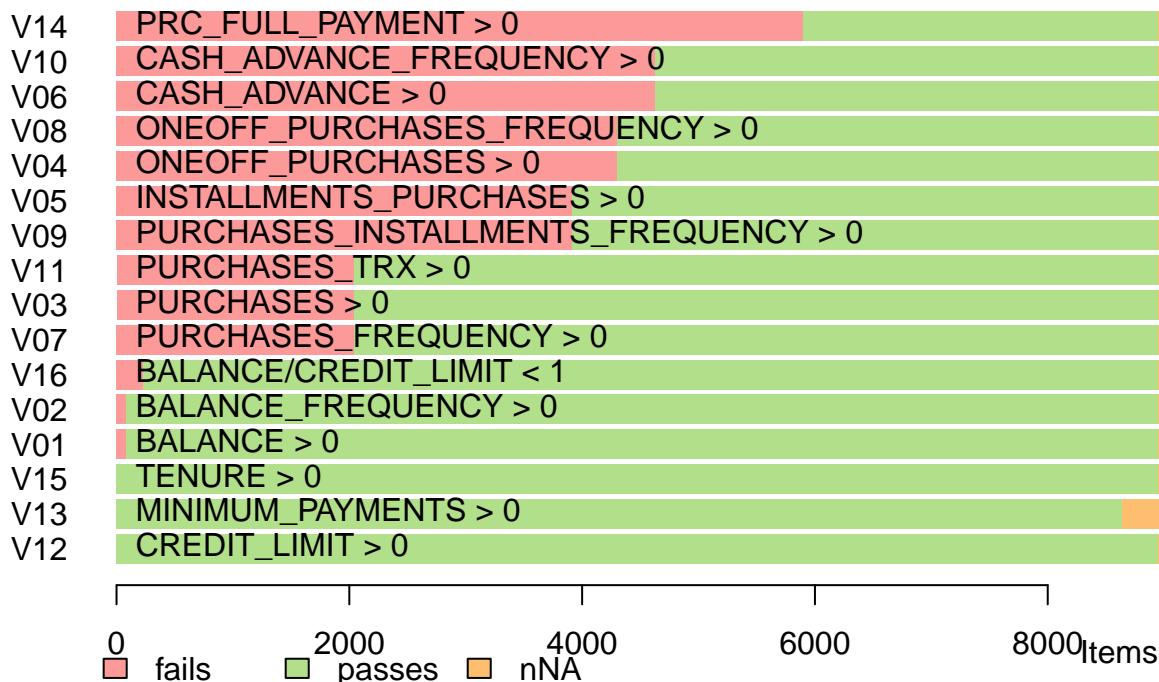
## 2           BALANCE_FREQUENCY > 0
## 3           PURCHASES > 0
## 4           ONEOFF_PURCHASES > 0
## 5           INSTALLMENTS_PURCHASES > 0
## 6           CASH_ADVANCE > 0
## 7           PURCHASES_FREQUENCY > 0
## 8           ONEOFF_PURCHASES_FREQUENCY > 0
## 9   PURCHASES_INSTALLMENTS_FREQUENCY > 0
## 10          CASH_ADVANCE_FREQUENCY > 0
## 11          PURCHASES_TRX > 0
## 12          CREDIT_LIMIT > 0
## 13          MINIMUM_PAYMENTS > 0
## 14          PRC_FULL_PAYMENT > 0
## 15          TENURE > 0
## 16          BALANCE/CREDIT_LIMIT < 1

```

There appear to be a great number of zeroes. Let's look at that a little more closely:

```
barplot(cf, main="Checks on the data set")
```

Checks on the data set



Lots of zero! Though, using our limited knowledge of this domain, we can understand that this could be perfectly reasonable. In the first few lines of the above chart, we can see that a lot of people never make the full payment, a lot of people never get a cash advance, and some people didn't even make a purchase. Seems entirely reasonable. We can note too that the “zeroes” in PURCHASES_TRX, PURCHASES, and PURCHASES_FREQUENCY appear to be identical. In line with what we'd expect!

Now let's include zero, and see how it all stacks up for Greater OR Equal to zero:

```
cf <- check_that(df, BALANCE >= 0, BALANCE_FREQUENCY >= 0, PURCHASES >= 0, ONEOFF_PURCHASES >= 0, INSTA
```

```
summary(cf)
```

```
##   name items passes fails nNA error warning
```

```

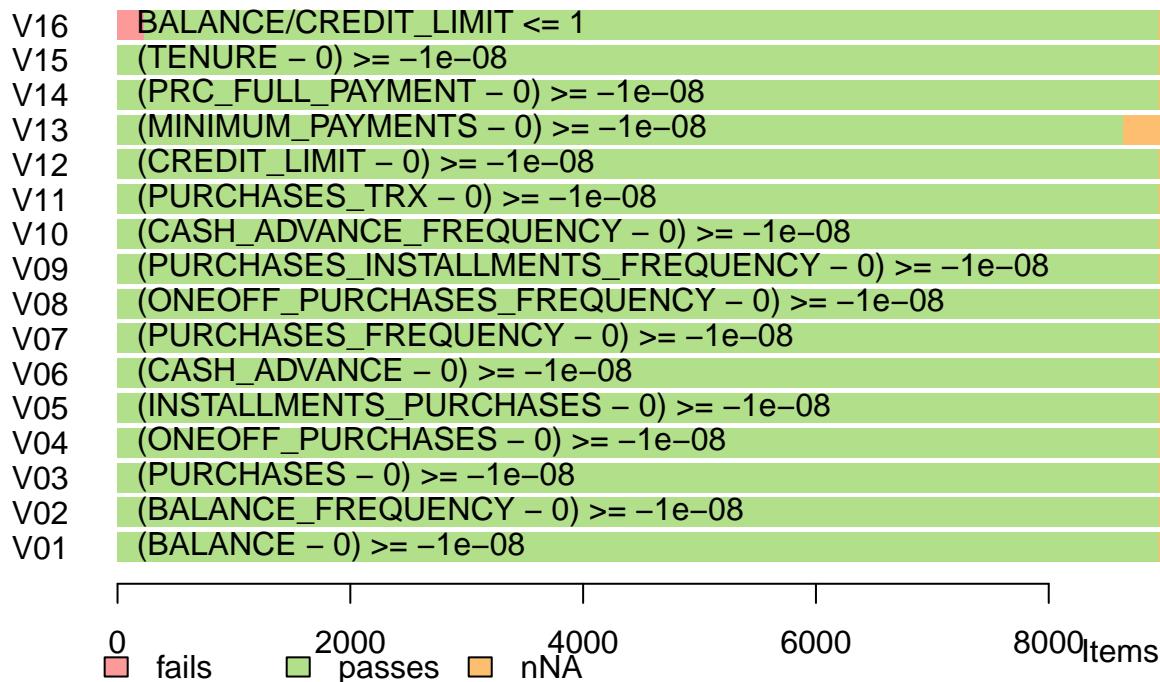
## 1 V01 8950 8950 0 0 FALSE FALSE
## 2 V02 8950 8950 0 0 FALSE FALSE
## 3 V03 8950 8950 0 0 FALSE FALSE
## 4 V04 8950 8950 0 0 FALSE FALSE
## 5 V05 8950 8950 0 0 FALSE FALSE
## 6 V06 8950 8950 0 0 FALSE FALSE
## 7 V07 8950 8950 0 0 FALSE FALSE
## 8 V08 8950 8950 0 0 FALSE FALSE
## 9 V09 8950 8950 0 0 FALSE FALSE
## 10 V10 8950 8950 0 0 FALSE FALSE
## 11 V11 8950 8950 0 0 FALSE FALSE
## 12 V12 8950 8949 0 1 FALSE FALSE
## 13 V13 8950 8637 0 313 FALSE FALSE
## 14 V14 8950 8950 0 0 FALSE FALSE
## 15 V15 8950 8950 0 0 FALSE FALSE
## 16 V16 8950 8722 227 1 FALSE FALSE
##                                         expression
## 1                               (BALANCE - 0) >= -1e-08
## 2                               (BALANCE_FREQUENCY - 0) >= -1e-08
## 3                               (PURCHASES - 0) >= -1e-08
## 4                               (ONEOFF_PURCHASES - 0) >= -1e-08
## 5                               (INSTALLMENTS_PURCHASES - 0) >= -1e-08
## 6                               (CASH_ADVANCE - 0) >= -1e-08
## 7                               (PURCHASES_FREQUENCY - 0) >= -1e-08
## 8                               (ONEOFF_PURCHASES_FREQUENCY - 0) >= -1e-08
## 9 (PURCHASES_INSTALLMENTS_FREQUENCY - 0) >= -1e-08
## 10 (CASH_ADVANCE_FREQUENCY - 0) >= -1e-08
## 11 (PURCHASES_TRX - 0) >= -1e-08
## 12 (CREDIT_LIMIT - 0) >= -1e-08
## 13 (MINIMUM_PAYMENTS - 0) >= -1e-08
## 14 (PRC_FULL_PAYMENT - 0) >= -1e-08
## 15 (TENURE - 0) >= -1e-08
## 16 BALANCE/CREDIT_LIMIT <= 1

```

Our only serious problems remaining are some N/As in Minimum Payment. We also see that a few people have a balance greater than their credit limit - this would not constitute a problem with the data. Simply a problem the customer in question is having. (Their balance is higher than their credit limit)

```
barplot(cf,main="Checks on the data set")
```

Checks on the data set



As a last look at the data, let's look at the variance of the individual columns.

```

var(df$BALANCE)

## [1] 4332775
var(df$BALANCE_FREQUENCY)

## [1] 0.05612351
var(df$PURCHASES)

## [1] 4565208
var(df$ONEOFF_PURCHASES)

## [1] 2755228
var(df$INSTALLMENTS_PURCHASES)

## [1] 817827.4
var(df$CASH_ADVANCE)

## [1] 4398096
var(df$PURCHASES_FREQUENCY)

## [1] 0.1610985
var(df$ONEOFF_PURCHASES_FREQUENCY)

## [1] 0.08900441
var(df$PURCHASES_INSTALLMENTS_FREQUENCY)

```

```

## [1] 0.1579647
var(df$CASH_ADVANCE_FREQUENCY)

## [1] 0.04004857
var(df$CASH_ADVANCE_TRX)

## [1] 46.5758
var(df$PURCHASES_TRX)

## [1] 617.9027
var(df$CREDIT_LIMIT)

## [1] NA
var(df$PAYMENTS)

## [1] 8381394
var(df$MINIMUM_PAYMENTS)

## [1] NA
var(df$PRC_FULL_PAYMENT)

## [1] 0.08555578
var(df$TENURE)

## [1] 1.791129

```

As expected, there is a big difference in the variance, between the different types of data. Dollar amounts are huge, and have huge variance. Transaction frequencies are small, and have small variance. As our ideal investigation into this dataset would involve clustering and predicting based on both types of values, we anticipate scaling some or all of the data columns we will be using, so that our algorithm can give comparable consideration to all of them.

DATA PREPARATION

Select Data

We saw that MINIMUM_PAYMENTS is riddled with N/As, so let's remove that column. This should be compatible with our Business Objectives, as we're hoping to get people to spend, rather than optimize their repayments.

```
df$MINIMUM_PAYMENTS <- NULL
```

And let's check that worked:

```
cf <- check_that(df, BALANCE >= 0, BALANCE_FREQUENCY >= 0, PURCHASES >= 0, ONEOFF_PURCHASES >= 0, INSTA
```

```

##      name items passes fails nNA error warning
## 1    V01   8950     8950      0    0 FALSE   FALSE
## 2    V02   8950     8950      0    0 FALSE   FALSE
## 3    V03   8950     8950      0    0 FALSE   FALSE
## 4    V04   8950     8950      0    0 FALSE   FALSE
## 5    V05   8950     8950      0    0 FALSE   FALSE

```

```

## 6   V06  8950  8950    0    0 FALSE  FALSE
## 7   V07  8950  8950    0    0 FALSE  FALSE
## 8   V08  8950  8950    0    0 FALSE  FALSE
## 9   V09  8950  8950    0    0 FALSE  FALSE
## 10  V10  8950  8950    0    0 FALSE  FALSE
## 11  V11  8950  8950    0    0 FALSE  FALSE
## 12  V12  8950  8949    0    1 FALSE  FALSE
## 13  V13    0     0    0    0 TRUE   FALSE
## 14  V14  8950  8950    0    0 FALSE  FALSE
## 15  V15  8950  8950    0    0 FALSE  FALSE
## 16  V16  8950  8722   227   1 FALSE  FALSE
##                                         expression
## 1                               (BALANCE - 0) >= -1e-08
## 2                               (BALANCE_FREQUENCY - 0) >= -1e-08
## 3                               (PURCHASES - 0) >= -1e-08
## 4                               (ONEOFF_PURCHASES - 0) >= -1e-08
## 5                               (INSTALLMENTS_PURCHASES - 0) >= -1e-08
## 6                               (CASH_ADVANCE - 0) >= -1e-08
## 7                               (PURCHASES_FREQUENCY - 0) >= -1e-08
## 8                               (ONEOFF_PURCHASES_FREQUENCY - 0) >= -1e-08
## 9   (PURCHASES_INSTALLMENTS_FREQUENCY - 0) >= -1e-08
## 10  (CASH_ADVANCE_FREQUENCY - 0) >= -1e-08
## 11  (PURCHASES_TRX - 0) >= -1e-08
## 12  (CREDIT_LIMIT - 0) >= -1e-08
## 13  (MINIMUM_PAYMENTS - 0) >= -1e-08
## 14  (PRC_FULL_PAYMENT - 0) >= -1e-08
## 15  (TENURE - 0) >= -1e-08
## 16  BALANCE/CREDIT_LIMIT <= 1

```

There were some N/As in CREDIT_LIMIT, but we'd really like to use that column in our clusterings and predictions. let's remove those rows so we can use that column.

```
df <- na.omit(df)
```

Let's check that we haven't misunderstood that N/A, or the code we used to remove it, and accidentally ransacked the data we're using with that move - looks good.

```
cf <- check_that(df, BALANCE >= 0, BALANCE_FREQUENCY >= 0, PURCHASES >= 0, ONEOFF_PURCHASES >= 0, INSTA
summary(cf)
```

```

##      name items passes fails nNA error warning
## 1   V01  8949   8949    0    0 FALSE  FALSE
## 2   V02  8949   8949    0    0 FALSE  FALSE
## 3   V03  8949   8949    0    0 FALSE  FALSE
## 4   V04  8949   8949    0    0 FALSE  FALSE
## 5   V05  8949   8949    0    0 FALSE  FALSE
## 6   V06  8949   8949    0    0 FALSE  FALSE
## 7   V07  8949   8949    0    0 FALSE  FALSE
## 8   V08  8949   8949    0    0 FALSE  FALSE
## 9   V09  8949   8949    0    0 FALSE  FALSE
## 10  V10  8949   8949    0    0 FALSE  FALSE
## 11  V11  8949   8949    0    0 FALSE  FALSE
## 12  V12  8949   8949    0    0 FALSE  FALSE
## 13  V13    0     0    0    0 TRUE   FALSE
## 14  V14  8949   8949    0    0 FALSE  FALSE
## 15  V15  8949   8949    0    0 FALSE  FALSE

```

```

## 16 V16 8949 8722 227 0 FALSE FALSE
##                                expression
## 1          (BALANCE - 0) >= -1e-08
## 2          (BALANCE_FREQUENCY - 0) >= -1e-08
## 3          (PURCHASES - 0) >= -1e-08
## 4          (ONEOFF_PURCHASES - 0) >= -1e-08
## 5          (INSTALLMENTS_PURCHASES - 0) >= -1e-08
## 6          (CASH_ADVANCE - 0) >= -1e-08
## 7          (PURCHASES_FREQUENCY - 0) >= -1e-08
## 8          (ONEOFF_PURCHASES_FREQUENCY - 0) >= -1e-08
## 9          (PURCHASES_INSTALLMENTS_FREQUENCY - 0) >= -1e-08
## 10         (CASH_ADVANCE_FREQUENCY - 0) >= -1e-08
## 11         (PURCHASES_TRX - 0) >= -1e-08
## 12         (CREDIT_LIMIT - 0) >= -1e-08
## 13         (MINIMUM_PAYMENTS - 0) >= -1e-08
## 14         (PRC_FULL_PAYMENT - 0) >= -1e-08
## 15         (TENURE - 0) >= -1e-08
## 16         BALANCE/CREDIT_LIMIT <= 1

```

Clean Data

During the Data Exploration phase, we discovered that of the 9,000 rows of data, there are a handful that have much higher values than all the rest. While we are certainly interested in these customers, they can be easily found without expensive machine learning. Simple code can be used to scrape off the customers with a credit limit, or purchases, higher than some amount.

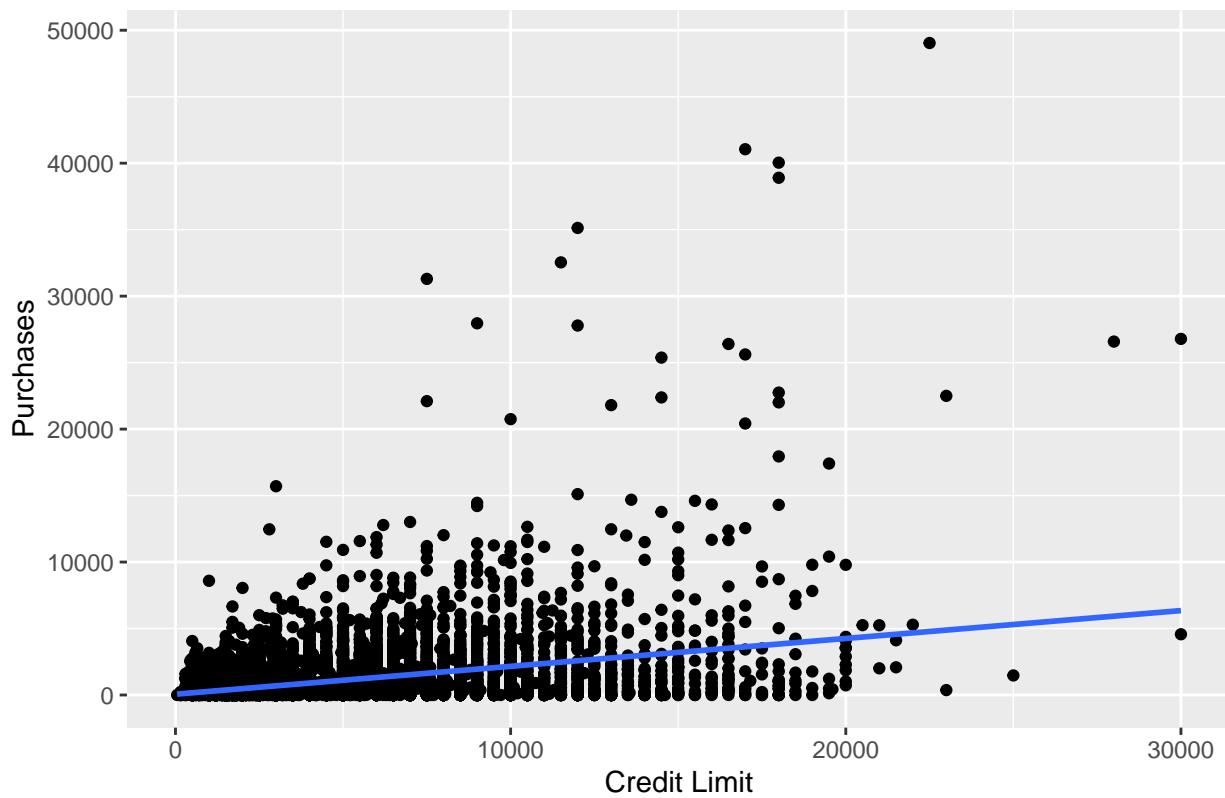
As we're hoping to come up with customer segments that are a significant portion of our customer base, perhaps containing hundreds or thousands of customers, let's look at "capping" some of these extreme values, to see if it moderates these few extreme outliers and suggests more patterns at a smaller scale.

```
df_capped <- df
```

First, here's the graph that drew this to our attention:

```
df %>% ggplot(aes(y=PURCHASES, x=CREDIT_LIMIT)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) +
  labs(title="Purchases vs. Credit Limit", x="Credit Limit", y="Purchases")
```

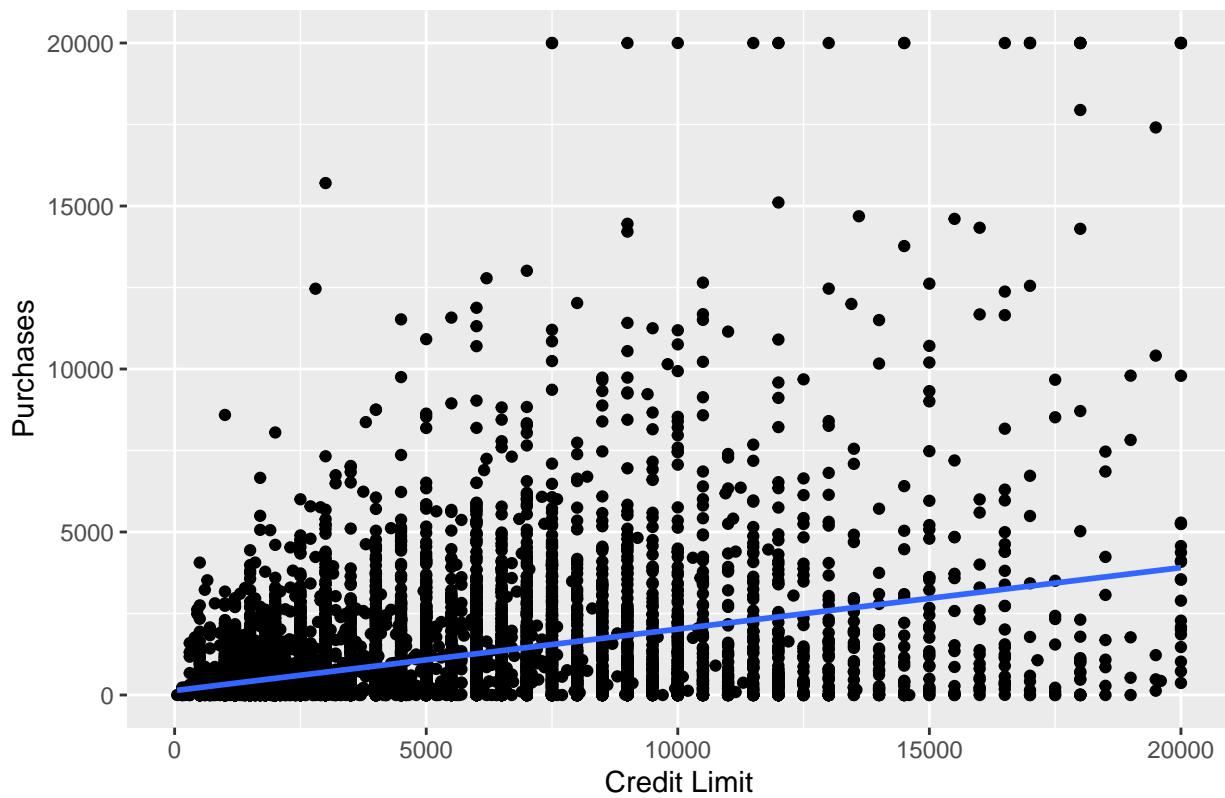
Purchases vs. Credit Limit



Let's cap those two variables at \$20,000 and see if it seems a little less crazy:

```
df_capped$CREDIT_LIMIT[df$CREDIT_LIMIT > 20000 ] <- 20000  
df_capped$PURCHASES [df$PURCHASES > 20000 ] <- 20000  
df_capped %>% ggplot(aes(y=PURCHASES, x=CREDIT_LIMIT)) +  
  geom_point() +  
  geom_smooth(method = lm, se = FALSE) +  
  labs(title="Purchases vs. Credit Limit", x="Credit Limit", y="Purchases")
```

Purchases vs. Credit Limit

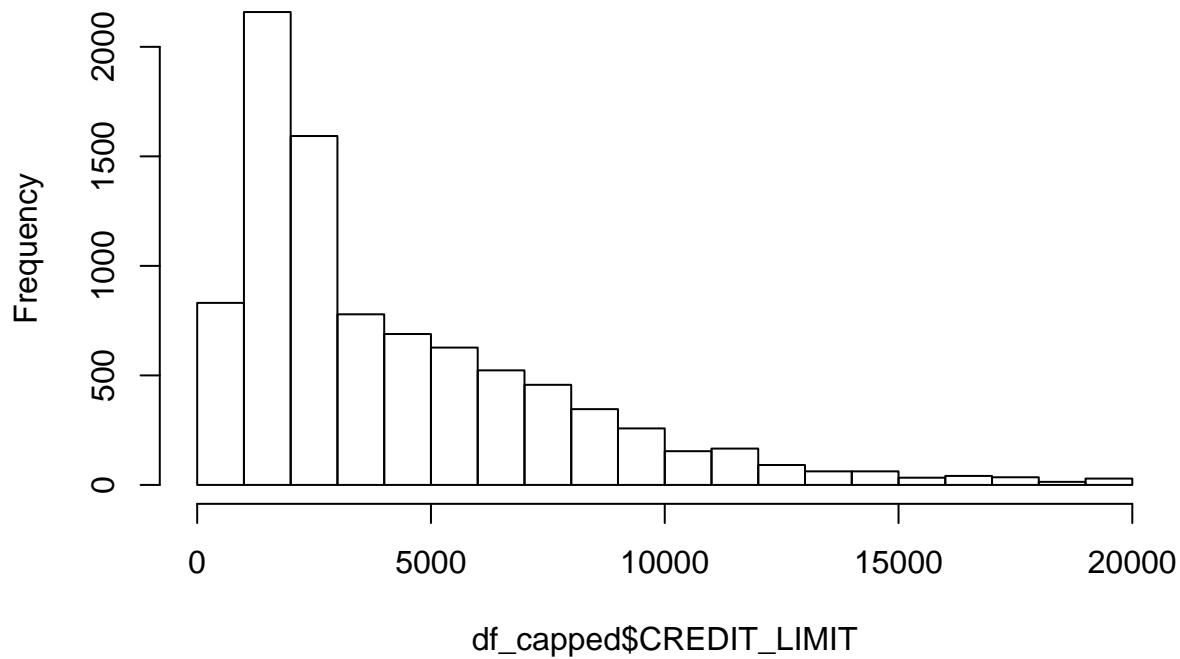


That seems much better, both for the purchases, and the credit limit.

Let's look at one variable at a time, to ensure we haven't created a big distortion out at the high end:

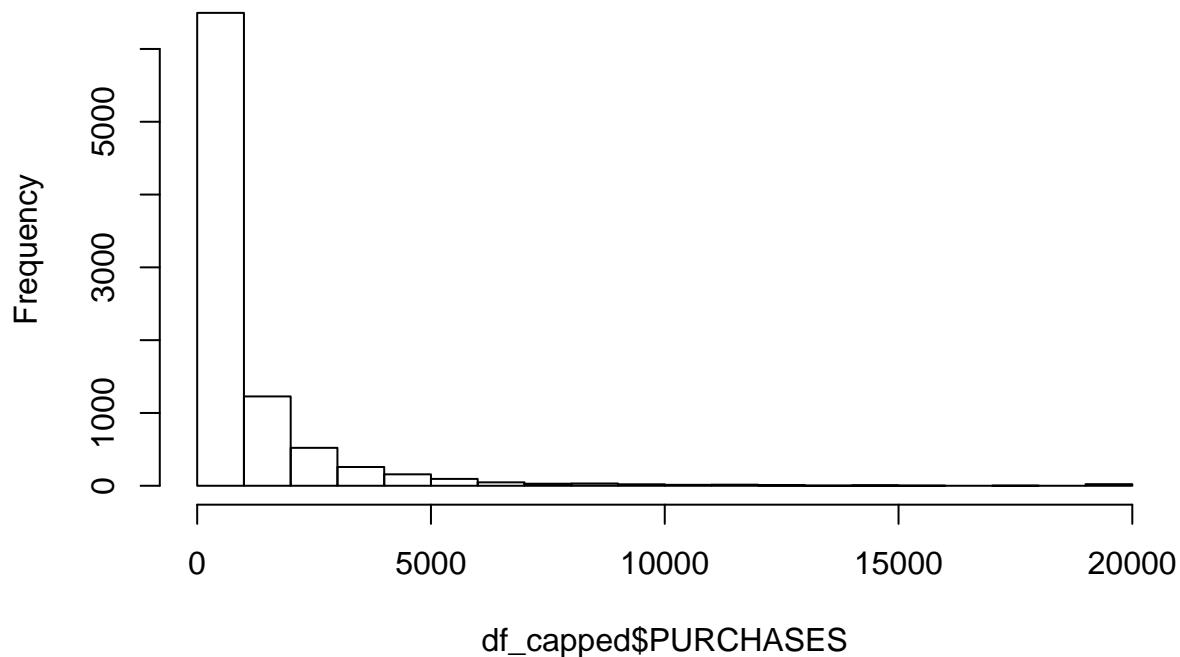
```
hist(df_capped$CREDIT_LIMIT)
```

Histogram of df_capped\$CREDIT_LIMIT



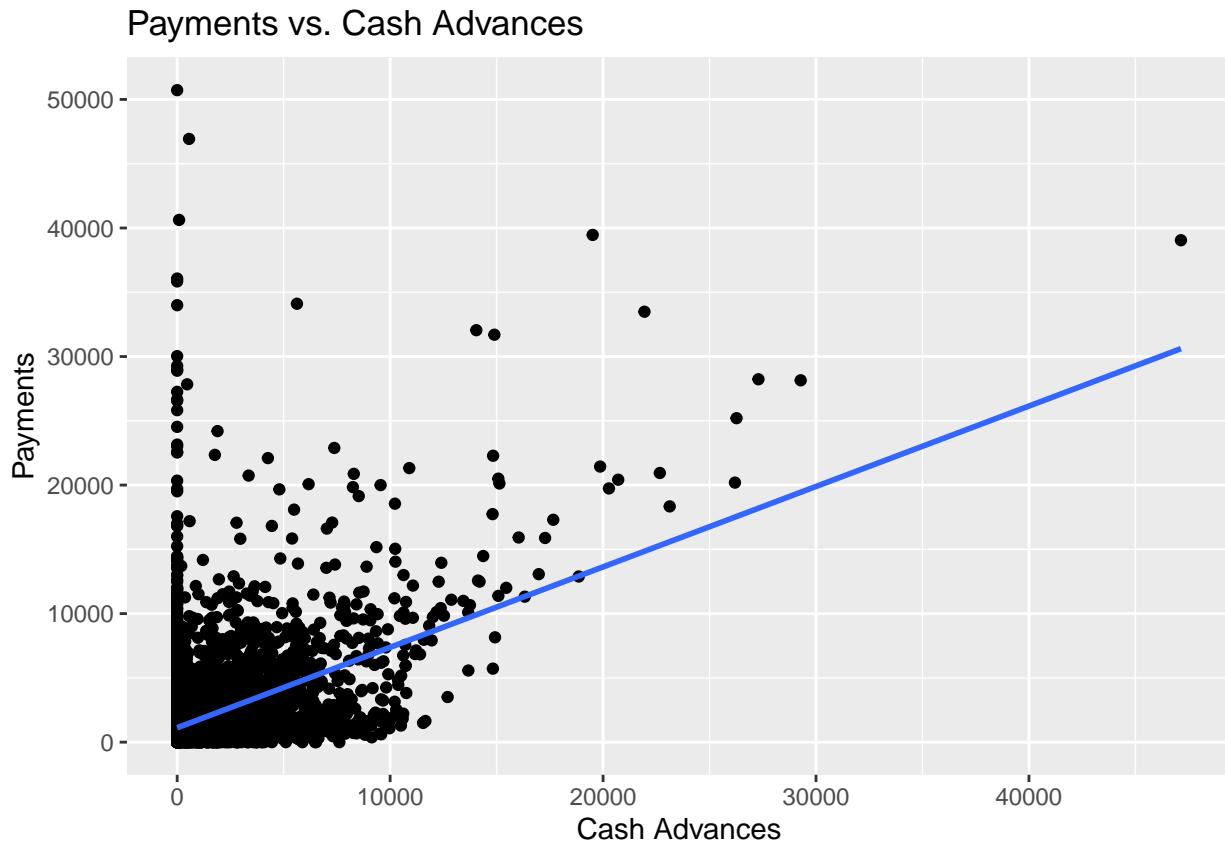
```
hist(df_capped$PURCHASES)
```

Histogram of df_capped\$PURCHASES



Seems okay. The distortions appear small. Let's look at the distribution of some of the other big-value columns that we're most interested in, to see how they look.

```
df %>% ggplot(aes(y=PAYMENTS, x=CASH_ADVANCE)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) +
  labs(title="Payments vs. Cash Advances", x="Cash Advances", y="Payments")
```



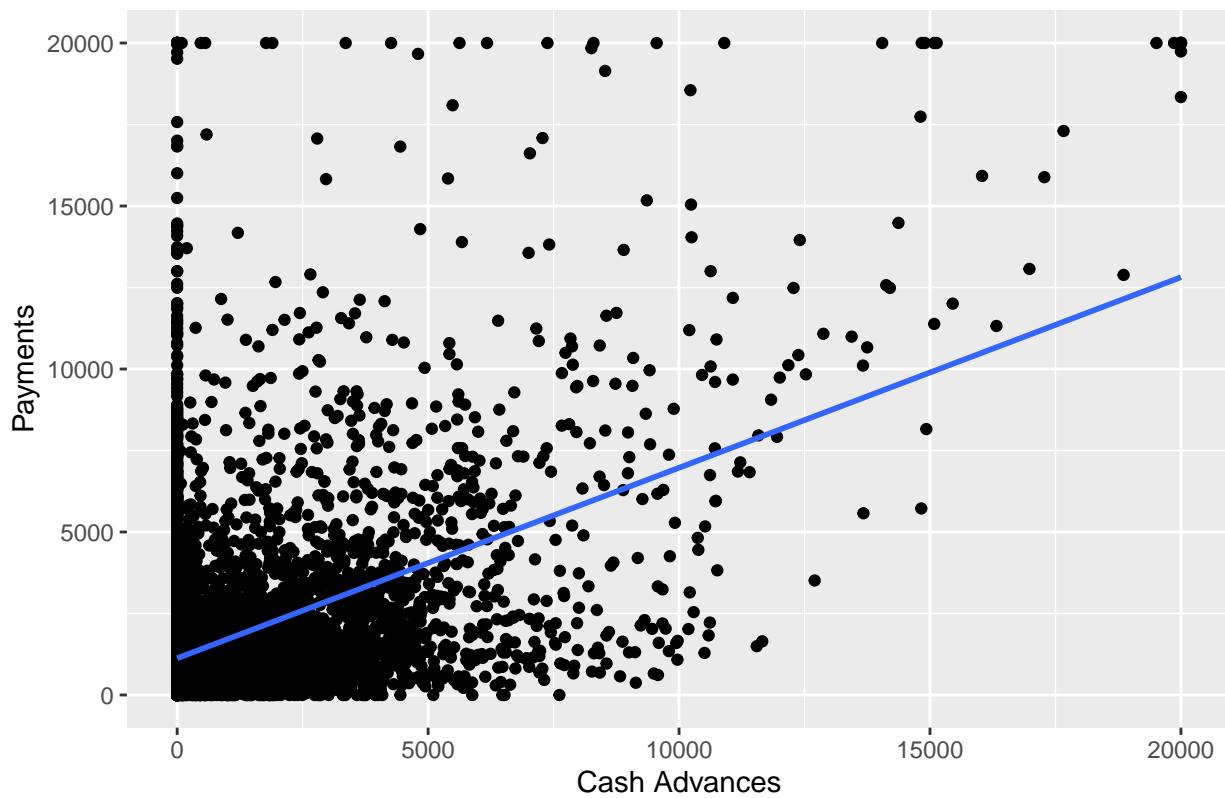
Wow, that's some serious outliers. Let's cap them.

We'll continue to use 20,000, to keep things simple, unless something suggests either that we're affecting large numbers of data points, or not budging even a handful.

```
df_capped$PAYMENTS[df$PAYMENTS > 20000] <- 20000
df_capped$CASH_ADVANCE[df$CASH_ADVANCE > 20000] <- 20000

df_capped %>% ggplot(aes(y=PAYMENTS, x=CASH_ADVANCE)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) +
  labs(title="Payments vs. Cash Advances", x="Cash Advances", y="Payments")
```

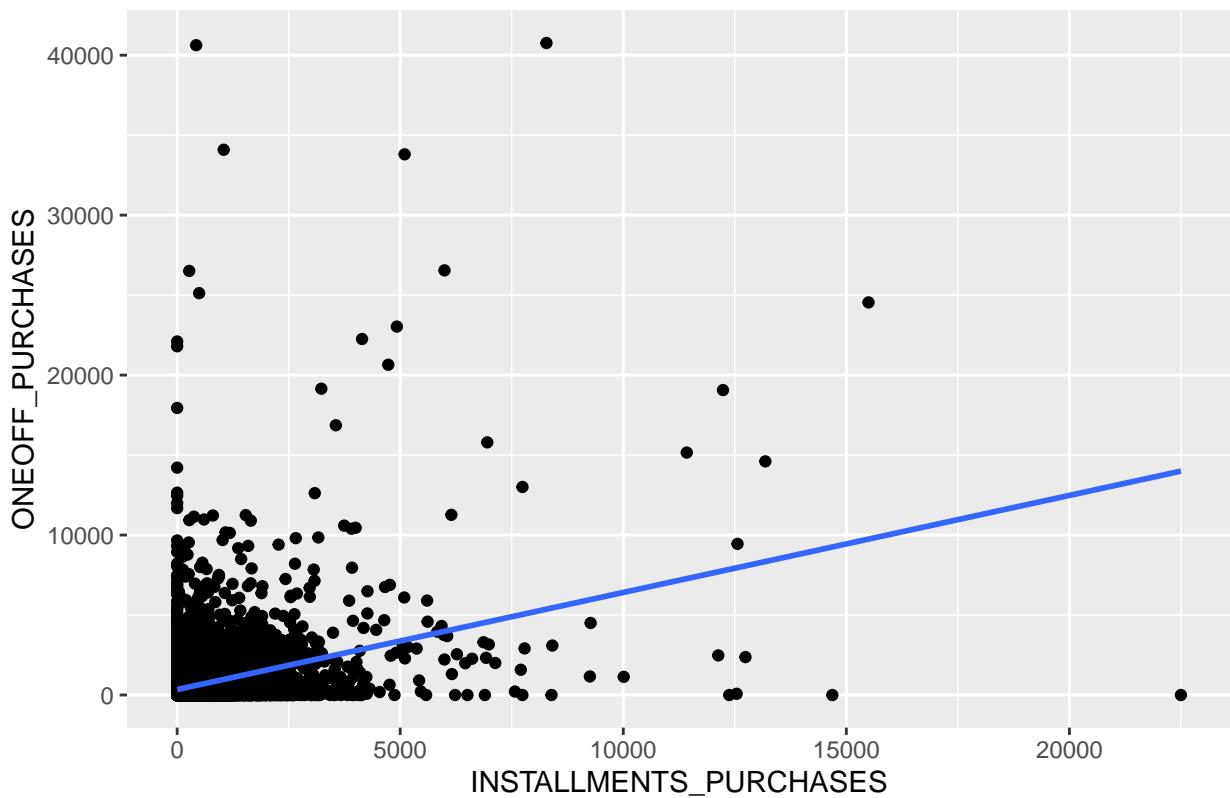
Payments vs. Cash Advances



Two more columns we're very interested in using need to be check for this:

```
df %>% ggplot(aes(y=ONEOFF_PURCHASES, x=INSTALLMENTS_PURCHASES)) +  
  geom_point() +  
  geom_smooth(method = lm, se = FALSE) +  
  labs(title="ONEOFF_PURCHASES vs. INSTALLMENTS_PURCHASES", x="INSTALLMENTS_PURCHASES", y="ONEOFF_PURCHASES")
```

ONEOFF_PURCHASES vs. INSTALLMENTS_PURCHASES

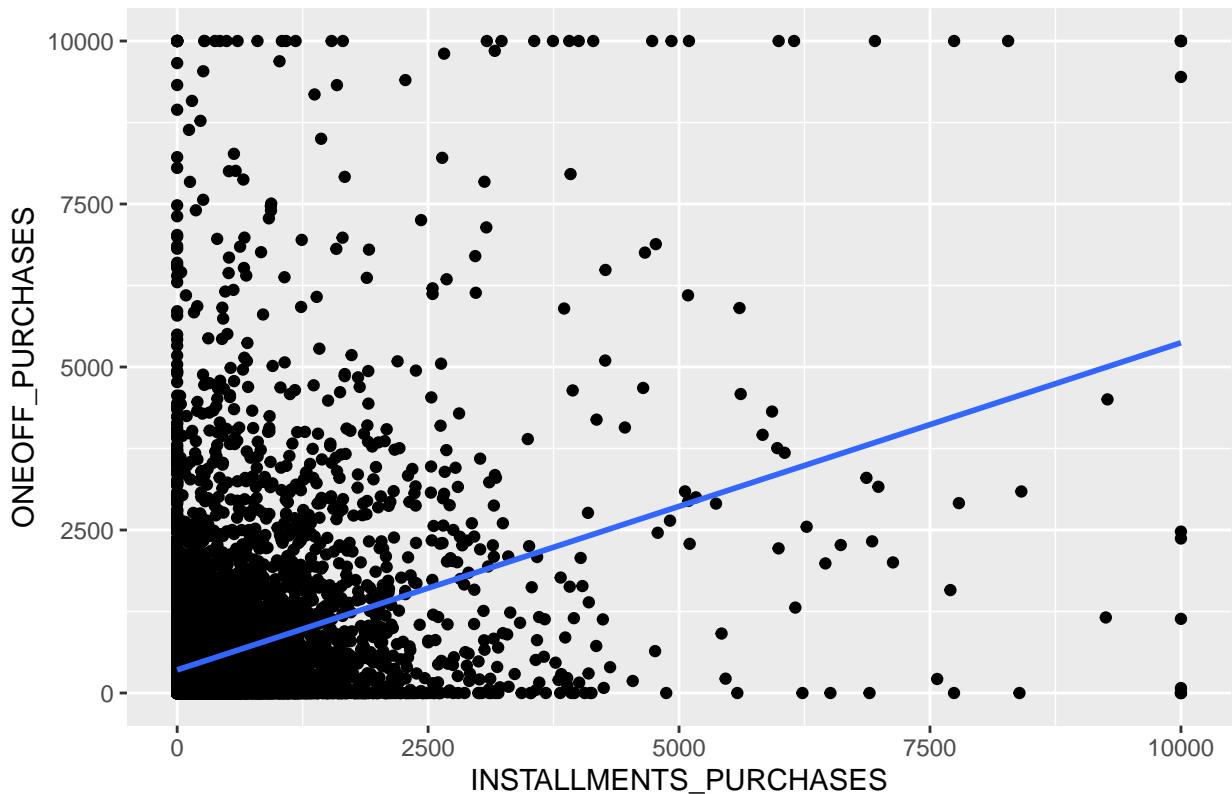


A very few extreme outliers. In this case, 20,000 will hardly touch these - let's cap these two at 10,000.

```
df_capped$ONEOFF_PURCHASES[df$ONEOFF_PURCHASES > 10000] <- 10000
df_capped$INSTALLMENTS_PURCHASES[df$INSTALLMENTS_PURCHASES > 10000] <- 10000
```

```
df_capped %>% ggplot(aes(y=ONEOFF_PURCHASES, x=INSTALLMENTS_PURCHASES)) +
  geom_point() +
  geom_smooth(method = lm, se = FALSE) +
  labs(title="ONEOFF_PURCHASES vs. INSTALLMENTS_PURCHASES", x="INSTALLMENTS_PURCHASES", y="ONEOFF_PURCHASES")
```

ONEOFF_PURCHASES vs. INSTALLMENTS_PURCHASES



Let's look at the range of the parameters, to make sure they've wound up how we were expecting:

```
range(df_capped$BALANCE)
## [1] 0.00 19043.14

range(df_capped$BALANCE_FREQUENCY)
## [1] 0 1

range(df_capped$PURCHASES)
## [1] 0 20000

range(df_capped$ONEOFF_PURCHASES)
## [1] 0 10000

range(df_capped$INSTALLMENTS_PURCHASES)
## [1] 0 10000

range(df_capped$CASH_ADVANCE)
## [1] 0 20000

range(df_capped$PURCHASES_FREQUENCY)
## [1] 0 1

range(df_capped$ONEOFF_PURCHASES_FREQUENCY)
## [1] 0 1
```

```

range(df_capped$PURCHASES_INSTALLMENTS_FREQUENCY)

## [1] 0 1

range(df_capped$CASH_ADVANCE_FREQUENCY)

## [1] 0.0 1.5

range(df_capped$CASH_ADVANCE_TRX)

## [1] 0 123

range(df_capped$PURCHASES_TRX)

## [1] 0 358

range(df_capped$CREDIT_LIMIT)

## [1] 50 20000

range(df_capped$PAYMENTS)

## [1] 0 20000

#range(df$MINIMUM_PAYMENTS)      - doesn't work as "null".
range(df_capped$PRC_FULL_PAYMENT)

## [1] 0 1

range(df_capped$TENURE)

## [1] 6 12

```

Construct Data

We wish to try three different approaches to the inputs for our model.

- 1) A simple model, using only dollar values, capped to reduce the impact of outliers. (df_capped)
- 2) A preliminary “hack” of scaling the data, for easy interpretation of our scaled values, that will allow training and prediction based on both dollar amounts and frequencies of transactions.
- 3) More rigourously scaled data, which we may or may not use in our Shiny app due to time constraints and the thrill of mapping our sample input values into the distribution of the scaled data.

Let's do the basic, hackish scaling:

```

hackish_scaled_df <- df_capped
hackish_scaled_df$BALANCE <- df_capped$BALANCE/20000
hackish_scaled_df$PURCHASES <- df_capped$PURCHASES/20000
hackish_scaled_df$ONEOFF_PURCHASES <- df_capped$ONEOFF_PURCHASES/10000
hackish_scaled_df$INSTALLMENTS_PURCHASES <- df_capped$INSTALLMENTS_PURCHASES/10000
hackish_scaled_df$CASH_ADVANCE <- df_capped$CASH_ADVANCE/20000
hackish_scaled_df$CREDIT_LIMIT <- df_capped$CREDIT_LIMIT/20000

range(hackish_scaled_df$BALANCE)

## [1] 0.0000000 0.9521569

range(hackish_scaled_df$PURCHASES)

## [1] 0 1

```

```

range(hackish_scaled_df$ONEOFF_PURCHASES)

## [1] 0 1

range(hackish_scaled_df$INSTALLMENTS_PURCHASES)

## [1] 0 1

range(hackish_scaled_df$CASH_ADVANCE)

## [1] 0 1

range(hackish_scaled_df$CREDIT_LIMIT)

## [1] 0.0025 1.0000

# And the frequency variables, which are pre-scaled, for comparison:
range(hackish_scaled_df$PURCHASES_FREQUENCY)

## [1] 0 1

range(hackish_scaled_df$ONEOFF_PURCHASES_FREQUENCY)

## [1] 0 1

range(hackish_scaled_df$PURCHASES_INSTALLMENTS_FREQUENCY)

## [1] 0 1

range(hackish_scaled_df$CASH_ADVANCE_FREQUENCY)

## [1] 0.0 1.5

```

Now we will do a formal scaling. This can be compared to the hackish one above, both in terms of accuracy of clustering and prediction, and the ease with which we can map new, arbitrary samples into its distribution for future predictions.

```

scaled_df <- df_capped

scaled_df$BALANCE <- scale(scaled_df$BALANCE) [, 1]
scaled_df$BALANCE_FREQUENCY <- scale(df$BALANCE_FREQUENCY) [, 1]
scaled_df$PURCHASES <- scale(df$PURCHASES) [, 1]
scaled_df$ONEOFF_PURCHASES <- scale(df$ONEOFF_PURCHASES) [, 1]
scaled_df$INSTALLMENTS_PURCHASES<- scale(df$INSTALLMENTS_PURCHASES) [, 1]
scaled_df$CASH_ADVANCE<- scale(df$CASH_ADVANCE) [, 1]
scaled_df$PURCHASES_FREQUENCY<- scale(df$PURCHASES_FREQUENCY) [, 1]
scaled_df$ONEOFF_PURCHASES_FREQUENCY <- scale(df$ONEOFF_PURCHASES_FREQUENCY) [, 1]
scaled_df$PURCHASES_INSTALLMENTS_FREQUENCY<- scale(df$PURCHASES_INSTALLMENTS_FREQUENCY) [, 1]
scaled_df$CASH_ADVANCE_FREQUENCY <- scale(df$CASH_ADVANCE_FREQUENCY) [, 1]
scaled_df$cash_advance_trx <- scale(df$cash_advance_trx) [, 1]
scaled_df$PURCHASES_TRX <- scale(df$PURCHASES_TRX) [, 1]
scaled_df$CREDIT_LIMIT<- scale(df$CREDIT_LIMIT) [, 1]
scaled_df$PAYMENTS <- scale(df$PAYMENTS) [, 1]
# scaled_df$MINIMUM_PAYMENTS <- scale(df$MINIMUM_PAYMENTS) [, 1]      # We've nullified this column.
scaled_df$PRC_FULL_PAYMENT <- scale(df$PRC_FULL_PAYMENT) [, 1]
scaled_df$TENURE<- scale(df$TENURE) [, 1]
str(scaled_df)

## 'data.frame': 8949 obs. of 17 variables:
##   $ CUST_ID                      : chr "C10001" "C10002" "C10003" "C10004" ...

```

```

## $ BALANCE : num -0.732 0.787 0.447 0.049 -0.359 ...
## $ BALANCE_FREQUENCY : num -0.25 0.134 0.518 -1.018 0.518 ...
## $ PURCHASES : num -0.425 -0.47 -0.108 0.232 -0.462 ...
## $ ONEOFF_PURCHASES : num -0.357 -0.357 0.109 0.546 -0.347 ...
## $ INSTALLMENTS_PURCHASES : num -0.349 -0.455 -0.455 -0.455 -0.455 ...
## $ CASH_ADVANCE : num -0.467 2.605 -0.467 -0.369 -0.467 ...
## $ PURCHASES_FREQUENCY : num -0.807 -1.222 1.27 -1.014 -1.014 ...
## $ ONEOFF_PURCHASES_FREQUENCY : num -0.679 -0.679 2.673 -0.399 -0.399 ...
## $ PURCHASES_INSTALLMENTS_FREQUENCY: num -0.707 -0.917 -0.917 -0.917 -0.917 ...
## $ CASH_ADVANCE_FREQUENCY : num -0.675 0.574 -0.675 -0.259 -0.675 ...
## $ CASH_ADVANCE_TRX : num -0.476 0.11 -0.476 -0.33 -0.476 ...
## $ PURCHASES_TRX : num -0.511 -0.592 -0.109 -0.552 -0.552 ...
## $ CREDIT_LIMIT : num -0.96 0.689 0.826 0.826 -0.905 ...
## $ PAYMENTS : num -0.529 0.819 -0.384 -0.599 -0.364 ...
## $ PRC_FULL_PAYMENT : num -0.526 0.234 -0.526 -0.526 -0.526 ...
## $ TENURE : num 0.361 0.361 0.361 0.361 0.361 ...
## - attr(*, "na.action")= 'omit' Named int 5204
## ..- attr(*, "names")= chr "5204"
range(scaled_df$BALANCE)

## [1] -0.751662  8.396726
range(scaled_df$PURCHASES)

## [1] -0.4695577 22.4812220
range(scaled_df$ONEOFF_PURCHASES)

## [1] -0.3569366 24.1984941
range(scaled_df$INSTALLMENTS_PURCHASES)

## [1] -0.4545815 24.4243905
range(scaled_df$CASH_ADVANCE)

## [1] -0.4667793 22.0087908
range(scaled_df$CREDIT_LIMIT)

## [1] -1.2214  7.0093
# And the frequency variables, which are pre-scaled, for comparison:
range(scaled_df$PURCHASES_FREQUENCY)

## [1] -1.221860  1.269671
range(scaled_df$ONEOFF_PURCHASES_FREQUENCY)

## [1] -0.6786783  2.6731453
range(scaled_df$PURCHASES_INSTALLMENTS_FREQUENCY)

## [1] -0.9170383  1.5989932
range(scaled_df$CASH_ADVANCE_FREQUENCY)

## [1] -0.6752567  6.8197856

```

Oh, goodie! I'm having trouble interpreting that already. If we have trouble getting to the bottom of all that - we'll stick with the basic, student-model hackish version for our predictions.

MODELING

Select Modeling Technique

We have decided to try a few different models based on the K-means algorithm. Our models will differ by the inputs they take, and how those inputs are cleaned and scaled.

This will allow us to compare the resulting models easily, and determine which data preparation techniques are best suited to the knowledge discovery we seek.

Generate Test Design

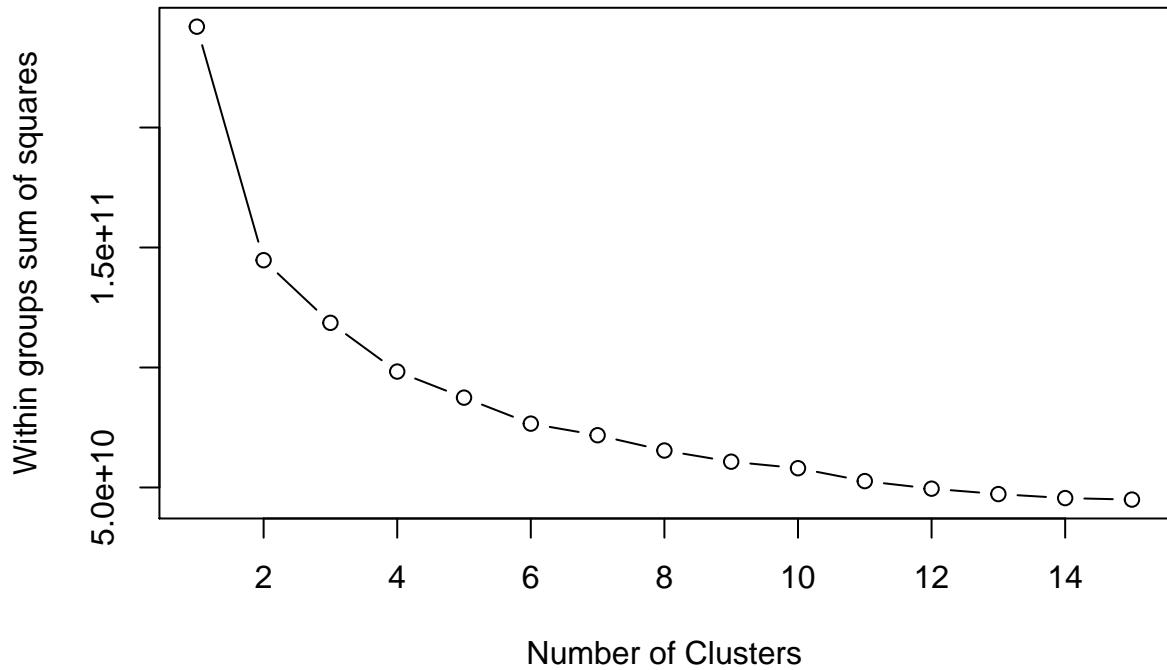
Model #1 - clean data with capped outliers.

The best thing about this model was that we were able to generate it quickly, and use it to build up the first implementations of our Shiny App. We do have high hopes that our more advanced models will yeild better results. We shall check that assumption!

Determine number of clusters

Run tests on various numbers of clusters to look for an “elbow” which will suggest how many useful clusters we can hope to get from this data.

```
cluster_data_1 <- df_capped[c(2,4:7, 14)]  
  
wss <- (nrow(cluster_data_1)-1)*sum(apply(cluster_data_1,2,var))  
for (i in 2:15) wss[i] <- sum(kmeans(cluster_data_1,  
  centers=i)$withinss)  
plot(1:15, wss, type="b", xlab="Number of Clusters",  
  ylab="Within groups sum of squares")
```



In terms of how many clusters we should choose, these graphs could be suggesting anywhere from 4 to 7, depending on how we read them.

Here's one with 6 clusters:

```
# K-Means Cluster Analysis
fit <- kmeans(cluster_data_1, 6) # 6 cluster solution
# get cluster means
aggregate(cluster_data_1, by=list(fit$cluster), FUN=mean)

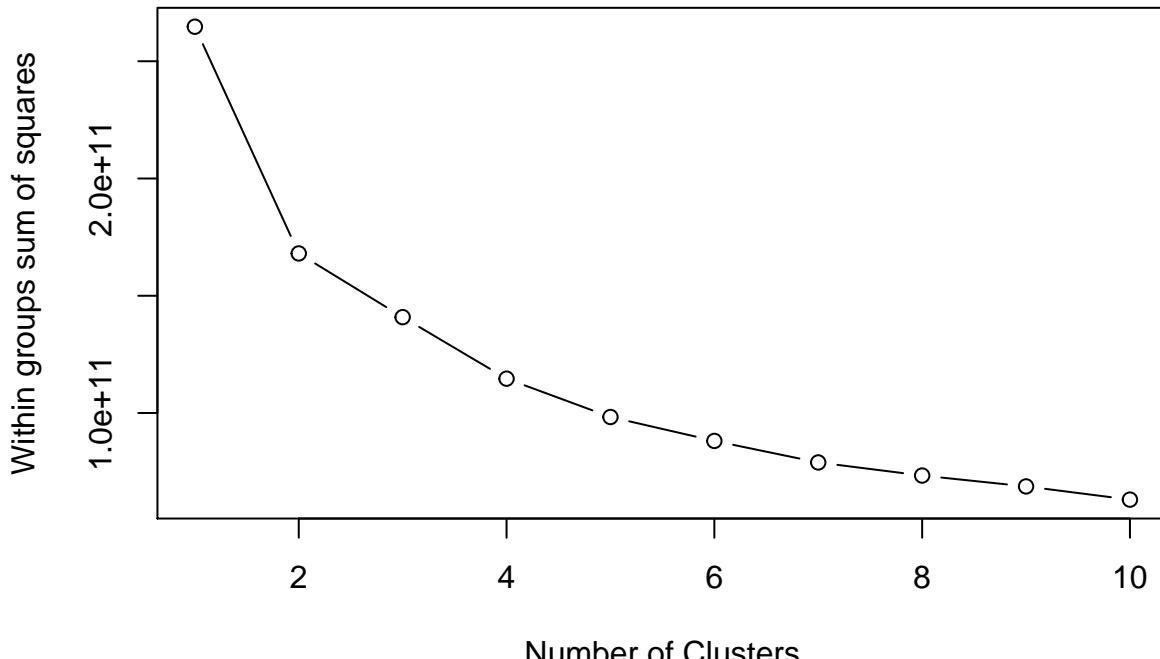
##   Group.1    BALANCE PURCHASES ONEOFF_PURCHASES INSTALLMENTS_PURCHASES
## 1      1  840.2367  1586.7179       961.0053        626.5743
## 2      2 4164.0772 10562.1790      6402.0114       3223.3283
## 3      3 3789.6287  567.5871      305.5582        262.0607
## 4      4 1895.3129 1688.8598      1072.5319       616.3683
## 5      5 7356.4406  855.5663      511.9015       343.8396
## 6      6  765.4780  474.0965      225.0695       249.2429
##   CASH_ADVANCE CREDIT_LIMIT
## 1      170.3353     6101.162
## 2      673.7984    11109.012
## 3      2931.4048     6011.828
## 4      461.6586    12147.064
## 5      6769.9186    11437.598
## 6      444.8720    2034.965

# append cluster assignment
cluster_data_1 <- data.frame(cluster_data_1, fit$cluster)
## Assess Model
```

Sure, but how does this look WITHOUT all that capping?

```
df0 <- read.csv("credit-card-cust-behav-data.csv", stringsAsFactors = FALSE, header = TRUE, encoding =
df0 <- na.omit(df0)
cld2 <- df0[c(2,4:7,14)]
```

```
wss <- (nrow(cld2)-1)*sum(apply(cld2,2,var))
for (i in 2:10) wss[i] <- sum(kmeans(cld2,
  centers=i)$withinss)
plot(1:10, wss, type="b", xlab="Number of Clusters",
  ylab="Within groups sum of squares")
```



And a clustering into 6, with the uncapped data, yeilds almost identical results.

```
# K-Means Cluster Analysis
fit <- kmeans(cld2, 6) # 6 cluster solution
# get cluster means
aggregate(cld2, by=list(fit$cluster), FUN=mean)

##   Group.1    BALANCE PURCHASES ONEOFF_PURCHASES INSTALLMENTS_PURCHASES
## 1      1 1879.4951    1403.5947       835.8505        567.7896
## 2      2 6361.6936    654.0964      382.8017        271.4090
## 3      3 1759.8680   1072.5659      615.9646        457.0468
## 4      4  793.0308    513.2348      246.6321        266.9216
## 5      5 3712.2097   7121.0150     4842.3952      2278.6198
## 6      6 5390.3896  27690.8658     21422.8846      6267.9812
##   CASH_ADVANCE CREDIT_LIMIT
## 1      464.1035    12448.244
## 2     6243.6069     9811.185
## 3      927.0586    6084.161
## 4      482.8431    2061.928
## 5      568.1552    9277.667
## 6     929.6892    16333.333

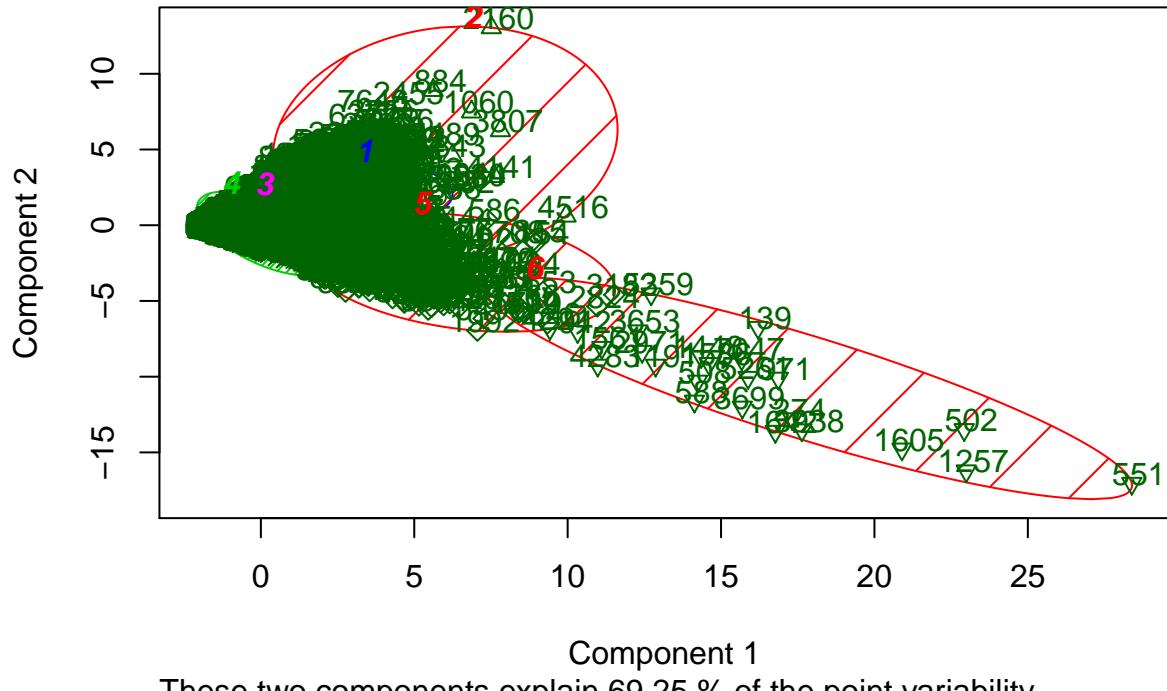
# append cluster assignment
cld2 <- data.frame(cld2, fit$cluster)
## Assess Model
```

This 6-cluster solution clearly nets us 6 groups we can identify - Big Spenders, Cash Advancers, Smaller Cash Advancers, Balance Carriers, Balance Payers, and Light Users.

It is interesting that whether or not we cap, with 6 clusters, we get two groups of “cash advancers”, which differ only in terms of credit limit.

```
clusplot(cld2, fit$cluster, color=TRUE, shade=TRUE,
         labels=2, lines=0)
```

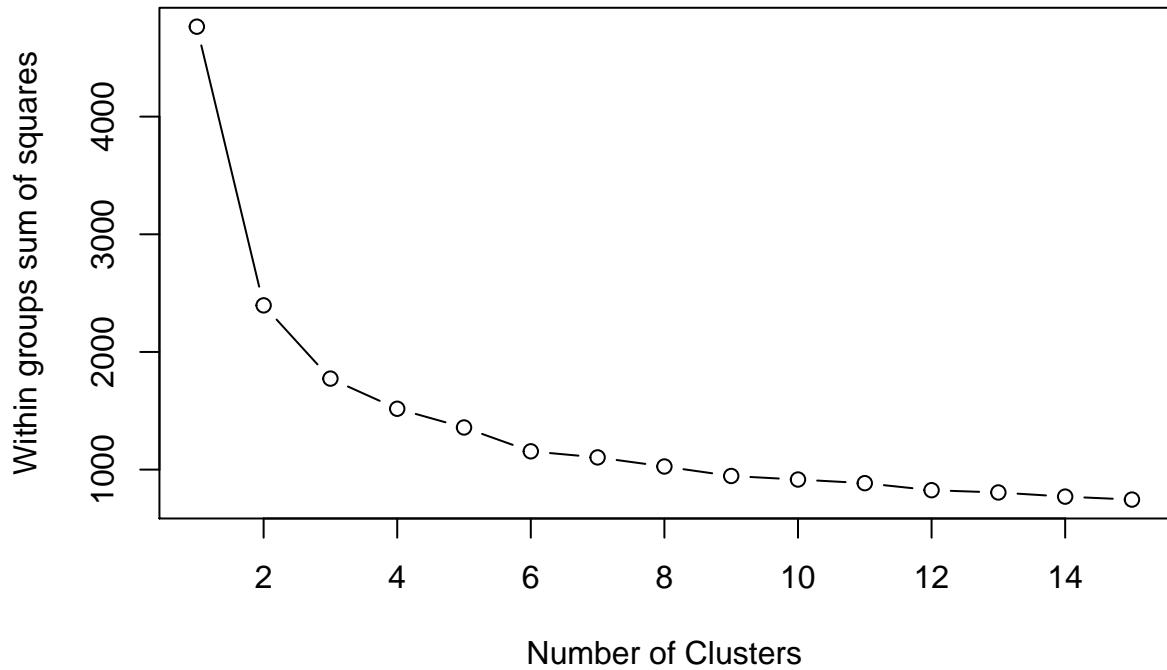
CLUSPLOT(cld2)



These two components explain 69.25 % of the point variability.

Model 2 - With hackishly scaled data

```
cluster_data_2 <- hackish_scaled_df[c(2,4:11, 14)]  
  
wss <- (nrow(cluster_data_2)-1)*sum(apply(cluster_data_2,2,var))  
for (i in 2:15) wss[i] <- sum(kmeans(cluster_data_2,  
           centers=i)$withinss)  
plot(1:15, wss, type="b", xlab="Number of Clusters",  
     ylab="Within groups sum of squares")
```



```

# K-Means Cluster Analysis
fit_2 <- kmeans(cluster_data_2, 6) # 6 cluster solution
# get cluster means
aggregate(cluster_data_2,by=list(fit_2$cluster),FUN=mean)

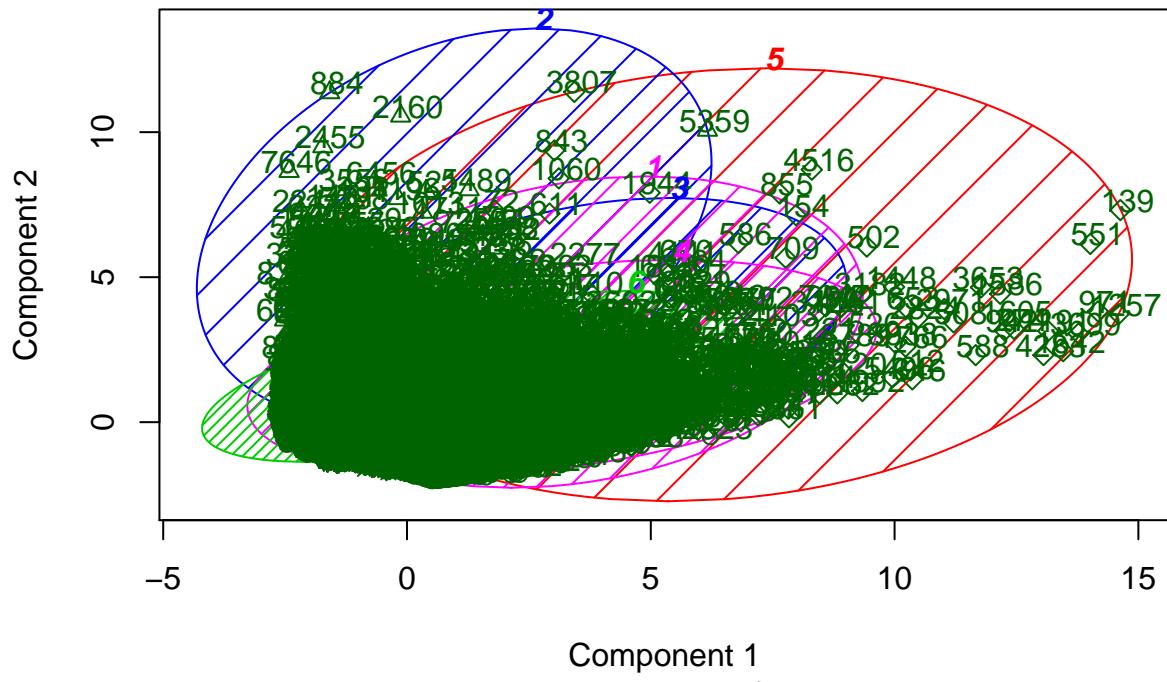
##   Group.1      BALANCE    PURCHASES ONEOFF_PURCHASES INSTALLMENTS_PURCHASES
## 1       1 0.05433250 0.051144899      0.01926841      0.082717717
## 2       2 0.19512014 0.009984959      0.01586191      0.002919724
## 3       3 0.07382462 0.090582823      0.16521636      0.012295215
## 4       4 0.05013546 0.034025443      0.02419485      0.042988239
## 5       5 0.12200038 0.211244629      0.26018529      0.151636150
## 6       6 0.05384008 0.011814957      0.02111173      0.002193311
##   CASH_ADVANCE PURCHASES_FREQUENCY ONEOFF_PURCHASES_FREQUENCY
## 1     0.02520970          0.95370860          0.07994499
## 2     0.19885101          0.08473351          0.05893067
## 3     0.02637644          0.78297722          0.75179043
## 4     0.01950801          0.54666717          0.09492730
## 5     0.03875876          0.97007256          0.77365825
## 6     0.03230485          0.09768933          0.07382961
##   PURCHASES_INSTALLMENTS_FREQUENCY CASH_ADVANCE_FREQUENCY CREDIT_LIMIT
## 1           0.91403511          0.06851220      0.1836167
## 2           0.02598998          0.49208639      0.3356844
## 3           0.11038313          0.08302058      0.2782671
## 4           0.45271620          0.06290382      0.1977867
## 5           0.84570220          0.10124346      0.3904970
## 6           0.02312163          0.10783468      0.1624301

# append cluster assignment
cluster_data_2 <- data.frame(cluster_data_2, fit_2$cluster)
## Assess Model

clusplot(cluster_data_2, fit_2$cluster, color=TRUE, shade=TRUE,
         labels=2, lines=0)

```

CLUSPLOT(cluster_data_2)



Component 1

These two components explain 56.36 % of the point variability.

And Let's try that with 8 clusters, to see if we can make use of that:

```
# K-Means Cluster Analysis
fit_2_8 <- kmeans(cluster_data_2, 8) # 8 cluster solution
# get cluster means
aggregate(cluster_data_2,by=list(fit_2_8$cluster),FUN=mean)

##      Group.1      BALANCE PURCHASES ONEOFF_PURCHASES INSTALLMENTS_PURCHASES
## 1          1 0.12200038 0.211244629      0.26018529      0.151636150
## 2          2 0.05384008 0.011814957      0.02111173      0.002193311
## 3          3 0.05808952 0.042809612      0.02902159      0.055072285
## 4          4 0.04303042 0.026178910      0.01988332      0.032194063
## 5          5 0.07382462 0.090582823      0.16521636      0.012295215
## 6          6 0.12690762 0.006846774      0.01222389      0.001475236
## 7          7 0.05433250 0.051144899      0.01926841      0.082717717
## 8          8 0.29326426 0.014500177      0.02109628      0.004998053
##      CASH_ADVANCE PURCHASES_FREQUENCY ONEOFF_PURCHASES_FREQUENCY
## 1      0.03875876           0.97007256           0.77365825
## 2      0.03230485           0.09768933           0.07382961
## 3      0.02520289           0.66274530           0.11993162
## 4      0.01442101           0.44297939           0.07259199
## 5      0.02637644           0.78297722           0.75179043
## 6      0.13531814           0.07863250           0.06203891
## 7      0.02520970           0.95370860           0.07994499
## 8      0.29026205           0.09351162           0.05445852
##      PURCHASES_INSTALLMENTS_FREQUENCY CASH_ADVANCE_FREQUENCY CREDIT_LIMIT
## 1                  0.84570220           0.10124346      0.3904970
## 2                  0.02312163           0.10783468      0.1624301
## 3                  0.55883825           0.08472862      0.2095712
```

```

## 4                      0.35792178          0.04340864          0.1872601
## 5                      0.11038313          0.08302058          0.2782671
## 6                      0.01758667          0.54447827          0.2140128
## 7                      0.91403511          0.06851220          0.1836167
## 8                      0.03808066          0.41670500          0.5107453

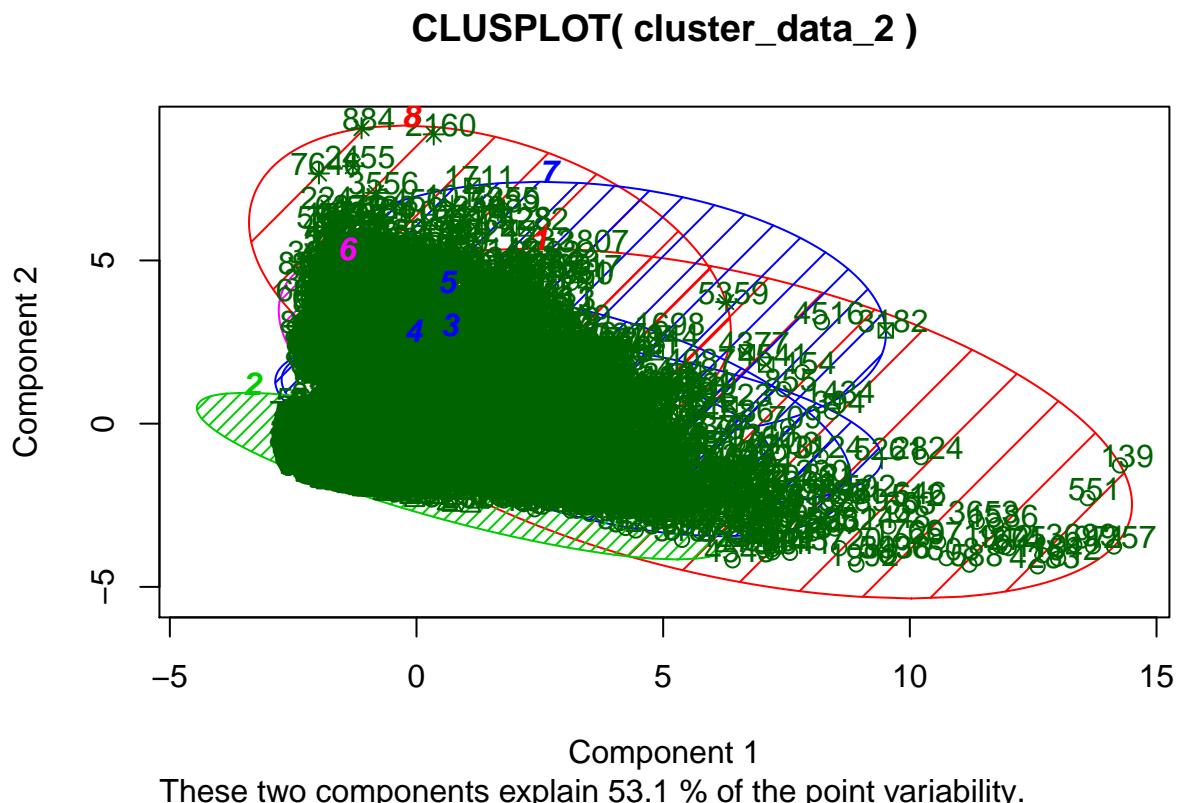
##    fit_2.cluster
## 1              5
## 2              6
## 3              4
## 4              4
## 5              3
## 6              2
## 7              1
## 8              2

# append cluster assignment
cluster_data_2 <- data.frame(cluster_data_2, fit_2_8$cluster)
## Save Model
save(fit_2_8, file = 'CreditCardBehaviour8Clusters.rda')

```

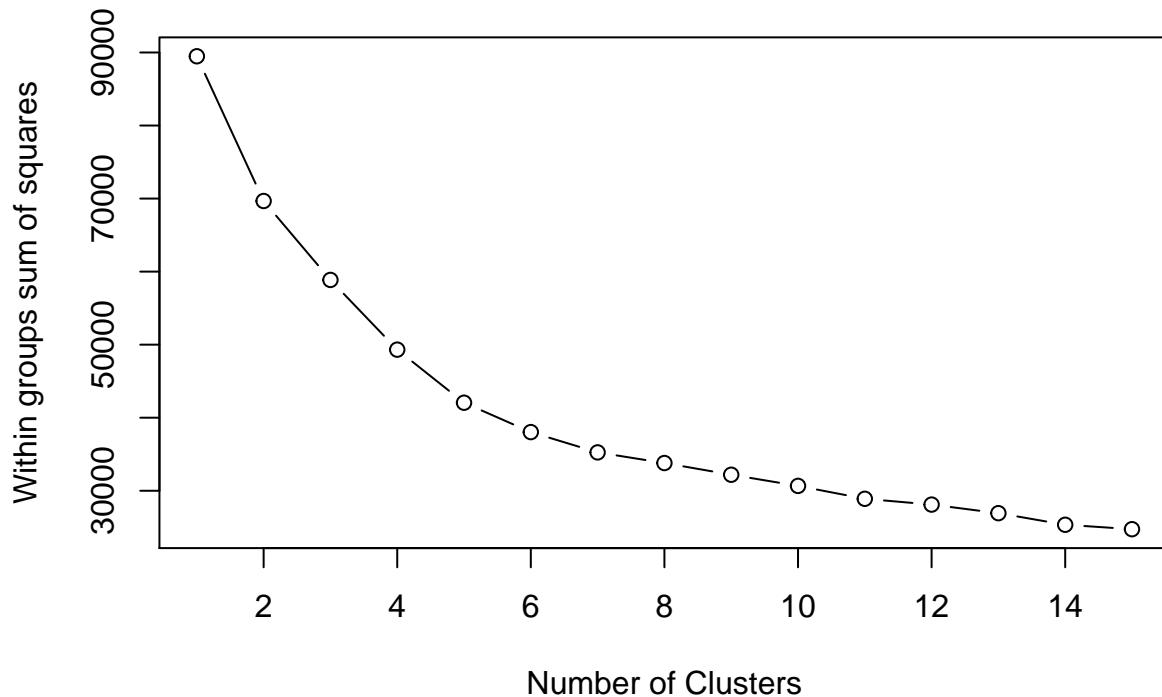
Differernt way of looking at clusters.

```
clusplot(cluster_data_2, fit_2$cluster, color=TRUE, shade=TRUE,  
        labels=2, lines=0)
```



Model 3 - With formally scaled data

```
cluster_data_3 <- scaled_df[c(2,4:11, 14)]  
  
wss <- (nrow(cluster_data_3)-1)*sum(apply(cluster_data_3,2,var))  
for (i in 2:15) wss[i] <- sum(kmeans(cluster_data_3,  
  centers=i)$withinss)  
plot(1:15, wss, type="b", xlab="Number of Clusters",  
 ylab="Within groups sum of squares")
```



```
# K-Means Cluster Analysis  
fit_3 <- kmeans(cluster_data_3, 6) # 6 cluster solution  
# get cluster means  
aggregate(cluster_data_3, by=list(fit_3$cluster), FUN=mean)
```

```
##   Group.1      BALANCE    PURCHASES ONEOFF_PURCHASES INSTALLMENTS_PURCHASES  
## 1       1  1.6749207 10.63231257      9.7005066      7.3149920  
## 2       2  2.3187049 -0.12553398     -0.1001982     -0.1128896  
## 3       3 -0.3770985 -0.31138562     -0.1972150     -0.3738458  
## 4       4  0.1682206  1.05034324      1.0160165      0.6163812  
## 5       5  0.3145946 -0.40547308     -0.2942079     -0.4182823  
## 6       6 -0.3900179 -0.05770902     -0.2505818      0.3242015  
##   CASH_ADVANCE PURCHASES_FREQUENCY ONEOFF_PURCHASES_FREQUENCY  
## 1   -0.1058713          1.0058615          1.7900804  
## 2    2.4403901          -0.3261965         -0.1389410  
## 3   -0.3317747          -0.7006330         -0.2967531  
## 4   -0.2833188           1.0706772          1.9024903  
## 5    0.5715518          -0.9813747         -0.4843687  
## 6   -0.3532189           0.9845943         -0.3837720  
##   PURCHASES_INSTALLMENTS_FREQUENCY CASH_ADVANCE_FREQUENCY CREDIT_LIMIT  
## 1                  0.9607913          -0.5282953      3.00912962
```

```

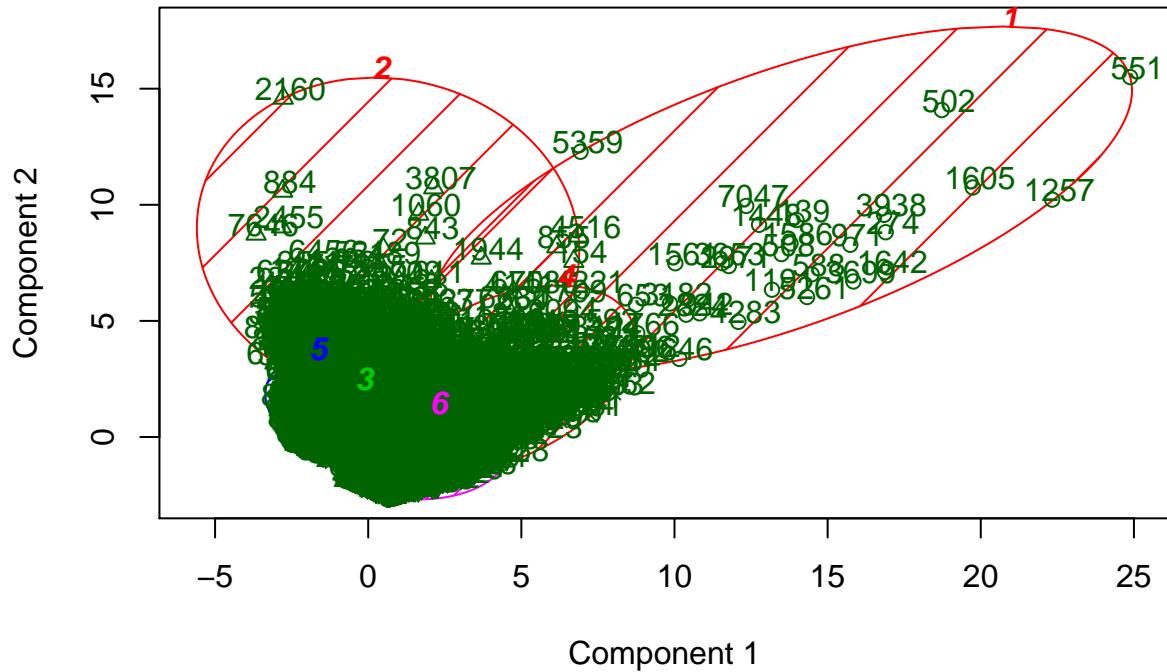
## 2 -0.2798139 1.7204331 1.52409476
## 3 -0.6799667 -0.3821199 -0.34480199
## 4 0.4954976 -0.3491326 0.72100068
## 5 -0.8173481 1.2168178 -0.09354128
## 6 1.1743109 -0.4456686 -0.32100574

# append cluster assignment
cluster_data_3 <- data.frame(cluster_data_3, fit_3$cluster)
## Assess Model

clusplot(cluster_data_3, fit_3$cluster, color=TRUE, shade=TRUE,
         labels=2, lines=0)

```

CLUSPLOT(cluster_data_3)



Component 1

These two components explain 56.86 % of the point variability.

And Let's try that with 11 clusters, to see if we can make use of that:

```

# K-Means Cluster Analysis
fit_3_11 <- kmeans(cluster_data_3, 11) # 11 cluster solution
# get cluster means
aggregate(cluster_data_3, by=list(fit_3_11$cluster), FUN=mean)

```

	Group.1	BALANCE	PURCHASES	ONEOFF_PURCHASES	INSTALLMENTS_PURCHASES
## 1	1	1.70291649	13.1473920	14.1141157	5.15614584
## 2	2	-0.24442132	0.3108626	0.4630957	-0.11587624
## 3	3	1.54280864	4.4284241	1.8510025	7.07569890
## 4	4	-0.50691575	-0.1878307	-0.2839620	0.07799801
## 5	5	-0.04759961	0.3224026	-0.1494373	1.03568127
## 6	6	2.29511903	-0.1635101	-0.1272388	-0.15298015
## 7	7	-0.38055597	-0.3237528	-0.2136773	-0.37284065
## 8	8	0.63551731	1.8479884	1.7754251	1.10705638
## 9	9	-0.06916501	-0.4317587	-0.3168565	-0.43884208

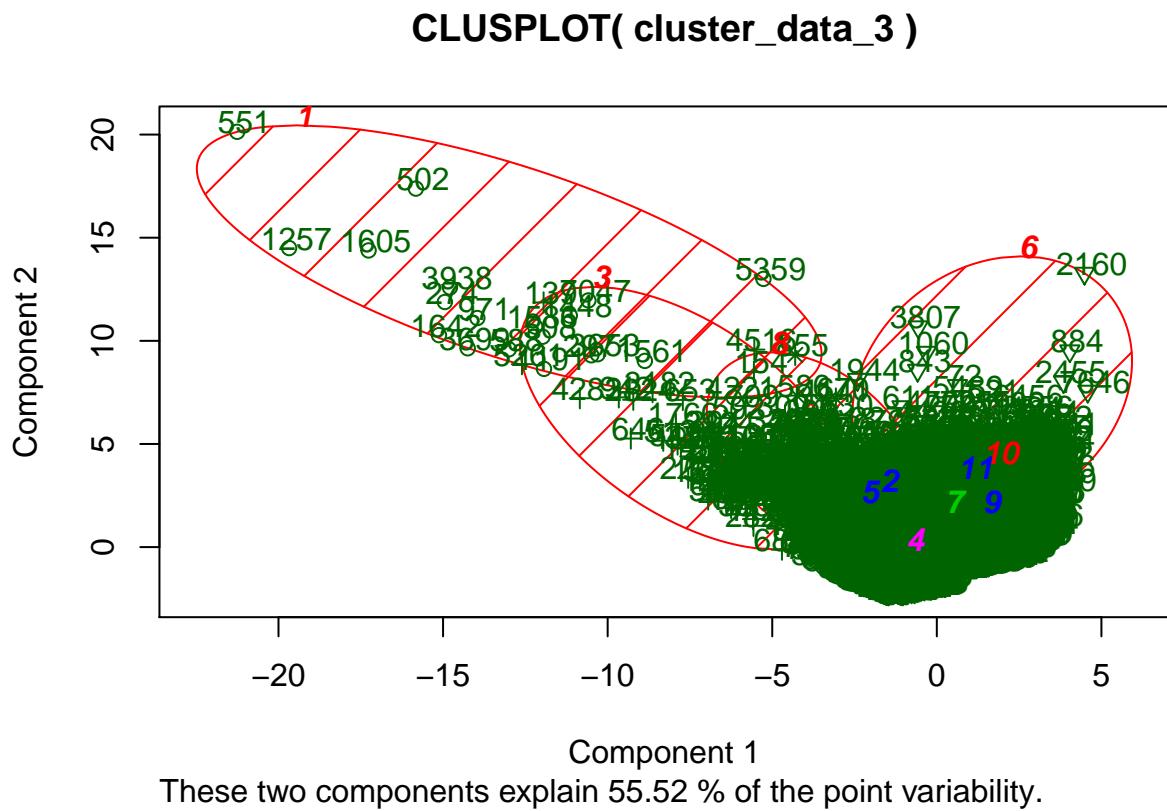
```

## 10      10  0.25253702 -0.3105141      -0.2128622      -0.34313298
## 11      11  0.85743030 -0.4046682      -0.2918610      -0.42071547
##   CASH_ADVANCE PURCHASES_FREQUENCY ONEOFF_PURCHASES_FREQUENCY
## 1    0.06516462      0.9976785      2.1368537
## 2   -0.36015411      0.9284430      1.9686427
## 3   -0.08532871      1.1256710      1.2355965
## 4   -0.38751747      0.9636146      -0.4472664
## 5   -0.27292472      1.0374092      -0.1890803
## 6    2.43737620     -0.3422893      -0.1579093
## 7   -0.32915945     -0.7438549      -0.3635702
## 8   -0.18721491      1.1421597      1.7986222
## 9    0.27158330     -1.0986574      -0.5718217
## 10   0.76766928     -0.5085803      -0.1526789
## 11   0.83473642     -1.0279929      -0.5246229
##   PURCHASES_INSTALLMENTS_FREQUENCY CASH_ADVANCE_FREQUENCY CREDIT_LIMIT
## 1                      0.75612251      -0.4462417      3.1825054
## 2                     -0.07923163      -0.4392303      0.2556730
## 3                      1.36534414      -0.3125934      1.8443351
## 4                      1.14589739      -0.4891249      -0.5929776
## 5                      1.25010229      -0.3768079      0.5136645
## 6                     -0.29805503      1.7207975      1.5052698
## 7                     -0.67638115      -0.3789989      -0.3551891
## 8                      1.05566701      -0.2518005      1.1678620
## 9                     -0.87568880      1.0628086      -0.5580208
## 10                     -0.54330416      2.6034568      -0.1515324
## 11                     -0.83913644      0.6247314      0.5653652
##   fit_3.cluster
## 1      1.000000
## 2      3.855670
## 3      3.290323
## 4      6.000000
## 5      5.919129
## 6      2.000000
## 7      3.000000
## 8      3.988550
## 9      5.000000
## 10     5.013559
## 11     4.978000

# append cluster assignment
cluster_data_3 <- data.frame(cluster_data_3, fit_3_11$cluster)
## Save Model
#save(fit_3_11 , file = 'CreditCardBehaviour11Clusters.rda')

clusplot(cluster_data_3, fit_3_11$cluster, color=TRUE, shade=TRUE,
         labels=2, lines=0)

```



EVALUATION

Evaluation? Coming soon.

Evaluate Results

Approved Models

Determine Next Steps

DEPLOYMENT

Plan Deployment

Plan Monitoring and Maintenance

Produce Final Report ————— KNIT THIS .Rmd FILE WHEN DONE!!!!

Review Project