

Stat 412

Yolanda Jin

24/11/2021

Contents

0. Cleaning Up Data and Creating Data Groups	2
a) Remove NA Data	2
b) Split the data to Testing/Training & Minority, Majority Groups	3
Creating Indicator Variables for Extreme Data	15
1. Question/Goal	17
2. EDA (Still need to Clean up for better graphs)	17
3. Balancing Data (Creating New Datasets)	19
a) Boostrapping Minority Data (cs_train_min_add) (1000 obs)	19
b) Resampling Majority Data (new_train.final 1 & 2)	21
c) Using Clustering to select from Majority Group	23
4. Applying Different Methods (with and without Balanced Data)	24
Method 1: PCA	24
1-a) PCA with Unbalanced Data	24
PCA using all 22 variables	24
PCA using original 10 variables	26
PCA Using Dummy variables	27
1 b) PCA with Balanced Data	28
Method 2: Using GLM Original Predictors vs Extreme Binned Predictors + MonthlyIncome	28
2 a) GLM Original Predictors with Unbalanced Data	28
2 b) GLM Extreme Binned Predictors with Unbalanced Data	33
2 c) GLM Original Predictors with Balanced Data	36
Method 3: Ridge	36
3 a) Ridge with Unbalanced Data	36
3 b) Ridge with Balanced Data	36
Method 4: Lasso	36
4 a) Lasso with Unbalanced Data	36

4 b) Lasso with Balanced Data	37
Method 4: Random Forest	37
4 a) Random Forest with Unbalanced Data	37
4 b) Random Forest with Balanced Data	38
5. Evaluating Model/Comparing results	38
5-1 Evaluation on Final Model using Training Data	38
5-2 Marginal Model Plots	39
Need to Update!!	39
5-3 Goodness of Fit using Hosmer-Lemeshow Test	40
AUC	42
Need to fix the Prediction function!!!	42
Maybe don't need session	45
(0) Ensemble Learning Used by Paper	45
Not Sure if we use this!!	45
Cluster Attempt	46
Maybe don't need session (END)	47

0. Cleaning Up Data and Creating Data Groups

a) Remove NA Data

- Split out data with NA in any predictors
- split non-NA group to minority (default) vs majority (non-default) group
- Additional TODO: can check out characteristics of NA group so we can replace the NA values

```
# Read Credit Scoring Data Training Set
#cs_train <- cs_training
cs_train = read.csv("cs-training.csv")
train <- cs_train
raw_data <- cs_train      #150k rows

# Remove NA cases otherwise cannot predict
cs_train.omit <- na.omit(raw_data)  # 150k -> 120,269 obs
Predictor_Variables <- subset.data.frame(cs_train.omit, select = c(RevolvingUtilizationOfUnsecuredLines))

```

```
summary(cs_train)
```

```
##          X      SeriousDlqin2yrs RevolvingUtilizationOfUnsecuredLines
##  Min.   : 1   Min.   :0.00000   Min.   : 0.00
##  1st Qu.: 37501 1st Qu.:0.00000   1st Qu.: 0.03
##  Median : 75001 Median :0.00000   Median : 0.15
##  Mean   : 75001 Mean   :0.06684   Mean   : 6.05
```

```

## 3rd Qu.:112500   3rd Qu.:0.00000   3rd Qu.:    0.56
## Max.    :150000   Max.    :1.00000   Max.    :50708.00
##
##      age      NumberOfTime30.59DaysPastDueNotWorse  DebtRatio
##  Min.    : 0.0      Min.    : 0.000                  Min.    : 0.0
##  1st Qu.: 41.0     1st Qu.: 0.000                  1st Qu.: 0.2
##  Median : 52.0     Median : 0.000                  Median : 0.4
##  Mean   : 52.3     Mean   : 0.421                  Mean   : 353.0
##  3rd Qu.: 63.0     3rd Qu.: 0.000                  3rd Qu.: 0.9
##  Max.   :109.0     Max.   :98.000                  Max.   :329664.0
##
##  MonthlyIncome   NumberOfOpenCreditLinesAndLoans  NumberOfTimes90DaysLate
##  Min.    :    0      Min.    : 0.000                  Min.    : 0.000
##  1st Qu.: 3400    1st Qu.: 5.000                  1st Qu.: 0.000
##  Median : 5400    Median : 8.000                  Median : 0.000
##  Mean   : 6670    Mean   : 8.453                  Mean   : 0.266
##  3rd Qu.: 8249    3rd Qu.:11.000                  3rd Qu.: 0.000
##  Max.   :3008750   Max.   :58.000                  Max.   :98.000
##  NA's    :29731
##
##  NumberRealEstateLoansOrLines  NumberOfTime60.89DaysPastDueNotWorse
##  Min.    : 0.000          Min.    : 0.0000
##  1st Qu.: 0.000          1st Qu.: 0.0000
##  Median : 1.000          Median : 0.0000
##  Mean   : 1.018          Mean   : 0.2404
##  3rd Qu.: 2.000          3rd Qu.: 0.0000
##  Max.   :54.000          Max.   :98.0000
##
##  NumberOfDependents
##  Min.    : 0.000
##  1st Qu.: 0.000
##  Median : 0.000
##  Mean   : 0.757
##  3rd Qu.: 1.000
##  Max.   :20.000
##  NA's    :3924

```

b) Split the data to Testing/Training & Minority, Majority Groups

```

# Sample 60% of data for training Purpose
set.seed(1)
n = nrow(cs_train.omit)
na.idx = raw_data$X[-cs_train.omit$X] # indexes of data with NA values that we removed
n.idx = sample(n, n*0.6) # Indexes for test train split

cs_train = cs_train.omit[n.idx,] # Training data 72,161 obs
cs_test = cs_train.omit[-n.idx,] # Testing data 48,108 obs.
cs_NA = raw_data[na.idx,] # data with NA value 29,731 obs

## 1. Separate minority data vs majority data vs NA data total 72,161 obs
cs_train_min <- cs_train[cs_train$SeriousDlqin2yrs==1,] # (omit) 10,026 -> (training) 5,009 obs
cs_train_maj <- cs_train[cs_train$SeriousDlqin2yrs==0,] # (omit) 139,974 -> (training) 67,152 obs

```

```
str(cs_train)
```

```
## 'data.frame': 72161 obs. of 12 variables:  
## $ X : int 30484 74216 54077 86784 14388 31442 40844 145293 17360  
## $ SeriousDlqin2yrs : int 0 0 0 0 0 0 0 0 1 0 ...  
## $ RevolvingUtilizationOfUnsecuredLines: num 0.09474 0.05277 0.73665 0.02366 0.00582 ...  
## $ age : int 47 68 36 36 44 40 59 65 50 49 ...  
## $ NumberOfTime30.59DaysPastDueNotWorse: int 0 0 0 0 0 1 0 0 0 0 ...  
## $ DebtRatio : num 0.3684 0.2276 0.0945 0.2171 0.1853 ...  
## $ MonthlyIncome : int 6250 11884 2000 5600 5100 2946 3950 10500 18381 10796  
## $ NumberOfOpenCreditLinesAndLoans : int 8 13 6 9 24 17 7 23 10 7 ...  
## $ NumberOfTimes90DaysLate : int 0 0 0 0 0 0 0 0 0 0 ...  
## $ NumberRealEstateLoansOrLines : int 1 2 0 1 1 3 0 2 1 2 ...  
## $ NumberOfTime60.89DaysPastDueNotWorse: int 0 0 0 0 0 0 0 1 0 0 ...  
## $ NumberOfDependents : int 3 1 1 0 2 0 0 0 1 0 ...  
## - attr(*, "na.action")= 'omit' Named int [1:29731] 7 9 17 33 42 53 59 63 72 87 ...  
## ..- attr(*, "names")= chr [1:29731] "7" "9" "17" "33" ...
```

```
str(cs_train_min)
```

```
## 'data.frame': 4941 obs. of 12 variables:  
## $ X : int 17360 90565 146254 102792 49072 142766 73164 75761 7780  
## $ SeriousDlqin2yrs : int 1 1 1 1 1 1 1 1 1 1 ...  
## $ RevolvingUtilizationOfUnsecuredLines: num 0.667 0.942 1 1.026 0.659 ...  
## $ age : int 50 59 28 30 43 37 46 47 37 42 ...  
## $ NumberOfTime30.59DaysPastDueNotWorse: int 0 1 98 1 0 0 0 0 98 0 ...  
## $ DebtRatio : num 0.103 0.1653 0 0.2348 0.0897 ...  
## $ MonthlyIncome : int 18381 8008 1664 2763 3900 16666 4200 11000 6000 6450 ...  
## $ NumberOfOpenCreditLinesAndLoans : int 10 6 0 9 4 22 7 14 0 12 ...  
## $ NumberOfTimes90DaysLate : int 0 5 98 0 2 0 0 0 98 0 ...  
## $ NumberRealEstateLoansOrLines : int 1 0 0 0 0 3 1 1 0 2 ...  
## $ NumberOfTime60.89DaysPastDueNotWorse: int 1 0 98 2 1 0 1 0 98 0 ...  
## $ NumberOfDependents : int 1 0 0 3 2 0 1 0 1 0 ...  
## - attr(*, "na.action")= 'omit' Named int [1:29731] 7 9 17 33 42 53 59 63 72 87 ...  
## ..- attr(*, "names")= chr [1:29731] "7" "9" "17" "33" ...
```

```
str(cs_train_maj)
```

```
## 'data.frame': 67220 obs. of 12 variables:  
## $ X : int 30484 74216 54077 86784 14388 31442 40844 145293 10260  
## $ SeriousDlqin2yrs : int 0 0 0 0 0 0 0 0 0 0 ...  
## $ RevolvingUtilizationOfUnsecuredLines: num 0.09474 0.05277 0.73665 0.02366 0.00582 ...  
## $ age : int 47 68 36 36 44 40 59 65 49 53 ...  
## $ NumberOfTime30.59DaysPastDueNotWorse: int 0 0 0 0 0 1 0 0 0 0 ...  
## $ DebtRatio : num 0.3684 0.2276 0.0945 0.2171 0.1853 ...  
## $ MonthlyIncome : int 6250 11884 2000 5600 5100 2946 3950 10500 10796 3741 ...  
## $ NumberOfOpenCreditLinesAndLoans : int 8 13 6 9 24 17 7 23 7 7 ...  
## $ NumberOfTimes90DaysLate : int 0 0 0 0 0 0 0 0 0 0 ...  
## $ NumberRealEstateLoansOrLines : int 1 2 0 1 1 3 0 2 2 2 ...  
## $ NumberOfTime60.89DaysPastDueNotWorse: int 0 0 0 0 0 0 0 0 0 0 ...  
## $ NumberOfDependents : int 3 1 1 0 2 0 0 0 0 1 ...
```

```

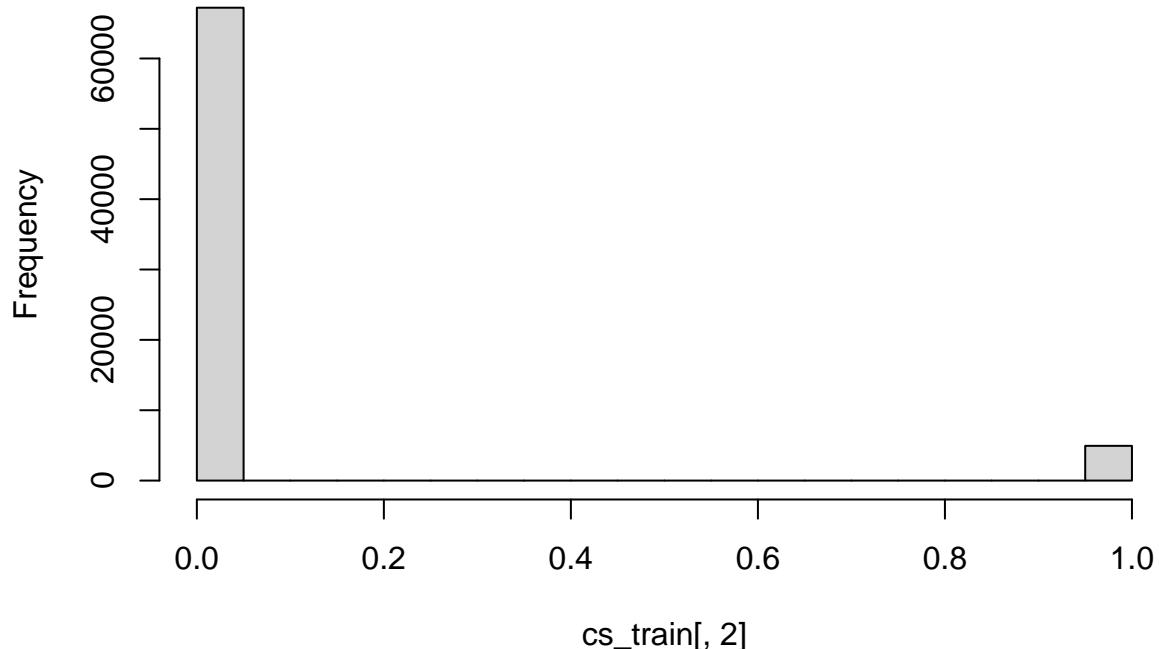
## - attr(*, "na.action")= 'omit' Named int [1:29731] 7 9 17 33 42 53 59 63 72 87 ...
## ..- attr(*, "names")= chr [1:29731] "7" "9" "17" "33" ...

# Split x and y variables
train.x = cs_train[,-which(names(cs_train) == "SeriousDlqin2yrs")]
train.x = cs_train[,-c(which(names(cs_train) == "SeriousDlqin2yrs"), which(names(cs_train) == "X"), which(names(cs_train) == "Y"))]
train.y = cs_train$SeriousDlqin2yrs
test.x = cs_test[,-which(names(cs_test) == "SeriousDlqin2yrs")]
test.x = cs_test[,-c(which(names(cs_test) == "SeriousDlqin2yrs"), which(names(cs_test) == "X"))]
test.y = cs_test$SeriousDlqin2yrs

hist(cs_train[,2]) #Response Variable

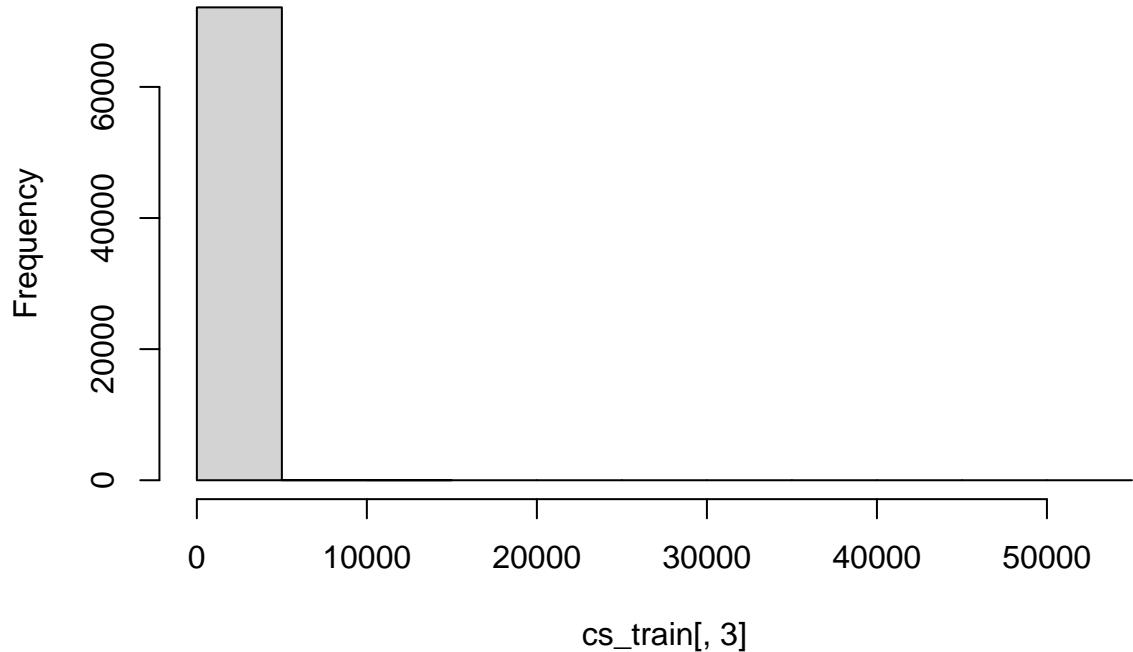
```

Histogram of cs_train[, 2]



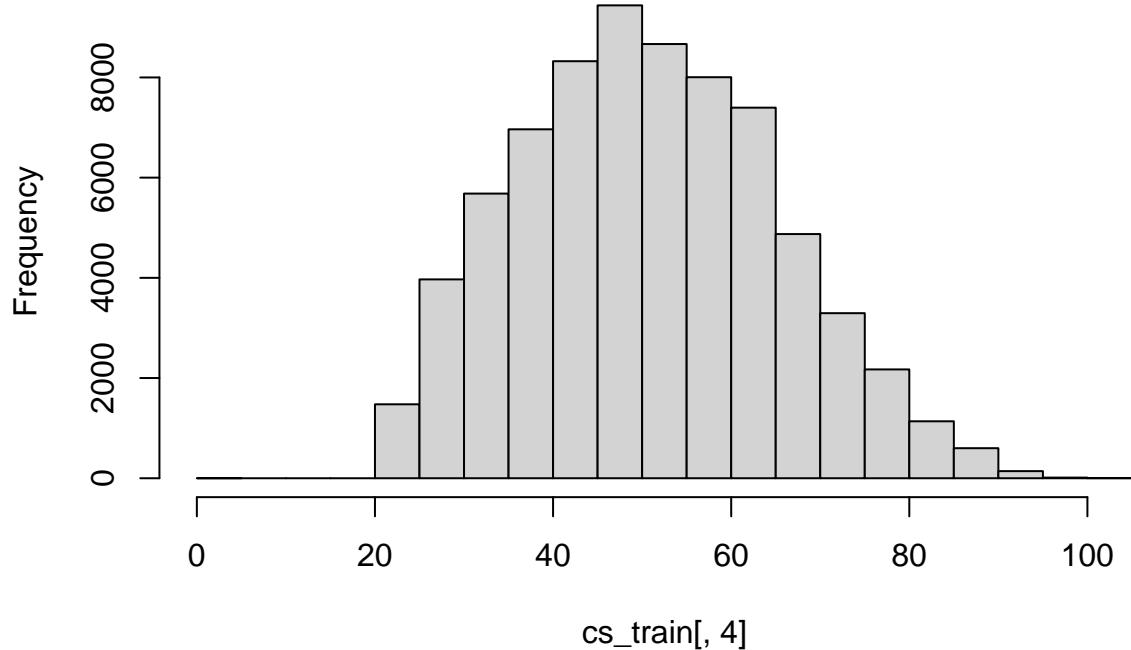
```
hist(cs_train[,3]) #RevolvingUtilizationOfUnsecuredLines
```

Histogram of cs_train[, 3]



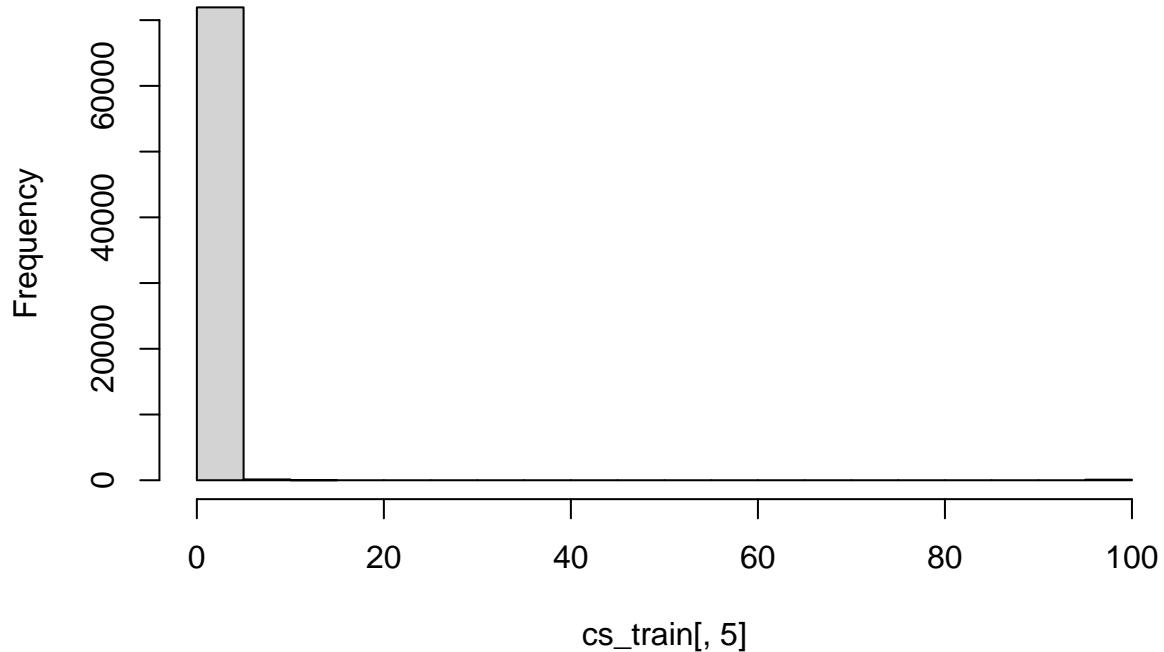
```
hist(cs_train[,4]) #age
```

Histogram of cs_train[, 4]



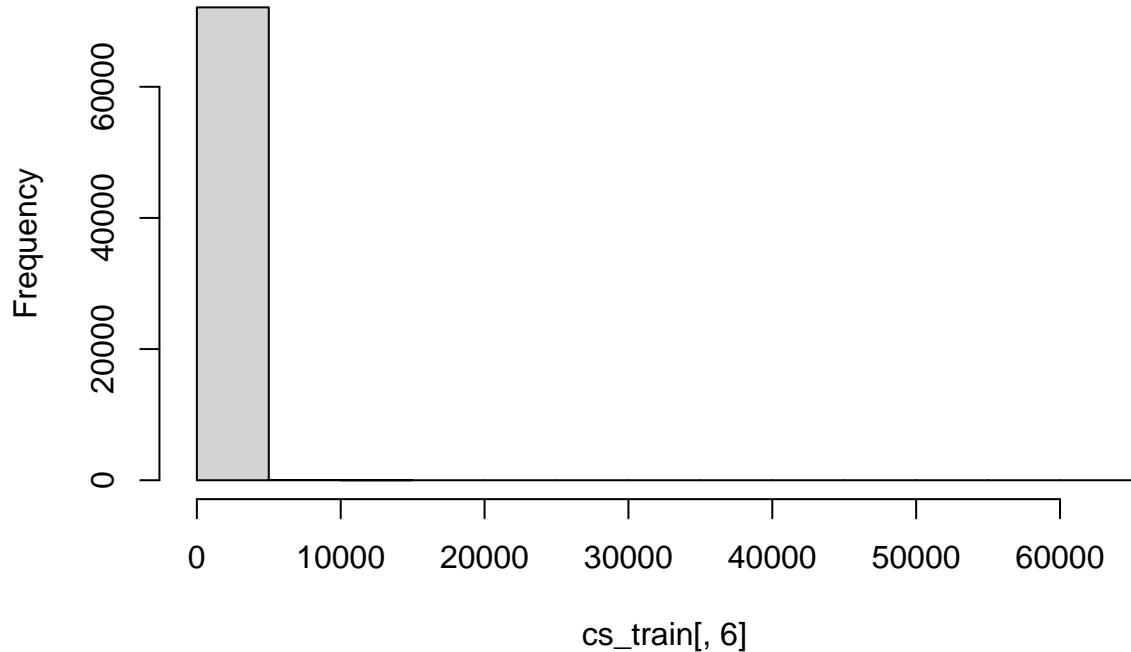
```
hist(cs_train[,5]) #Number of Time 30.59 Days Past Due Not Worse
```

Histogram of cs_train[, 5]



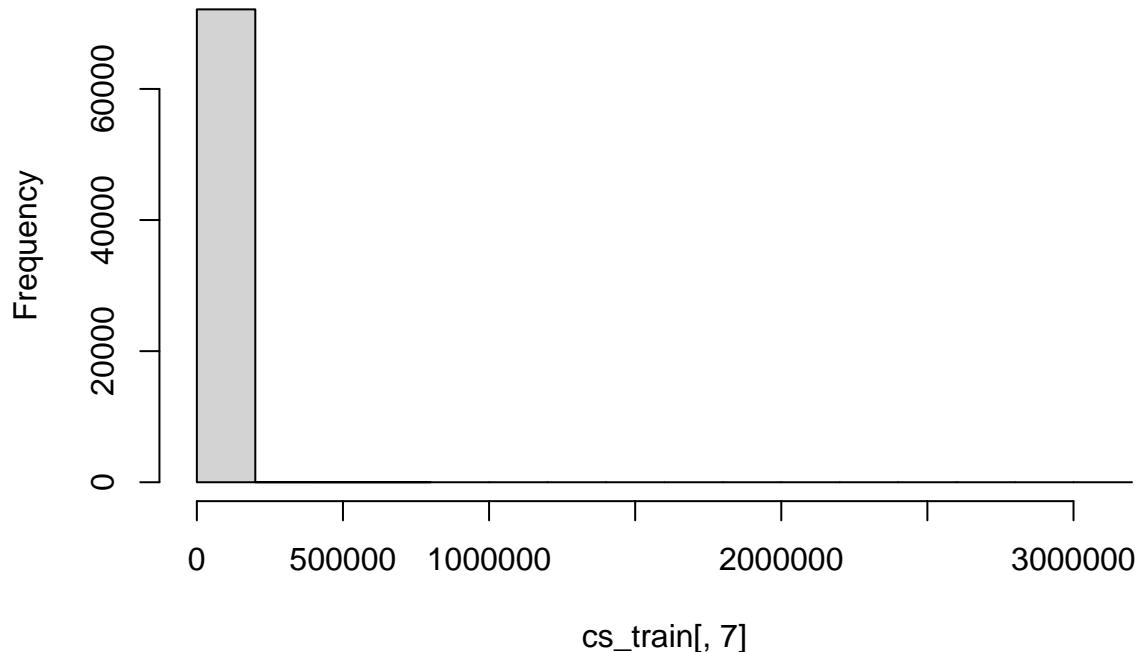
```
hist(cs_train[, 6]) #DebtRatio
```

Histogram of cs_train[, 6]



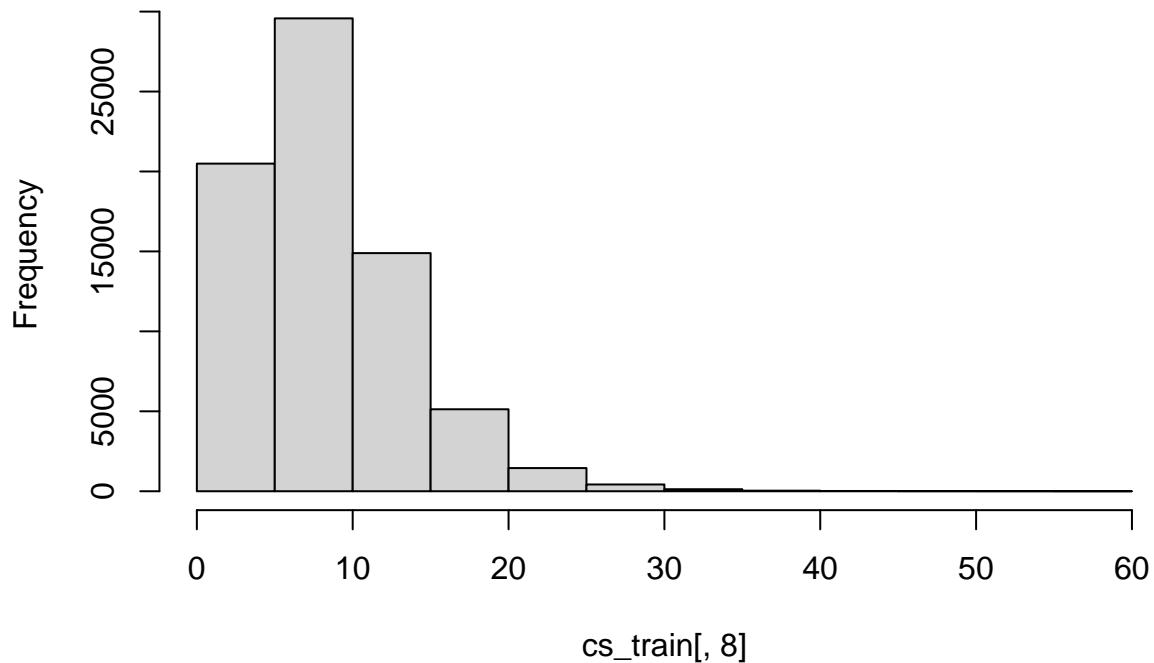
```
hist(cs_train[,7]) #MonthlyIncome
```

Histogram of cs_train[, 7]



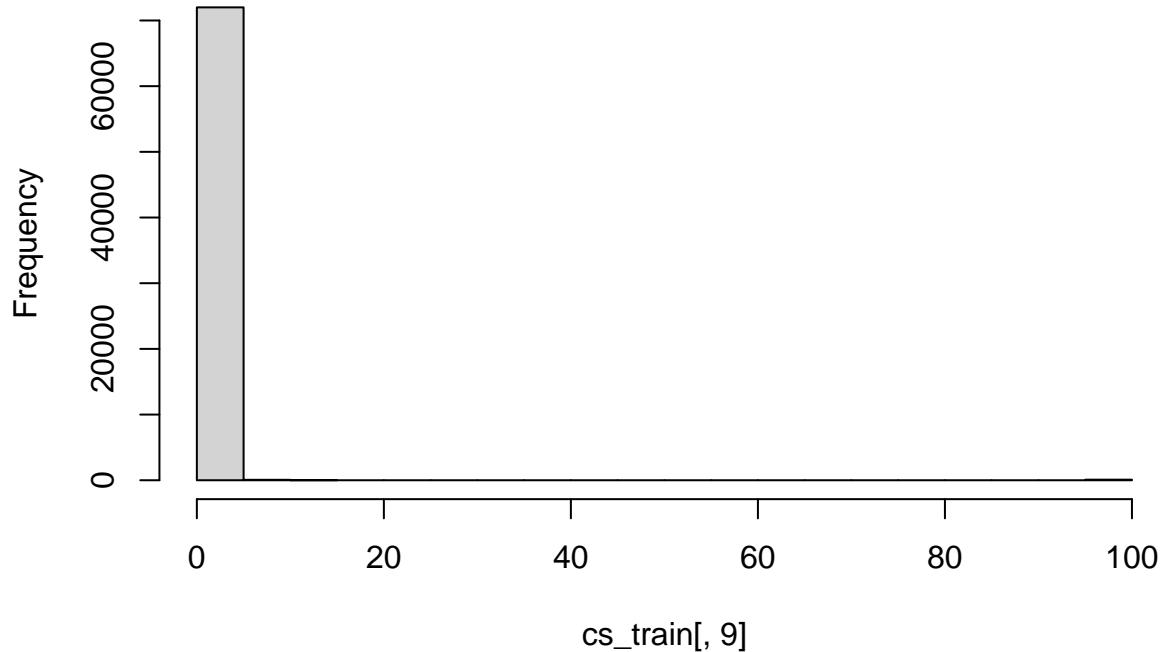
```
hist(cs_train[,8]) #NumberOfOpenCreditLinesAndLoans
```

Histogram of cs_train[, 8]



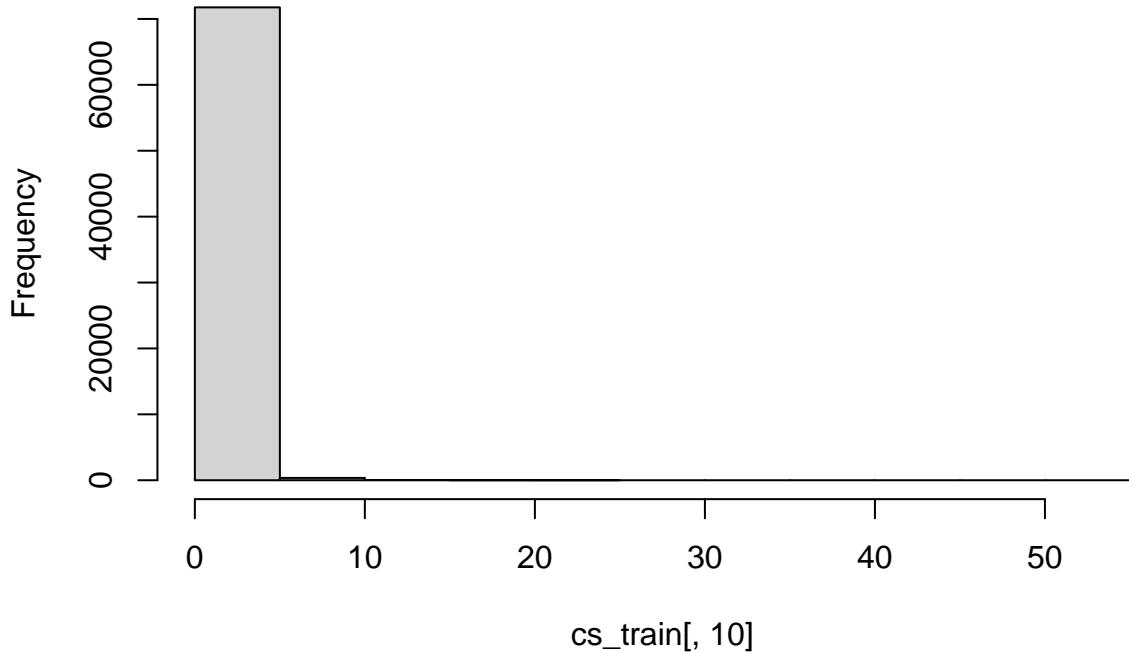
```
hist(cs_train[, 9]) #NumberOfTimes90DaysLate
```

Histogram of cs_train[, 9]



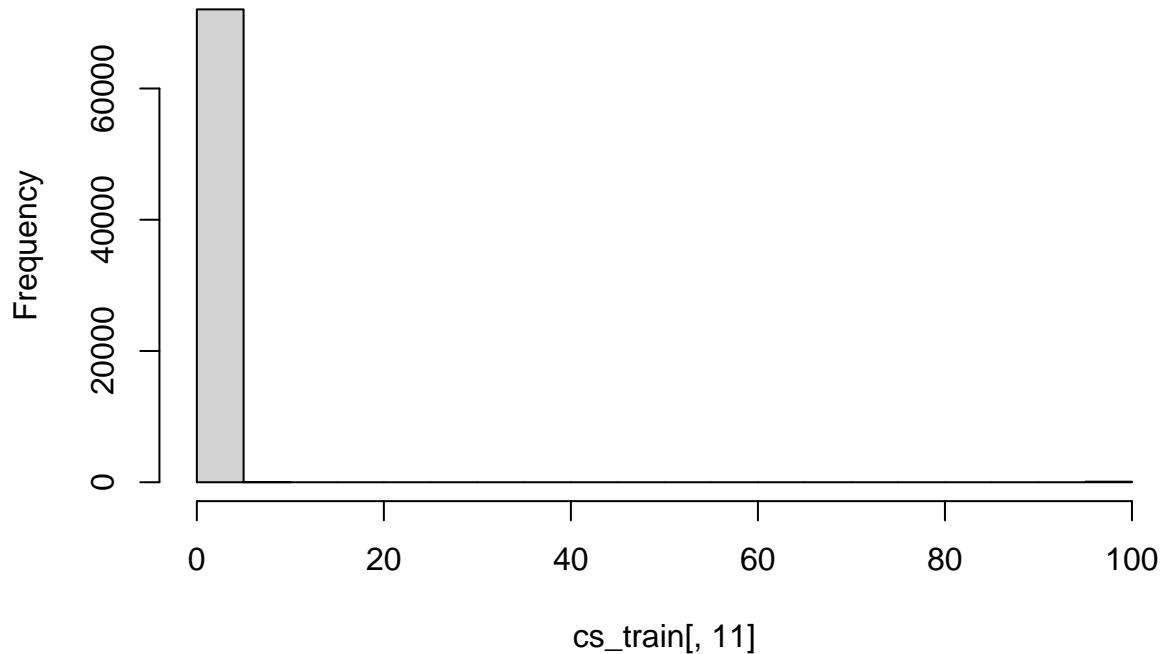
```
hist(cs_train[, 10]) #NumberRealEstateLoansOrLines
```

Histogram of cs_train[, 10]



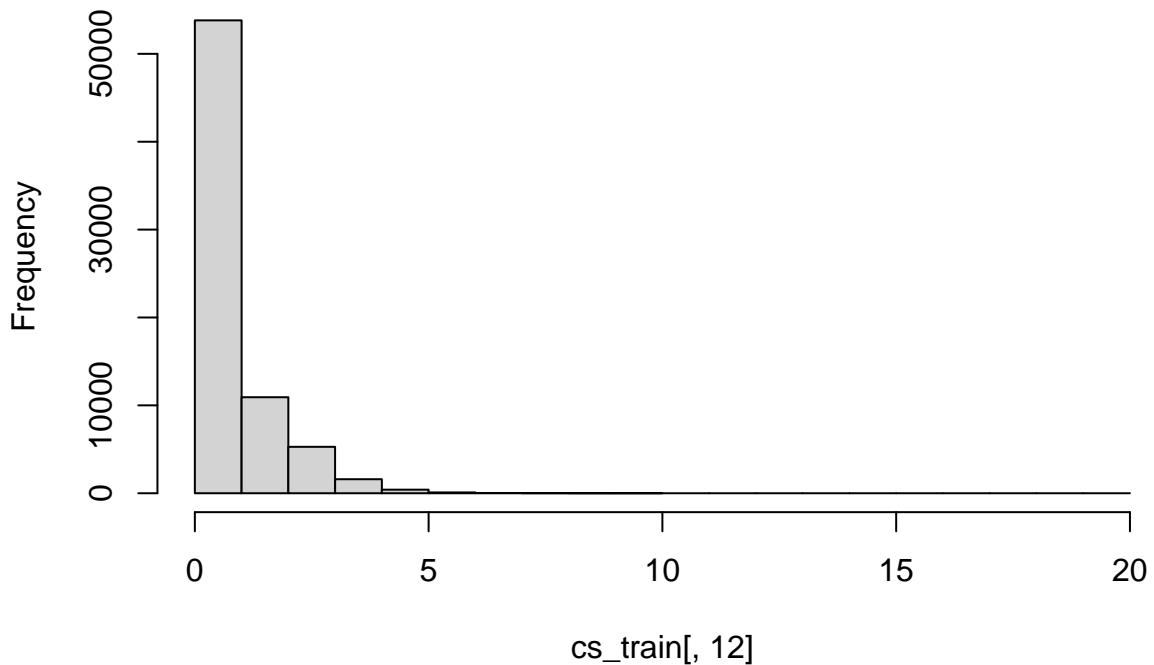
```
hist(cs_train[, 11]) #NumberOfTime60.89DaysPastDueNotWorse
```

Histogram of cs_train[, 11]



```
hist(cs_train[, 12]) #Number of Dependents` ``
```

Histogram of cs_train[, 12]



Creating Indicator Variables for Extreme Data

```
Jimmys_Bad_Variables <- cs_train
Jimmys_Bad_Variables$HighRevolving <- 0
Jimmys_Bad_Variables$HighRevolving[quantile(Jimmys_Bad_Variables[,3],0.95)] <- 1

Jimmys_Bad_Variables$Many_LessThan2MonthsLate <- 0
Jimmys_Bad_Variables$Many_LessThan2MonthsLate[quantile(Jimmys_Bad_Variables[,5],0.95)] <- 1

Jimmys_Bad_Variables$HighDebtRatio <- 0
Jimmys_Bad_Variables$HighDebtRatio[quantile(Jimmys_Bad_Variables[,6],0.95)] <- 1

Jimmys_Bad_Variables$Rich <- 0
Jimmys_Bad_Variables$Rich[quantile(Jimmys_Bad_Variables[,6],0.95)] <- 1

Jimmys_Bad_Variables$Many_CurrentLoans <- 0
Jimmys_Bad_Variables$Many_CurrentLoans[quantile(Jimmys_Bad_Variables[,8],0.95)] <- 1

Jimmys_Bad_Variables$Many_AtLeastThreeMonthsLate <- 0
Jimmys_Bad_Variables$Many_AtLeastThreeMonthsLate[quantile(Jimmys_Bad_Variables[,9],0.95)] <- 1

Jimmys_Bad_Variables$Many_HouseLoans <- 0
Jimmys_Bad_Variables$Many_HouseLoans[quantile(Jimmys_Bad_Variables[,10],0.95)] <- 1
```

```

Jimmys_Bad_Variables$Many_TwoToThreeMonthsLate <- 0
Jimmys_Bad_Variables$Many_TwoToThreeMonthsLate[quantile(Jimmys_Bad_Variables[,11],0.95)] <- 1

Jimmys_Bad_Variables$NA_Dependents_are_Mean <- Jimmys_Bad_Variables[,12] #Need to Change NA values to M
Jimmys_Bad_Variables$NA_Dependents_are_Mean[is.na(Jimmys_Bad_Variables[,12])] <- mean(Jimmys_Bad_Variables[,12])
Jimmys_Bad_Variables$Many_Dependents <- 0
Jimmys_Bad_Variables$Many_Dependents[quantile(Jimmys_Bad_Variables$NA_Dependents_are_Mean,0.95)] <- 1

train1 <- Jimmys_Bad_Variables[Jimmys_Bad_Variables[,3]< quantile(Jimmys_Bad_Variables[,3],0.95),]

head(Jimmys_Bad_Variables)

##          X SeriousDlqin2yrs RevolvingUtilizationOfUnsecuredLines age
## 30484    30484           0                      0.09474044 47
## 74216    74216           0                      0.05277307 68
## 54077    54077           0                      0.73665267 36
## 86784    86784           0                      0.02365700 36
## 14388   14388           0                      0.00581900 44
## 31442   31442           0                      0.11523783 40
##          NumberOfTime30.59DaysPastDueNotWorse DebtRatio MonthlyIncome
## 30484                   0 0.36842105            6250
## 74216                   0 0.22759781           11884
## 54077                   0 0.09445277            2000
## 86784                   0 0.21710409            5600
## 14388                   0 0.18525779            5100
## 31442                   1 2.49168646            2946
##          NumberOfOpenCreditLinesAndLoans NumberOfTimes90DaysLate
## 30484                  8                      0
## 74216                 13                     0
## 54077                  6                      0
## 86784                  9                      0
## 14388                 24                     0
## 31442                 17                     0
##          NumberRealEstateLoansOrLines NumberOfTime60.89DaysPastDueNotWorse
## 30484                  1                      0
## 74216                  2                      0
## 54077                  0                      0
## 86784                  1                      0
## 14388                  1                      0
## 31442                  3                      0
##          NumberOfDependents HighRevolving Many_LessThan2MonthsLate HighDebtRatio
## 30484                  3                      0                      0                      1
## 74216                  1                      0                      1                      0
## 54077                  1                      0                      0                      0
## 86784                  0                      0                      0                      0
## 14388                  2                      0                      0                      0
## 31442                  0                      0                      0                      0
##          Rich Many_CurrentLoans Many_AtLeastThreeMonthsLate Many_HouseLoans
## 30484     1              0                      1                      0
## 74216     0              0                      0                      0
## 54077     0              0                      0                      1
## 86784     0              0                      0                      0
## 14388     0              0                      0                      0

```

```

## 31442      0          0          0          0
##      Many_TwoToThreeMonthsLate NA_Dependents_are_Mean Many_Dependents
## 30484           1           3           0
## 74216           0           1           0
## 54077           0           1           1
## 86784           0           0           0
## 14388           0           2           0
## 31442           0           0           0

```

```
Just_Bad_Variables <- subset.data.frame(Jimmys_Bad_Variables, select = c(X, SeriousDlqin2yrs, age, HighR
```

1. Question/Goal

- Comparing our model performance against paper's model
- Most important factors

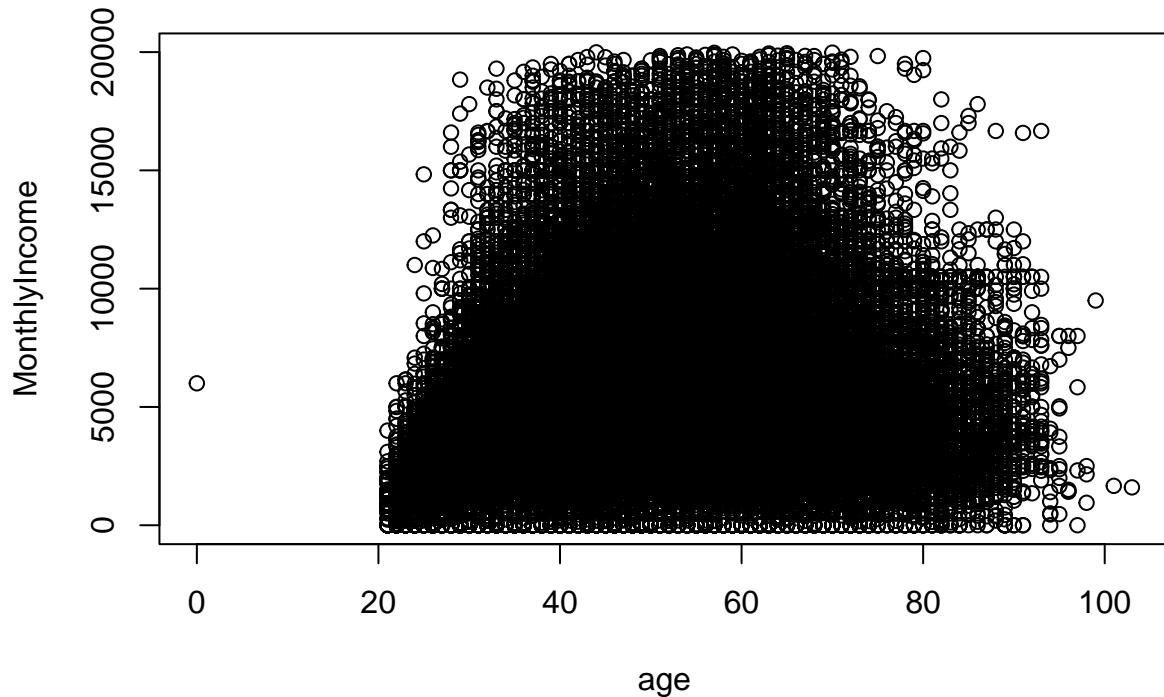
2. EDA (Still need to Clean up for better graphs)

- monthly income = 0
- 30k monthly income = NA
- Dependent = NA 4k
- age 0 remove - only 1 point
- age very old
- group by age group etc
- 13 - 101 or older (maybe cut at 100)
- 80 or more

```

#plot(SeriousDlqin2yrs ~ age, data = cs_train_maj)
cs_train_maj_g1 <- cs_train_maj[cs_train_maj$MonthlyIncome<500000,] # less than 500k monthly income
cs_train_maj_g1 <- cs_train_maj[cs_train_maj$MonthlyIncome<20000,] # less than 20k monthly income
plot(MonthlyIncome ~ age, data = cs_train_maj_g1)

```



```

train <- cs_train
##Shelly EDA Code
tapply(cs_train$RevolvingUtilizationOfUnsecuredLines, train$SeriousDlqin2yrs, median) #use this; data v

##          0          1
## 0.1544904 0.8066074

##Next Steps
# ***need final data set which excludes outliers for which we will create final model from
# do box plots for age group and income; fix axes so that box is readable -Shelly
# histogram of groupings like income by age groups, different colors -Yolanda
# sampling of data --> randomly select instead of adding additional data for response variable -Shelly
# ClassDiscovery package and DataExplorer --> provide initial graphs and analyses of data; also initial

#Remove outliers; maybe only keep anything greater than 10
tapply(train$RevolvingUtilizationOfUnsecuredLines, train$SeriousDlqin2yrs, mean) #d

##          0          1
## 6.566204 3.184407

tapply(train$age, train$SeriousDlqin2yrs, median) #median age of 45 for people defaulting

```

```

## 0 1
## 51 46

#tapply(train$`NumberOfTime30-59DaysPastDueNotWorse`, train$SeriousDlqin2yrs, mean) #2.4 times for those with delinquency
tapply(train$Number0fTimes90DaysLate, train$SeriousDlqin2yrs, mean) #2.1 times for those with delinquency

##          0          1
## 0.111931 1.523173

tapply(train$DebtRatio, train$SeriousDlqin2yrs, median) #0.43 for delinquency incidence vs 0.36 for non-delinquency

##          0          1
## 0.2924788 0.3600771

#tapply(train$NumIncome, train$SeriousDlqin2yrs, median) #delinquency monthly income of $3.8K vs $4.4K
tapply(train$Number0fTimes90DaysLate, train$SeriousDlqin2yrs, mean)

##          0          1
## 0.111931 1.523173

#summary(train$NumIncome)

```

3. Balancing Data (Creating New Datasets)

- Use clustering and pick one sub-group from majority
 - can try Age/Income
 - try Income/debt ratio
- Use bagging algorithm to create more minority data
- Use NA indicator 0 or 1 (maybe use mean/median)

a) Bootstrapping Minority Data (cs_train_min_add) (1000 obs)

- Currently created 1000 additional data, can add more

```

# Bootstrapping Minority Data
set.seed(1) # for reproducibility
# set number of minority data to reproduce
n_add <- 1000

n_min <- nrow(cs_train_min)
n_min

## [1] 4941

```

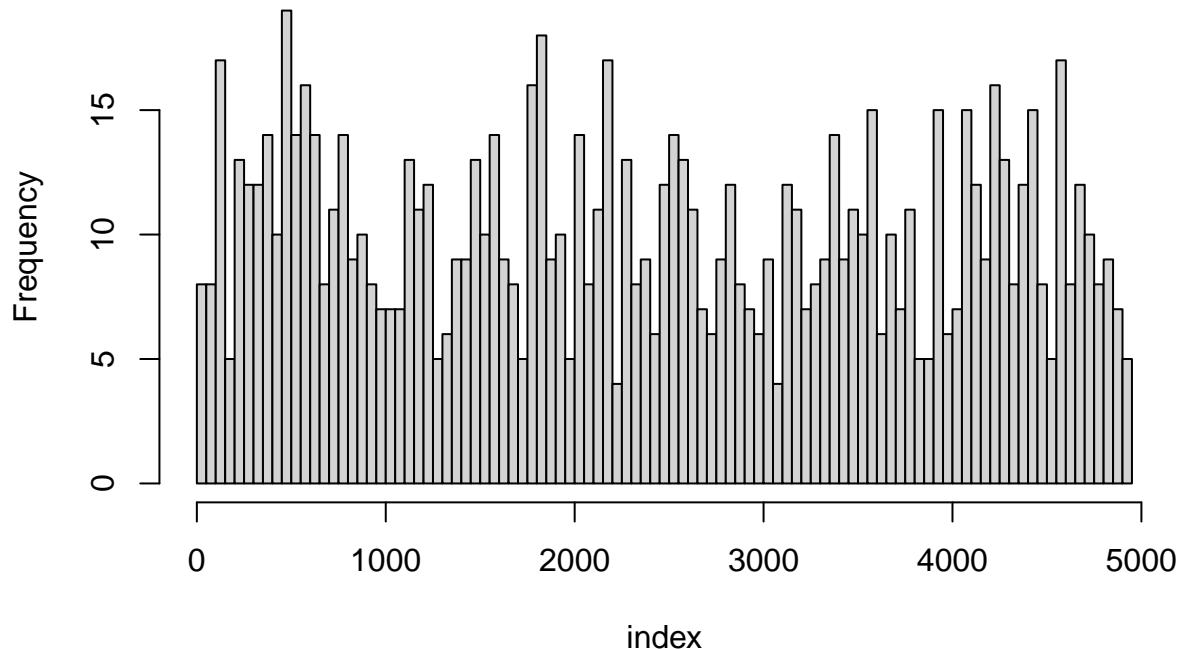
```

index <- sample(n_min, n_add, replace = TRUE)

#plot(density(index), main="") # show density curve of the index we randomized
hist(index, breaks = 100)

```

Histogram of index



```

min(index)

## [1] 15

max(index)

## [1] 4939

length(index)

## [1] 1000

# We add the additional data for future analysis
cs_train_min_add <- cs_train_min[index,]
head(cs_train_min_add)

##          X SeriousDlqin2yrs RevolvingUtilizationOfUnsecuredLines age

```

```

## 88154    88154          1          0.03545848 33
## 64256    64256          1          0.98335183 37
## 105760   105760         1          1.46506986 37
## 137307   137307         1          0.11727646 34
## 20062    20062          1          0.67475736 53
## 47668    47668          1          0.99989018 41
##           NumberOfTime30.59DaysPastDueNotWorse  DebtRatio MonthlyIncome
## 88154                  3 0.29635182      2000
## 64256                  2 0.32223968      7000
## 105760                 0 0.21990439      2300
## 137307                 0 0.08461026      1500
## 20062                  0 2.16232772      3264
## 47668                  0 0.07576094      4500
##           NumberOfOpenCreditLinesAndLoans  NumberOfTimes90DaysLate
## 88154                  10          0
## 64256                   5          0
## 105760                  5          3
## 137307                   9          0
## 20062                   21          0
## 47668                   4          0
##           NumberOfRealEstateLoansOrLines  NumberOfTime60.89DaysPastDueNotWorse
## 88154                  0          1
## 64256                  3          0
## 105760                 0          2
## 137307                 0          0
## 20062                  4          0
## 47668                  0          0
##           NumberOfDependents
## 88154                  0
## 64256                  1
## 105760                  2
## 137307                  4
## 20062                  0
## 47668                  0

```

b) Resampling Majority Data (new_train.final 1 & 2)

- Method 1: use Full omit Data to sample majority data (new_train.final1) (8357*2)
- Method 2: use Training Data to sample majority data (new_train.final2) (4941*2)

```

####SHELLY - Balance Response Variable
set.seed(2) # for reproducibility

## Method 1: Balance Full omit Data (new_train.final1)
new_train <- filter(cs_train_omit, cs_train_omit$SeriousDlqin2yrs == 0) #112k rows

#Check Response Variable Balance in cs_train_omit data
table(cs_train_omit$SeriousDlqin2yrs) #8,357 1's

```

```

##
##      0      1
## 111912  8357

```

```

n_add <- sum(cs_train_omit$SeriousDlqin2yrs == 1) # Number of Default
n_add

## [1] 8357

new_train2 <- new_train[sample(1:nrow(new_train)),] # 112k rows without replacement
new_train3 <- new_train2[1:n_add,]

#####***Merge Data Together (use new_train5 variable)
new_train4 <- filter(cs_train_omit, cs_train_omit$SeriousDlqin2yrs == 1)
#nrow(new_train4)
new_train.final1 <- rbind(new_train3, new_train4)
nrow(new_train.final1) # 16,714 rows as there should be (8357*2)

## [1] 16714

table(new_train.final1$SeriousDlqin2yrs) #correct; equal # of 0's and 1's

## 
##      0      1
## 8357 8357

## Method 2: Balance Training Data (new_train.final2)

#Check Response Variable Balance in cs_train_omit data
table(cs_train$SeriousDlqin2yrs) # 0: 67220 1: 4941

## 
##      0      1
## 67220 4941

n_add <- nrow(cs_train_min) # Number of Default

# Randomly Sample data from Non-Default group (sample = min group obs#)
new_train <- cs_train_maj[sample(1:nrow(cs_train_maj), n_add, replace=FALSE),] # 67220 rows without replacement

#####***Merge Data Together (use new_train.final2 variable)
new_train.final2 <- rbind(new_train, cs_train_min)
nrow(new_train.final2) # 10k rows as there should be (4941*2)

## [1] 9882

table(new_train.final2$SeriousDlqin2yrs) #correct; equal # of 0's and 1's

## 
##      0      1
## 4941 4941

```

c) Using Clustering to select from Majority Group

Select one cluster from majority group and combine with minority group data to form new Balanced Training Data

```
training_dataset <- data.frame(train.x, train.y)
kmeans_model_train <- kmeans(training_dataset, centers = 5) # centers because paper had 5 clusters
kmeans_model_train$centers #what each cluster's average values are for each variable. The most obvious
```

```
##   RevolvingUtilizationOfUnsecuredLines      age
## 1                      0.25563963 51.85714
## 2                      0.30094326 53.69880
## 3                      0.00732813 52.00000
## 4                     14.17070502 53.63628
## 5                     4.13445042 50.60142
##   NumberOfTime30.59DaysPastDueNotWorse  DebtRatio MonthlyIncome
## 1                      0.2857143  0.003406329  580671.714
## 2                      0.2228916  0.076526033  91885.494
## 3                      0.0000000  0.001470045 3008750.000
## 4                      0.2515422  0.296182799 12837.435
## 5                      0.4164810  35.479464445  4525.353
##   NumberOfOpenCreditLinesAndLoans NumberOfTimes90DaysLate
## 1                      9.142857          0.000000000
## 2                     11.415663          0.06626506
## 3                     10.000000          0.000000000
## 4                     10.787612          0.05615007
## 5                     8.178820          0.25217019
##   NumberRealEstateLoansOrLines NumberOfTime60.89DaysPastDueNotWorse
## 1                     1.0000000          0.28571429
## 2                     2.1927711          0.01807229
## 3                     1.0000000          0.000000000
## 4                     1.6671912          0.05470225
## 5                     0.8784692          0.22136860
##   NumberOfDependents      train.y
## 1             1.2857143 0.00000000
## 2             1.1566265 0.06024096
## 3             3.0000000 0.000000000
## 4             1.1644215 0.04418985
## 5             0.7642823 0.07538190
```

#Cluster 3 and 4 are different than Cluster 1 and 2 in terms of the Response Variable.

```
testing_data <- data.frame(test.x,test.y)
kmeans_model_test <- kmeans(testing_data, centers = 5) #testing to see if test data has similar clusters
kmeans_model_test$centers #Clusters are grouped as (1,2), (3,5) and (4) for Response Variables.
```

```
##   RevolvingUtilizationOfUnsecuredLines      age
## 1                      3.3340523 50.35998
## 2                      7.2909033 52.99487
## 3                      0.1501357 60.16667
## 4                     27.6093312 54.01161
## 5                     0.2950418 54.23077
```

```

##   NumberOfTime30.59DaysPastDueNotWorse    DebtRatio MonthlyIncome
## 1                               0.4511763 38.266839034      3910.045
## 2                               0.2605648  0.309173812      9940.086
## 3                               0.3333333  0.002840504     1103530.500
## 4                               0.2585139  0.234410533      25034.900
## 5                               0.3974359  0.077906671      96138.115
##   NumberOfOpenCreditLinesAndLoans NumberOfTimes90DaysLate
## 1                           7.916175          0.29827306
## 2                          10.209088          0.06744563
## 3                          11.666667          0.00000000
## 4                          11.632353          0.03715170
## 5                          10.666667          0.08974359
##   NumberRealEstateLoansOrLines NumberOfTime60.89DaysPastDueNotWorse
## 1                           0.8020238          0.26264883
## 2                           1.4717300          0.06452451
## 3                           1.3333333          0.00000000
## 4                           2.0890093          0.05263158
## 5                           2.3333333          0.08974359
##   NumberOfDependents      test.y
## 1                           0.718358  0.08213362
## 2                           1.074132  0.04985394
## 3                           0.500000  0.00000000
## 4                           1.330495  0.05495356
## 5                           1.243590  0.05128205

```

4. Applying Different Methods (with and without Balanced Data)

- Logistic regression (compare the different link functions)
- look maybe merge Lasso with logistic regression
- RF

Method 1: PCA

1-a) PCA with Unbalanced Data

PCA using all 22 variables

```

pca_including_dummy_variables <- prcomp(na.omit(Jimmys_Bad_Variables[,-c(which(names(Jimmys_Bad_Variables)
summary(pca_including_dummy_variables) #First PC has like 99.66% of variance

```

```

## Importance of components:
##                  PC1       PC2       PC3       PC4       PC5       PC6       PC7
## Standard deviation 1.420e+04 426.3224 2.96e+02 14.49 5.925 5.059 1.586
## Proportion of Variance 9.987e-01 0.0009 4.30e-04 0.00 0.000 0.000 0.000
## Cumulative Proportion 9.987e-01 0.9996 1.00e+00 1.00 1.000 1.000 1.000
##                  PC8       PC9       PC10      PC11      PC12      PC13      PC14
## Standard deviation 1.027 0.565 0.346 0.007445 0.005265 0.003723 0.003722
## Proportion of Variance 0.000 0.000 0.000 0.000000 0.000000 0.000000 0.000000

```

```

## Cumulative Proportion 1.000 1.000 1.000 1.000000 1.000000 1.000000 1.000000
##                               PC15      PC16      PC17      PC18      PC19
## Standard deviation     3.59e-15 1.287e-17 9.127e-18 3.089e-20 3.354e-23
## Proportion of Variance 0.00e+00 0.000e+00 0.000e+00 0.000e+00 0.000e+00
## Cumulative Proportion 1.00e+00 1.000e+00 1.000e+00 1.000e+00 1.000e+00
##                               PC20
## Standard deviation     2.404e-37
## Proportion of Variance 0.000e+00
## Cumulative Proportion 1.000e+00

pca_including_dummy_variables$rotation[,1] #I think it really likes Monthly Income as it's close to 1.

## RevolvingUtilizationOfUnsecuredLines                                age
##                               1.663543e-04                                3.754805e-05
## NumberOfTime30.59DaysPastDueNotWorse          DebtRatio
##                               -2.494353e-06                                -9.026497e-04
##                               MonthlyIncome      NumberOfOpenCreditLinesAndLoans
##                               9.999996e-01                                3.387604e-05
##                               NumberOfTimes90DaysLate    NumberRealEstateLoansOrLines
##                               -2.982336e-06                                1.014049e-05
## NumberOfTime60.89DaysPastDueNotWorse      NumberOfDependents
##                               -2.547678e-06                                5.400872e-06
##                               HighRevolving      Many_LessThan2MonthsLate
##                               0.000000e+00                                3.594719e-10
##                               HighDebtRatio      Rich
##                               -2.774522e-11                                -2.774522e-11
##                               Many_CurrentLoans      Many_AtLeastThreeMonthsLate
##                               3.130335e-09                                -2.774522e-11
##                               Many_HouseLoans      Many_TwoToThreeMonthsLate
##                               -3.198419e-10                                -2.774522e-11
## NA_Dependents_are_Mean      Many_Dependents
##                               5.400872e-06                                -3.198419e-10

glm_pca_including_dummy_variables <- glm(train.y ~ pca_including_dummy_variables$x[,1], family = "binomial")

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

summary(glm_pca_including_dummy_variables)

## 
## Call:
## glm(formula = train.y ~ pca_including_dummy_variables$x[, 1],
##      family = "binomial")
## 
## Deviance Residuals:
##      Min        1Q        Median       3Q        Max 
## -0.4302   -0.3975   -0.3784   -0.3507   5.2302 
## 
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)    
## (Intercept)              -2.638e+00  1.522e-02 -173.29  <2e-16 ***
## pca_including_dummy_variables$x[, 1] -4.536e-05  3.798e-06  -11.94  <2e-16 ***

```

```

## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 36033 on 72160 degrees of freedom
## Residual deviance: 35857 on 72159 degrees of freedom
## AIC: 35861
##
## Number of Fisher Scoring iterations: 6

```

PCA using original 10 variables

```
pca_original_data <- prcomp(train.x)
```

```
summary(pca_original_data) #PC1 has 99.66% variance explained
```

```

## Importance of components:
##                               PC1        PC2        PC3        PC4        PC5        PC6        PC7
## Standard deviation     1.420e+04 426.3224 2.96e+02 14.49 5.925 5.058 1.144
## Proportion of Variance 9.987e-01 0.0009 4.30e-04 0.00 0.000 0.000 0.000
## Cumulative Proportion 9.987e-01 0.9996 1.00e+00 1.00 1.000 1.000 1.000
##                               PC8        PC9        PC10
## Standard deviation     1.007 0.565 0.346
## Proportion of Variance 0.000 0.000 0.000
## Cumulative Proportion 1.000 1.000 1.000

```

```
pca_original_data$rotation[,1] #It also really likes MonthlyIncome and not anything else
```

## RevolvingUtilizationOfUnsecuredLines	age	
##	1.663543e-04	
## NumberOfTime30.59DaysPastDueNotWorse	3.754805e-05	
##	DebtRatio	
##	-2.494353e-06	
##	MonthlyIncome	
##	NumberOpenCreditLinesAndLoans	
##	9.999996e-01	3.387604e-05
##	NumberOfTimes90DaysLate	NumberRealEstateLoansOrLines
##	-2.982336e-06	1.014049e-05
##	NumberOfTime60.89DaysPastDueNotWorse	NumberOfDependents
##	-2.547678e-06	5.400872e-06

```
glm_pca_original_data <- glm(train.y ~ pca_original_data$x[,1], family = "binomial")
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
summary(glm_pca_original_data)
```

```

##
## Call:
## glm(formula = train.y ~ pca_original_data$x[, 1], family = "binomial")

```

```

## 
## Deviance Residuals:
##      Min       1Q   Median      3Q     Max
## -0.4302 -0.3975 -0.3784 -0.3507  5.2302
## 
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)           -2.638e+00  1.522e-02 -173.29 <2e-16 ***
## pca_original_data$x[, 1] -4.536e-05  3.798e-06 -11.94 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## (Dispersion parameter for binomial family taken to be 1)
## 
## Null deviance: 36033  on 72160  degrees of freedom
## Residual deviance: 35857  on 72159  degrees of freedom
## AIC: 35861
## 
## Number of Fisher Scoring iterations: 6

```

PCA Using Dummy variables

```

pca_dummy <- prcomp(Just_Bad_Variables[,-c(which(names(Just_Bad_Variables) == "SeriousDlqin2yrs")), which(names(Just_Bad_Variables) != "SeriousDlqin2yrs")])

summary(pca_dummy) #PC 1 is 100%?

## Importance of components:
##                  PC1        PC2        PC3        PC4        PC5        PC6
## Standard deviation    14.46  0.007445  0.005265  0.003723  0.003723 7.281e-16
## Proportion of Variance 1.00  0.000000  0.000000  0.000000  0.000000 0.000e+00
## Cumulative Proportion 1.00  1.000000  1.000000  1.000000  1.000000 1.000e+00
##                  PC7        PC8        PC9        PC10
## Standard deviation    6.661e-16 1.404e-19 8.404e-23 1.584e-35
## Proportion of Variance 0.000e+00 0.000e+00 0.000e+00 0.000e+00
## Cumulative Proportion 1.000e+00 1.000e+00 1.000e+00 1.000e+00

pca_dummy$rotation[,1] #Age is significant

##                age          HighRevolving
## 1.000000e+00 0.000000e+00
## Many_LessThan2MonthsLate HighDebtRatio
## 1.108186e-06 -2.834080e-07
## Rich            Many_CurrentLoans
## -2.834080e-07 5.117887e-07
## Many_AtLeastThreeMonthsLate Many_HouseLoans
## -2.834080e-07 -1.012338e-06
## Many_TwoToThreeMonthsLate Many_Dependents
## -2.834080e-07 -1.012338e-06

```

```
glm_pca_dummy <- glm(train.y ~ pca_dummy$x[, 1], family = "binomial")
summary(glm_pca_dummy)
```

```
##
## Call:
## glm(formula = train.y ~ pca_dummy$x[, 1], family = "binomial")
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -0.7208 -0.4135 -0.3553 -0.3049  2.8707
##
## Coefficients:
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept) -2.681737  0.015742 -170.36  <2e-16 ***
## pca_dummy$x[, 1] -0.028603  0.001092  -26.19  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 36033  on 72160  degrees of freedom
## Residual deviance: 35308  on 72159  degrees of freedom
## AIC: 35312
##
## Number of Fisher Scoring iterations: 5
```

1 b) PCA with Balanced Data

Method 2: Using GLM Original Predictors vs Extreme Binned Predictors + MonthlyIncome

We compared GLM original Predictors with extreme binned predictors using unbalanced data and see better performance in original numeric predictors.

We now try to improve the model using balanced data.

2 a) GLM Original Predictors with Unbalanced Data

```
## Original Predictors

model <- glm(Jimmys_Bad_Variables$SeriousDlqin2yrs ~ Jimmysts_Bad_Variables[, 3] + Jimmysts_Bad_Variables[, 4]

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

#To Change family link use family = quasi(variance = "mu^3", link = "log") change quasi
summary(model)

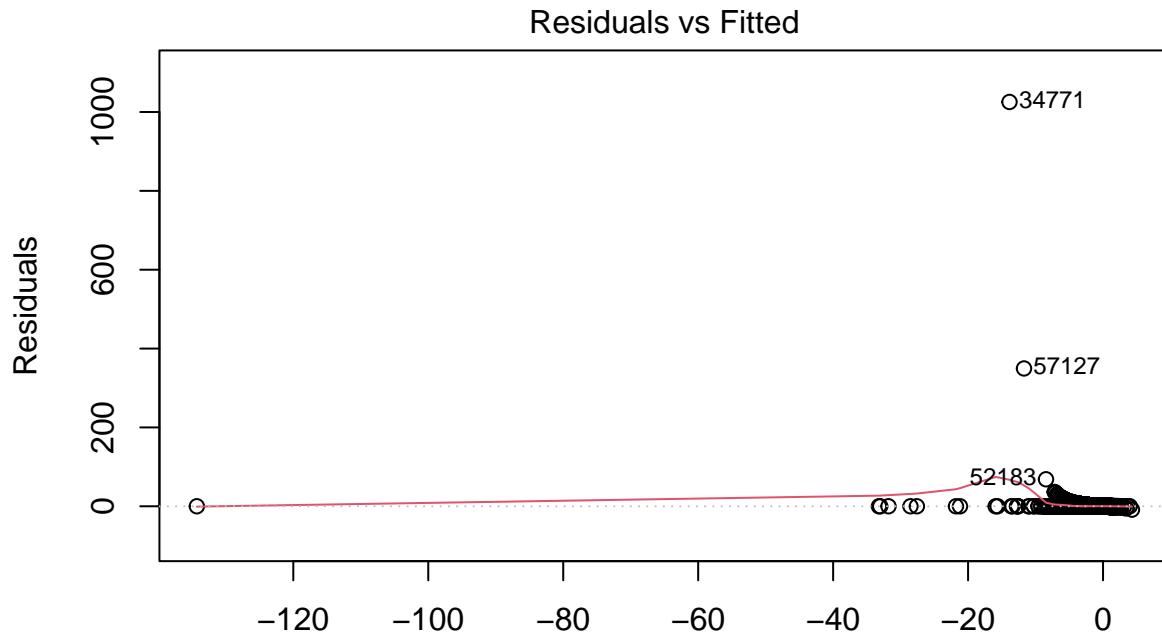
##
## Call:
## glm(formula = Jimmysts_Bad_Variables$SeriousDlqin2yrs ~ Jimmysts_Bad_Variables[, 3] + Jimmysts_Bad_Variables[, 4] + Jimmysts_Bad_Variables[, 5] +
```

```

##      Jimmys_Bad_Variables[, 6] + Jimmys_Bad_Variables[, 7] + Jimmys_Bad_Variables[, 
##      8] + Jimmys_Bad_Variables[, 9] + Jimmys_Bad_Variables[, 10] + 
##      Jimmys_Bad_Variables[, 11] + Jimmys_Bad_Variables[, 12],
##      family = "binomial", data = Jimmys_Bad_Variables)
##
## Deviance Residuals:
##      Min        1Q     Median        3Q       Max
## -2.9318  -0.3956  -0.3264  -0.2654   5.2661
##
## Coefficients:
##                               Estimate Std. Error z value Pr(>|z|)
## (Intercept)           -1.562e+00 6.016e-02 -25.957 < 2e-16 ***
## Jimmys_Bad_Variables[, 3] -8.778e-05 1.123e-04  -0.782   0.434
## Jimmys_Bad_Variables[, 4] -2.411e-02 1.203e-03 -20.045 < 2e-16 ***
## Jimmys_Bad_Variables[, 5]  4.991e-01 1.537e-02  32.464 < 2e-16 ***
## Jimmys_Bad_Variables[, 6] -7.277e-05 5.259e-05 -1.384   0.166
## Jimmys_Bad_Variables[, 7] -4.383e-05 4.296e-06 -10.202 < 2e-16 ***
## Jimmys_Bad_Variables[, 8] -1.905e-03 3.508e-03 -0.543   0.587
## Jimmys_Bad_Variables[, 9]  4.205e-01 2.195e-02  19.154 < 2e-16 ***
## Jimmys_Bad_Variables[, 10] 8.801e-02 1.403e-02   6.272 3.57e-10 ***
## Jimmys_Bad_Variables[, 11] -8.845e-01 2.524e-02 -35.043 < 2e-16 ***
## Jimmys_Bad_Variables[, 12]  1.061e-01 1.250e-02   8.493 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 36033  on 72160  degrees of freedom
## Residual deviance: 33374  on 72150  degrees of freedom
## AIC: 33396
##
## Number of Fisher Scoring iterations: 6

plot(model, which = 1) #Outliers skew residuals plot

```



```
glm(Jimmys_Bad_Variables$SeriousDlqin2yrs ~ Jimmys_Bad_Variables[, 3] + Jim ..
```

```
stepAIC(model)
```

```
## Start: AIC=33396.05
## Jimmys_Bad_Variables$SeriousDlqin2yrs ~ Jimmys_Bad_Variables[, 3] + Jimmys_Bad_Variables[, 4] + Jimmys_Bad_Variables[, 5] + Jimmys_Bad_Variables[, 6] + Jimmys_Bad_Variables[, 7] + Jimmys_Bad_Variables[, 8] + Jimmys_Bad_Variables[, 9] + Jimmys_Bad_Variables[, 10] + Jimmys_Bad_Variables[, 11] + Jimmys_Bad_Variables[, 12]

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```

##                                     Df Deviance   AIC
## - Jimmys_Bad_Variables[, 8]    1   33374 33394
## - Jimmys_Bad_Variables[, 3]    1   33375 33395
## <none>                           33374 33396
## - Jimmys_Bad_Variables[, 6]    1   33376 33396
## - Jimmys_Bad_Variables[, 10]   1   33412 33432
## - Jimmys_Bad_Variables[, 12]   1   33444 33464
## - Jimmys_Bad_Variables[, 7]    1   33504 33524
## - Jimmys_Bad_Variables[, 9]    1   33738 33758
## - Jimmys_Bad_Variables[, 4]    1   33796 33816
## - Jimmys_Bad_Variables[, 5]    1   34316 34336
## - Jimmys_Bad_Variables[, 11]   1   34539 34559

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

## 
## Step:  AIC=33394.35
## Jimmys_Bad_Variables$SeriousDlqin2yrs ~ Jimmys_Bad_Variables[, 3] + Jimmys_Bad_Variables[, 4] + Jimmys_Bad_Variables[, 5] +
## Jimmys_Bad_Variables[, 6] + Jimmys_Bad_Variables[, 7] + Jimmys_Bad_Variables[, 9] + Jimmys_Bad_Variables[, 10] + Jimmys_Bad_Variables[, 11] + Jimmys_Bad_Variables[, 12]

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##                                     Df Deviance   AIC
## - Jimmys_Bad_Variables[, 3]    1   33375 33393
## <none>                           33374 33394
## - Jimmys_Bad_Variables[, 6]    1   33377 33395
## - Jimmys_Bad_Variables[, 10]   1   33415 33433
## - Jimmys_Bad_Variables[, 12]   1   33444 33462
## - Jimmys_Bad_Variables[, 7]    1   33508 33526
## - Jimmys_Bad_Variables[, 9]    1   33744 33762
## - Jimmys_Bad_Variables[, 4]    1   33816 33834
## - Jimmys_Bad_Variables[, 5]    1   34323 34341
## - Jimmys_Bad_Variables[, 11]   1   34540 34558

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```

```

## 
## Step: AIC=33393.21
## Jimmys_Bad_Variables$SeriousDlqin2yrs ~ Jimmys_Bad_Variables[, 
##   4] + Jimmys_Bad_Variables[, 5] + Jimmys_Bad_Variables[, 6] + 
##   Jimmys_Bad_Variables[, 7] + Jimmys_Bad_Variables[, 9] + Jimmys_Bad_Variables[, 
##   10] + Jimmys_Bad_Variables[, 11] + Jimmys_Bad_Variables[, 
##   12]

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

##                                     Df Deviance    AIC
## <none>                               33375 33393
## - Jimmys_Bad_Variables[, 6]     1   33378 33394
## - Jimmys_Bad_Variables[, 10]    1   33416 33432
## - Jimmys_Bad_Variables[, 12]    1   33445 33461
## - Jimmys_Bad_Variables[, 7]     1   33509 33525
## - Jimmys_Bad_Variables[, 9]     1   33745 33761
## - Jimmys_Bad_Variables[, 4]     1   33816 33832
## - Jimmys_Bad_Variables[, 5]     1   34324 34340
## - Jimmys_Bad_Variables[, 11]    1   34541 34557

## 
## Call: glm(formula = Jimmys_Bad_Variables$SeriousDlqin2yrs ~ Jimmys_Bad_Variables[, 
##   4] + Jimmys_Bad_Variables[, 5] + Jimmys_Bad_Variables[, 6] + 
##   Jimmys_Bad_Variables[, 7] + Jimmys_Bad_Variables[, 9] + Jimmys_Bad_Variables[, 
##   10] + Jimmys_Bad_Variables[, 11] + Jimmys_Bad_Variables[, 
##   12], family = "binomial", data = Jimmys_Bad_Variables)
## 

## Coefficients:
##             (Intercept) Jimmys_Bad_Variables[, 4]
##                   -1.567e+00           -2.422e-02
## Jimmys_Bad_Variables[, 5] Jimmys_Bad_Variables[, 6]
##                   4.981e-01           -7.371e-05
## Jimmys_Bad_Variables[, 7] Jimmys_Bad_Variables[, 9]
##                   -4.418e-05          4.219e-01
## Jimmys_Bad_Variables[, 10] Jimmys_Bad_Variables[, 11]
##                   8.517e-02           -8.849e-01
## Jimmys_Bad_Variables[, 12]
##                   1.061e-01

## 
## Degrees of Freedom: 72160 Total (i.e. Null); 72152 Residual

```

```

## Null Deviance:      36030
## Residual Deviance: 33380      AIC: 33390

```

2 b) GLM Extreme Binned Predictors with Unbalanced Data

```

model2 <- glm(Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmysts_Bad_Variables[,13] + Jimmysts_Bad_Variables[,14] + Jimmysts_Bad_Variables[,15] + Jimmysts_Bad_Variables[,16] + Jimmysts_Bad_Variables[,17] + Jimmysts_Bad_Variables[,18] + Jimmysts_Bad_Variables[,19] + Jimmysts_Bad_Variables[,20] + Jimmysts_Bad_Variables[,22], family = "binomial", data = Jimmysts_Bad_Variables)

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

model2

##
## Call: glm(formula = Jimmysts_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome +
##           Jimmysts_Bad_Variables[, 13] + Jimmysts_Bad_Variables[, 14] +
##           Jimmysts_Bad_Variables[, 15] + Jimmysts_Bad_Variables[, 16] +
##           Jimmysts_Bad_Variables[, 17] + Jimmysts_Bad_Variables[, 18] +
##           Jimmysts_Bad_Variables[, 19] + Jimmysts_Bad_Variables[, 20] +
##           Jimmysts_Bad_Variables[, 22], family = "binomial", data = Jimmysts_Bad_Variables)
##
## Coefficients:
## (Intercept)          MonthlyIncome
## -2.336e+00          -4.536e-05
## Jimmysts_Bad_Variables[, 13] Jimmysts_Bad_Variables[, 14]
## NA                  -7.691e+00
## Jimmysts_Bad_Variables[, 15] Jimmysts_Bad_Variables[, 16]
## -7.946e+00          NA
## Jimmysts_Bad_Variables[, 17] Jimmysts_Bad_Variables[, 18]
## -5.862e+00          NA
## Jimmysts_Bad_Variables[, 19] Jimmysts_Bad_Variables[, 20]
## -8.139e+00          NA
## Jimmysts_Bad_Variables[, 22]
## NA
##
## Degrees of Freedom: 72160 Total (i.e. Null); 72155 Residual
## Null Deviance:      36030
## Residual Deviance: 35860      AIC: 35870

model2$rank # equals 7 which is how many variables are not NA

## [1] 6

#Certain Dummy Variables are a Singular Matrix Meaning that some of our variables can be constructed using
#Not Sure what to do, but I'll remove these variables in the meantime

model3 <- glm(Jimmysts_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmysts_Bad_Variables[,15] + Jimmysts_Bad_Variables[,16] + Jimmysts_Bad_Variables[,17] + Jimmysts_Bad_Variables[,18] + Jimmysts_Bad_Variables[,19] + Jimmysts_Bad_Variables[,20])

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred

```

```
model3
```

```
##  
## Call: glm(formula = Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome +  
##          Jimmys_Bad_Variables[, 15] + Jimmys_Bad_Variables[, 16] +  
##          Jimmys_Bad_Variables[, 17] + Jimmys_Bad_Variables[, 18] +  
##          Jimmys_Bad_Variables[, 19], family = "binomial", data = Jimmys_Bad_Variables)  
##  
## Coefficients:  
##              (Intercept)          MonthlyIncome  
##                  -2.336e+00           -4.536e-05  
##  Jimmys_Bad_Variables[, 15]  Jimmys_Bad_Variables[, 16]  
##                  -7.946e+00             NA  
##  Jimmys_Bad_Variables[, 17]  Jimmys_Bad_Variables[, 18]  
##                  -5.862e+00             NA  
##  Jimmys_Bad_Variables[, 19]  
##                  -8.139e+00  
##  
## Degrees of Freedom: 72160 Total (i.e. Null); 72156 Residual  
## Null Deviance: 36030  
## Residual Deviance: 35860 AIC: 35870
```

```
model4 <- glm(Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmys_Bad_Variables[,15] + Jimmy
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
model4
```

```
##  
## Call: glm(formula = Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome +  
##          Jimmys_Bad_Variables[, 15] + Jimmys_Bad_Variables[, 16] +  
##          Jimmys_Bad_Variables[, 17] + Jimmys_Bad_Variables[, 18] +  
##          Jimmys_Bad_Variables[, 19], family = binomial(link = "probit"),  
##          data = Jimmys_Bad_Variables)  
##  
## Coefficients:  
##              (Intercept)          MonthlyIncome  
##                  -1.385e+00           -1.658e-05  
##  Jimmys_Bad_Variables[, 15]  Jimmys_Bad_Variables[, 16]  
##                  -2.677e+00             NA  
##  Jimmys_Bad_Variables[, 17]  Jimmys_Bad_Variables[, 18]  
##                  -1.915e+00             NA  
##  Jimmys_Bad_Variables[, 19]  
##                  -2.748e+00  
##  
## Degrees of Freedom: 72160 Total (i.e. Null); 72156 Residual  
## Null Deviance: 36030  
## Residual Deviance: 35890 AIC: 35900
```

```
model5 <- glm(Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmys_Bad_Variables[,15] + Jimmy
```

```
## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```

```
model5
```

```
##  
## Call: glm(formula = Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome +  
## Jimmys_Bad_Variables[, 15] + Jimmys_Bad_Variables[, 16] +  
## Jimmys_Bad_Variables[, 17] + Jimmys_Bad_Variables[, 18] +  
## Jimmys_Bad_Variables[, 19], family = binomial(link = "cloglog"),  
## data = Jimmys_Bad_Variables)  
##  
## Coefficients:  
## (Intercept) MonthlyIncome  
## -2.3775977 -0.0000445  
## Jimmys_Bad_Variables[, 15] Jimmys_Bad_Variables[, 16]  
## -7.8205053 NA  
## Jimmys_Bad_Variables[, 17] Jimmys_Bad_Variables[, 18]  
## -5.7758409 NA  
## Jimmys_Bad_Variables[, 19]  
## -8.0096201  
##  
## Degrees of Freedom: 72160 Total (i.e. Null); 72156 Residual  
## Null Deviance: 36030  
## Residual Deviance: 35850 AIC: 35860
```

```
anova(model, model2, model3, model4, model5) #model 2 is the one with NA values
```

```
## Analysis of Deviance Table  
##  
## Model 1: Jimmys_Bad_Variables$SeriousDlqin2yrs ~ Jimmys_Bad_Variables[,  
## 3] + Jimmys_Bad_Variables[, 4] + Jimmys_Bad_Variables[, 5] +  
## Jimmys_Bad_Variables[, 6] + Jimmys_Bad_Variables[, 7] + Jimmys_Bad_Variables[,  
## 8] + Jimmys_Bad_Variables[, 9] + Jimmys_Bad_Variables[, 10] +  
## Jimmys_Bad_Variables[, 11] + Jimmys_Bad_Variables[, 12]  
## Model 2: Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmys_Bad_Variables[,  
## 13] + Jimmys_Bad_Variables[, 14] + Jimmys_Bad_Variables[,  
## 15] + Jimmys_Bad_Variables[, 16] + Jimmys_Bad_Variables[,  
## 17] + Jimmys_Bad_Variables[, 18] + Jimmys_Bad_Variables[,  
## 19] + Jimmys_Bad_Variables[, 20] + Jimmys_Bad_Variables[,  
## 22]  
## Model 3: Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmys_Bad_Variables[,  
## 15] + Jimmys_Bad_Variables[, 16] + Jimmys_Bad_Variables[,  
## 17] + Jimmys_Bad_Variables[, 18] + Jimmys_Bad_Variables[,  
## 19]  
## Model 4: Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmys_Bad_Variables[,  
## 15] + Jimmys_Bad_Variables[, 16] + Jimmys_Bad_Variables[,  
## 17] + Jimmys_Bad_Variables[, 18] + Jimmys_Bad_Variables[,  
## 19]  
## Model 5: Jimmys_Bad_Variables$SeriousDlqin2yrs ~ MonthlyIncome + Jimmys_Bad_Variables[,  
## 15] + Jimmys_Bad_Variables[, 16] + Jimmys_Bad_Variables[,  
## 17] + Jimmys_Bad_Variables[, 18] + Jimmys_Bad_Variables[,  
## 19]  
## Resid. Df Resid. Dev Df Deviance  
## 1 72150 33374  
## 2 72155 35856 -5 -2482.44
```

```

## 3      72156    35857 -1     -0.11
## 4      72156    35889  0     -32.45
## 5      72156    35854  0      35.02

```

```
#model 1 is the best but we can check model 5 using our dummy variables
```

2 c) GLM Original Predictors with Balanced Data

Method 3: Ridge

3 a) Ridge with Unbalanced Data

```

x <- data.frame(Jimmys_Bad_Variables$SeriousDlqin2yrs, Jimmysts_Bad_Variables$MonthlyIncome, Jimmysts_Bad_V
x <- na.omit(x) #Remember Monthly Income contains NA values
x_for_ridge <- data.matrix(x[,-1]) #Everything but Response Variable
ridge_model15 <- cv.glmnet(x_for_ridge,x[,1], alpha = 0, standardize = TRUE, nfolds = length(x))
ridge_model15

```

```

##
## Call: cv.glmnet(x = x_for_ridge, y = x[, 1], nfolds = length(x), alpha = 0,           standardize = TRUE)
##
## Measure: Mean-Squared Error
##
##      Lambda Index Measure          SE Nonzero
## min   0.845      20 0.06378 0.001039       6
## 1se   4.951      1 0.06379 0.001041       6

```

3 b) Ridge with Balanced Data

Method 4: Lasso

4 a) Lasso with Unbalanced Data

```

lasso_model15 <- cv.glmnet(x_for_ridge,x[,1], alpha = 1, standardize = TRUE, nfolds = length(x))
lasso_model15

```

```

##
## Call: cv.glmnet(x = x_for_ridge, y = x[, 1], nfolds = length(x), alpha = 1,           standardize = TRUE)
##
## Measure: Mean-Squared Error
##
##      Lambda Index Measure          SE Nonzero
## min 0.004951      1 0.0638 0.0005664       0
## 1se 0.004951      1 0.0638 0.0005664       0

```

4 b) Lasso with Balanced Data

Method 4: Random Forest

4 a) Random Forest with Unbalanced Data

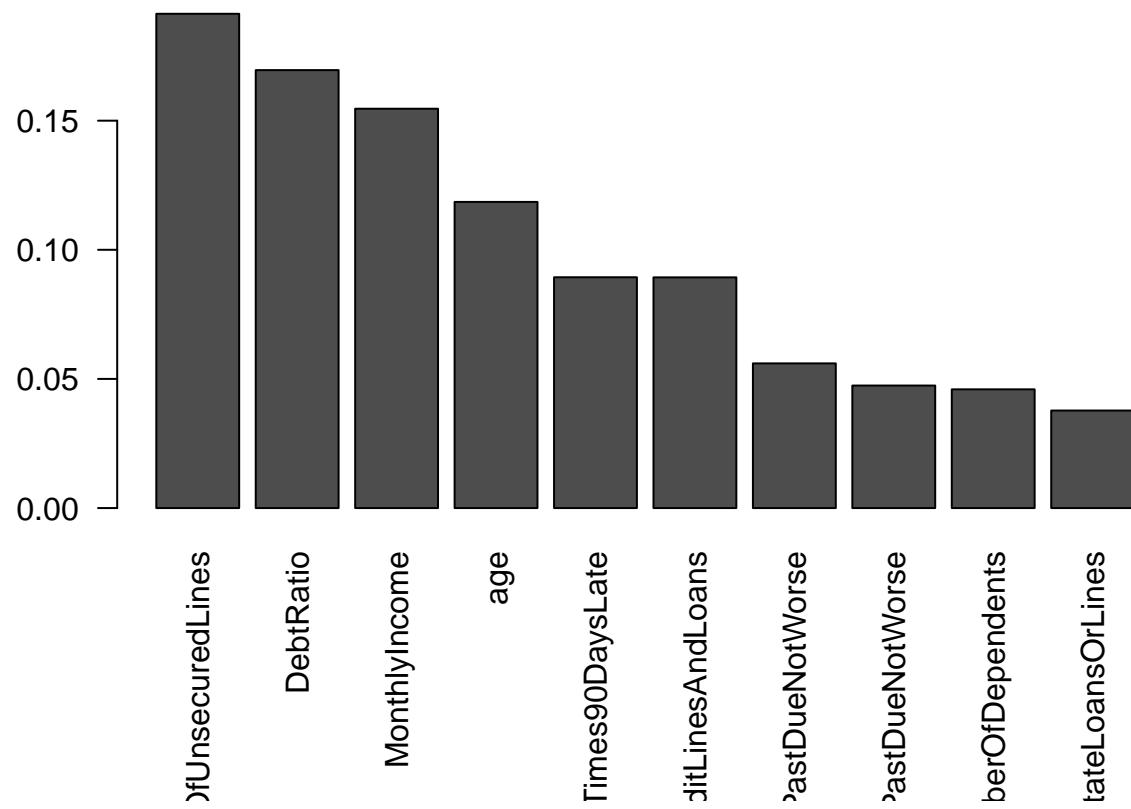
```
# Fits Random Forest
model.rf = randomForest(y = as.factor(train.y), x=train.x, xtest=test.x, ytest=as.factor(test.y), mtry = 3)

model.rf # Output shows confusion matrix for both train and test

##
## Call:
##   randomForest(x = train.x, y = as.factor(train.y), xtest = test.x,      ytest = as.factor(test.y), mtry = 3)
##   Type of random forest: classification
##   Number of trees: 500
##   No. of variables tried at each split: 3
##
##       OOB estimate of  error rate: 6.64%
## Confusion matrix:
##     0 1 class.error
## 0 66591 629 0.009357334
## 1 4166 775 0.843149160
##             Test set error rate: 6.7%
## Confusion matrix:
##     0 1 class.error
## 0 44307 385 0.008614517
## 1 2840 576 0.831381733

# Random Forest Output
var.imp = data.frame(importance(model.rf, type=2))
var.imp$Variables = row.names(var.imp)
varimp = var.imp[order(var.imp$MeanDecreaseGini, decreasing = T),]
par(mar = c(7.5,3,2,2))
giniplot = barplot(t(varimp[-2]/sum(varimp[-2])), las=2, cex.names=1, main="Gini Impurity Index Plot")
```

Gini Impurity Index Plot



4 b) Random Forest with Balanced Data

5. Evaluating Model/Comparing results

- mmp plots to see if model fits the data, as well as the pearson Chi-square test
- checking for outliers and influential points
- Some measure of pesudo R-square and accuracy of the model
- Use confusion matrix, ROC/AUC curve, AIC to evaluate the different models

5-1 Evaluation on Final Model using Training Data

```
## Set model as final model
model.final <- model

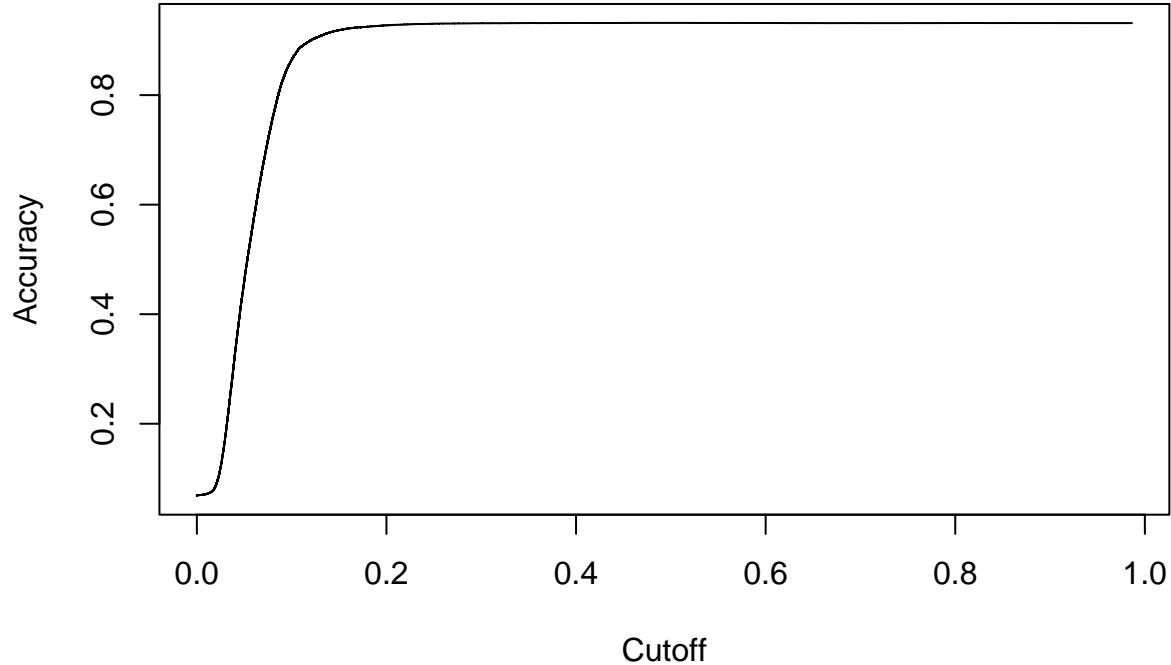
## Evaluation on Final model using Training Data
train_preds = predict(model.final, newdata=train.x, type="response")
head(train_preds[is.na(train_preds)])
```



```
## named numeric(0)

#train.x[c(7,9),]
#train_preds[is.na(train_preds)]
```

```
pred_compare = prediction(train_preds, train.y)
plot(performance(pred_compare, "acc"))
```



```
table(train.y, train_preds>0.2) # accuracy on train
```

```
##  
## train.y FALSE TRUE  
##      0 66148 1072  
##      1  4161   780
```

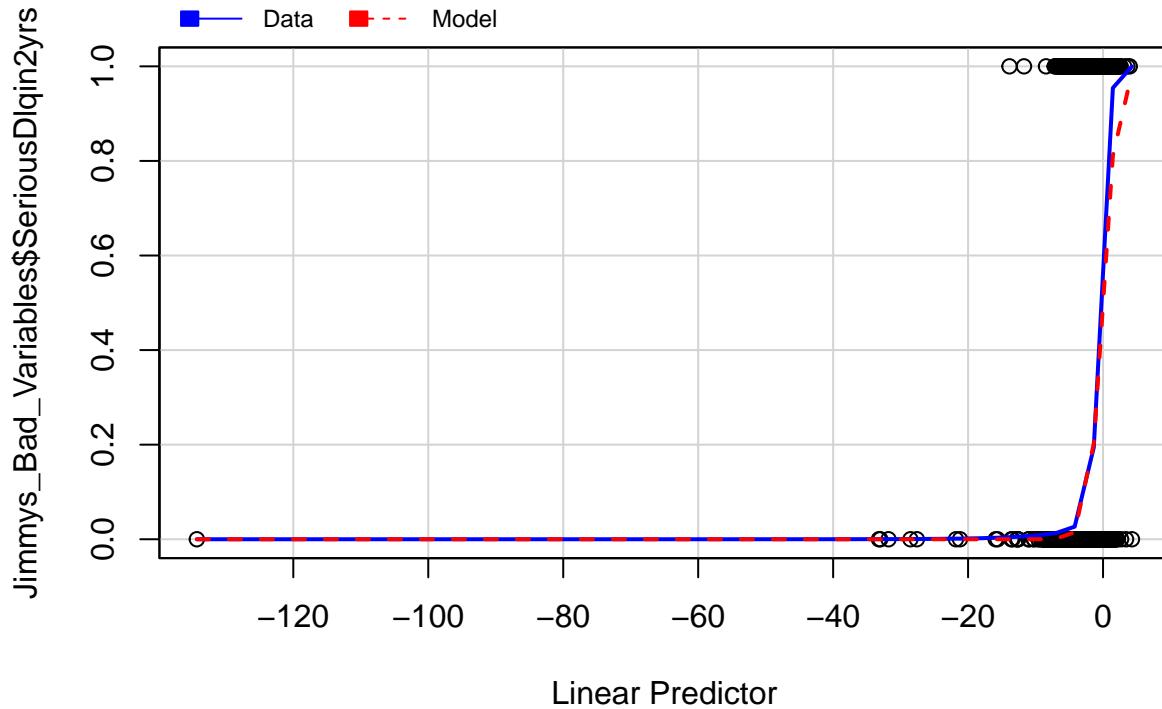
5-2 Marginal Model Plots

Need to Update!!

```
# residual plots
#residualPlots(model.final)

# Marginal Model Plots
library(car)
mmp(model.final)

## Warning: glm.fit: fitted probabilities numerically 0 or 1 occurred
```



5-3 Goodness of Fit using Hosmer-Lemeshow Test

The p-value is 0, meaning that we want to reject the null hypothesis that the model is adequate.

```
# Goodness of Fit using Hosmer-Lemeshow Test
linpred=predict(model.final)

cs_train_m <- mutate(cs_train, predprob=predict(model.final, type="response")) # cal p^_i
gdf <- group_by(cs_train_m, ntile(linpred, 1000)) # group up the data by eta_x into 100 groups
hldf <- summarise(gdf, y=sum(SeriousDlqin2yrs==1), ppred=mean(predprob), count=n())
head(hldf)

## # A tibble: 6 x 4
##   `ntile(linpred, 1000)`    y    ppred count
##   <int> <int>    <dbl> <int>
## 1 1       8 0.000387    73
## 2 2      13 0.00251     73
## 3 3      22 0.00485     73
## 4 4      14 0.00717     73
## 5 5      12 0.00920     73
## 6 6      15 0.0105      73

# We adjust the size of the bins until there's only one group with less than 5
hldf[hldf$count<5,]
```

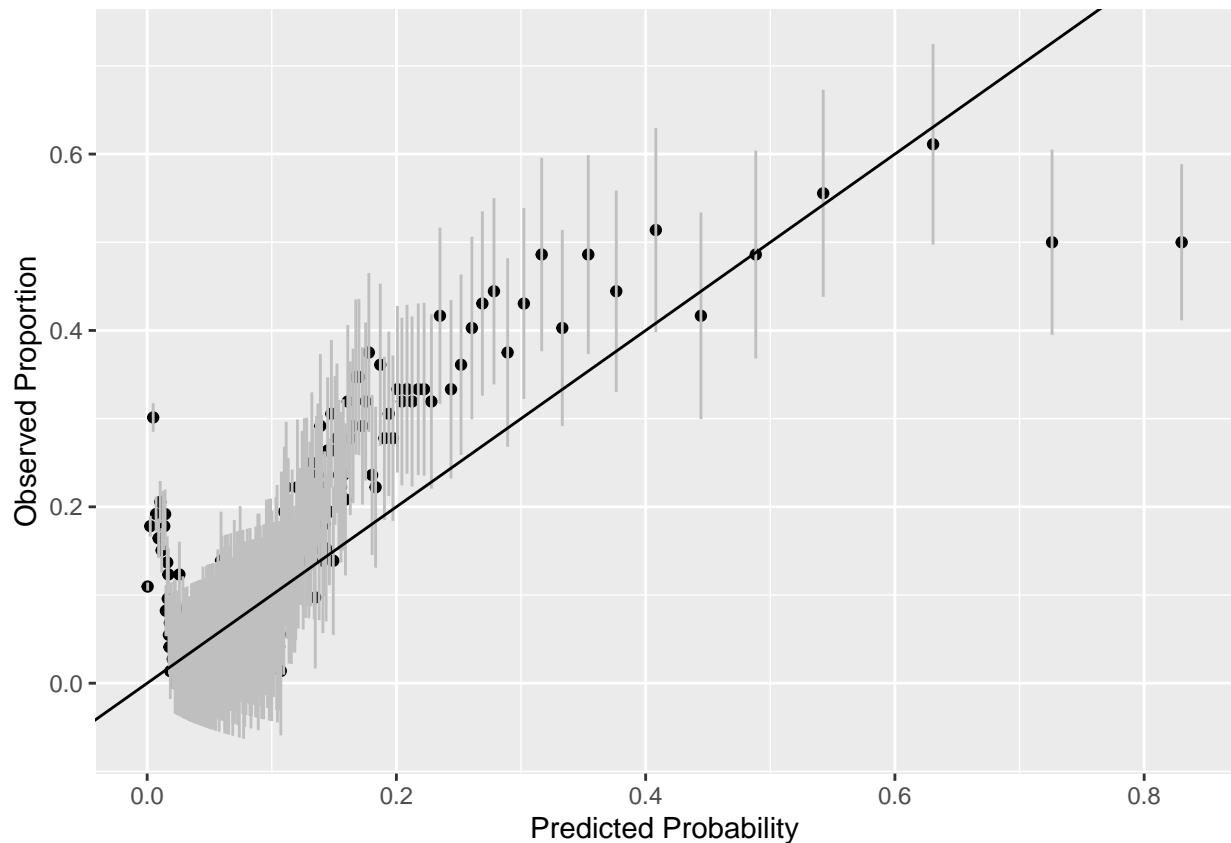
```

## # A tibble: 0 x 4
## # ... with 4 variables: ntile(linpred, 1000) <int>, y <int>, ppred <dbl>,
## #   count <int>

# Observed Proportion Confidence Interval vs Predicted Probability
hldf <- mutate(hldf, se.fit=sqrt(ppred*(1-ppred)/count))

ggplot(hldf,aes(x=ppred,y=y/count,ymin=y/count-2*se.fit, ymax=y/count+2*se.fit))+geom_point()+geom_linerange(color=grey(0.75))+geom_abline(intercept = 0,slope = 1)+xlab("Predicted Probability")+ylab("Observed Proportion")

```



```

# Hosmer-Lemeshow statistics
hlstat <- with(hldf, sum((y-count*ppred)^2/(count * ppred * (1-ppred))))/c(hlstat, nrow(hldf))

## [1] 7509.678 1000.000

# The p-value is given by:
1-pchisq(hlstat, nrow(hldf)-2)

## [1] 0

```

AUC

Need to fix the Prediction function!!!

Predict -> Currently gives 120k, expecting 80k

```
#Final Model(w/ interaction term ReugularMedicine * PhysicallyActive)
result_m2 = predict(model.final, newdata=test.x, type="response")

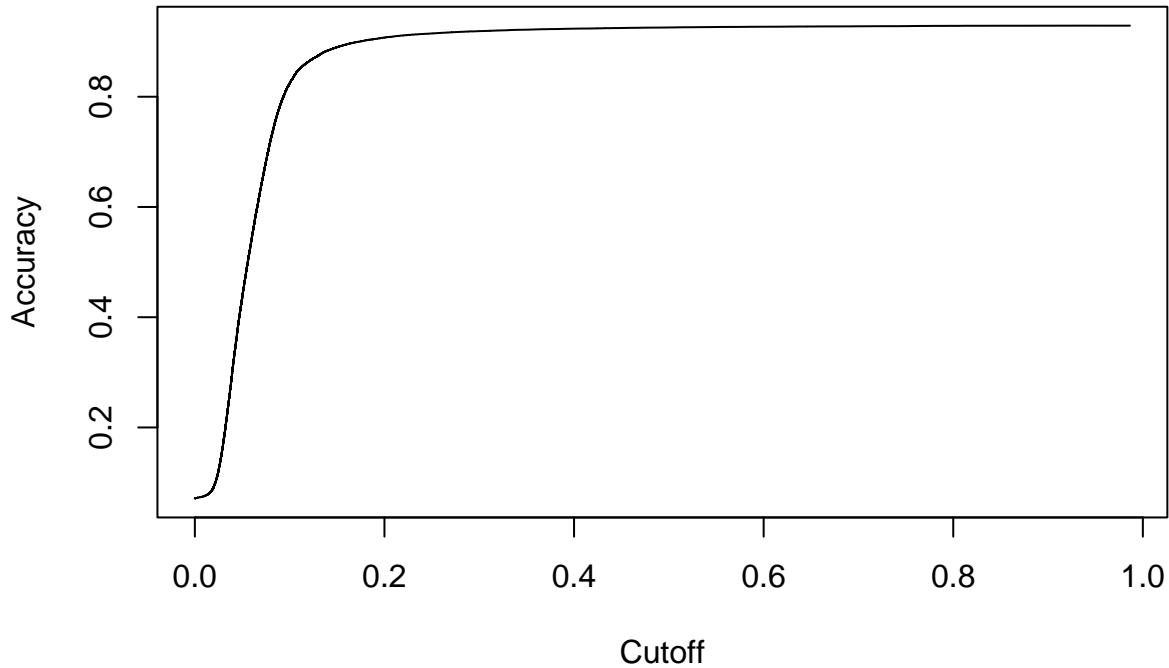
## Warning: 'newdata' had 48108 rows but variables found have 72161 rows

# Fix Later!!! Forced the length to be the same
result_m2 <- result_m2[1:length(test.y)]
head(test.y)

## [1] 0 0 0 0 0 0

pred_m2 = prediction(result_m2, test.y)

plot(performance(pred_m2, "acc")) #It seems like 0.52 cutoff has the highest accuracy
```



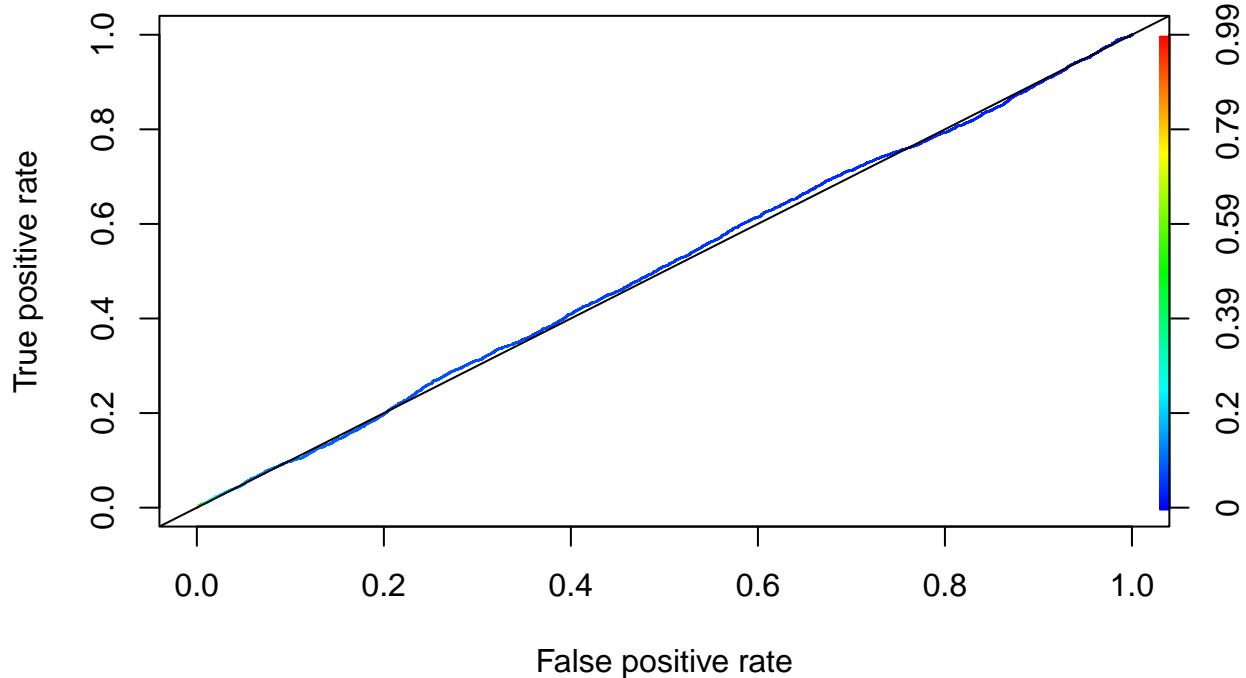
```

table(test.y, result_m2>0.2)

##
## test.y FALSE TRUE
##      0 43559 1133
##      1  3321   95

#Accuracy :
#Sensitivity :
#Specificity :
#The Specificity and accuracy improved a bit compared to the previous model without interaction term, s
plot(performance(pred_m2,"tpr","fpr"), colorize=T)
abline(0,1)

```



```

#Now we calculate the area under the curve (AUC) and accuracy of the model given above (glmModel2)
auc_ROCR2 <- performance(pred_m2, measure = "auc")
auc_ROCR2@y.values[[1]]

## [1] 0.5049617

#predict(kmeans_model_train, newdata=test.x, type="response")
prediction_result <- predict.glm(glm_pca_original_data, newdata = testing_data)

## Warning: 'newdata' had 48108 rows but variables found have 72161 rows

```

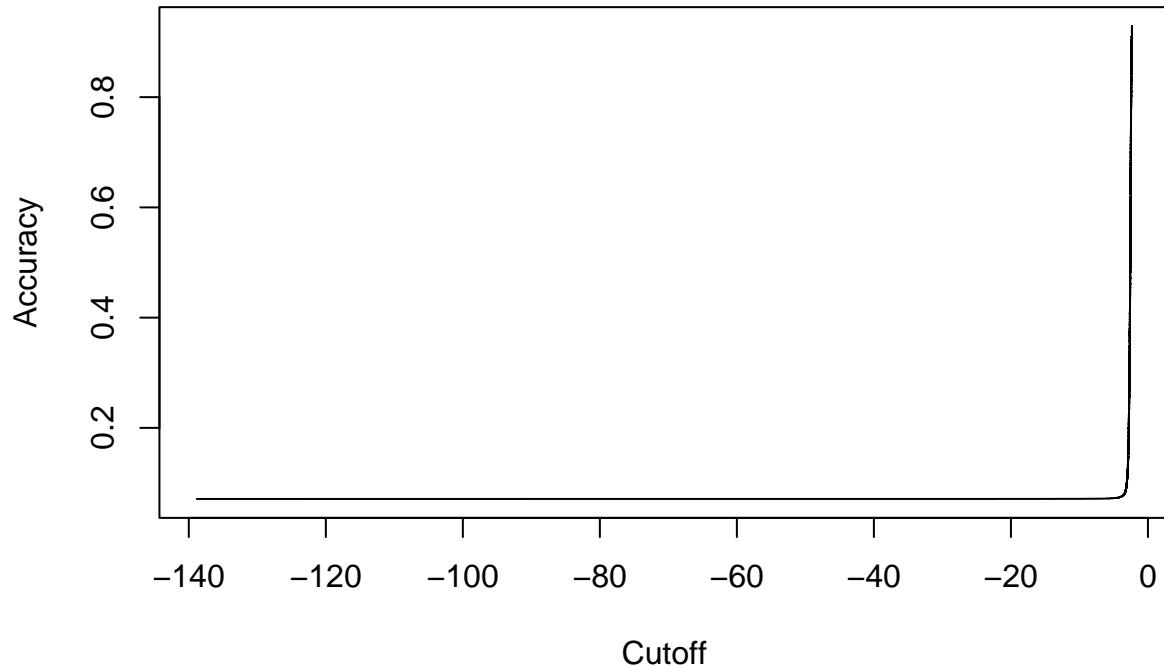
```

prediction_result <- prediction_result[1:length(test.y)]
head(test.y)

## [1] 0 0 0 0 0 0

prediction_pca = prediction(prediction_result, test.y)
plot(performance(prediction_pca, "acc"))

```



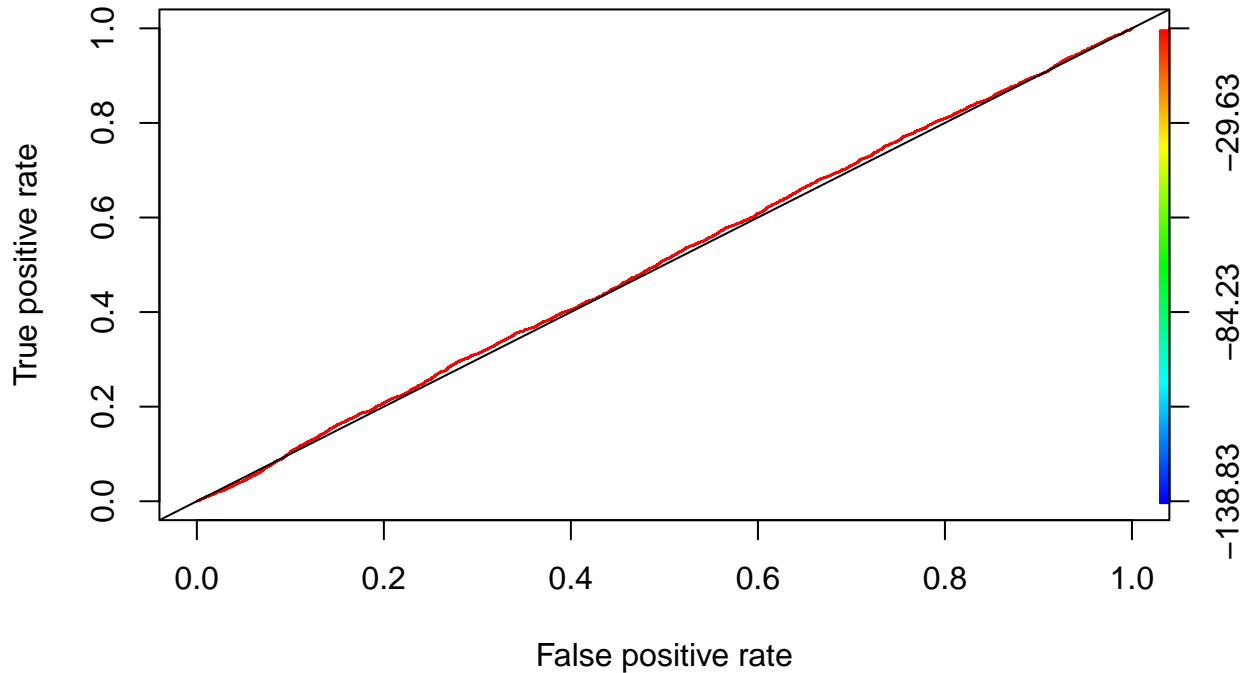
```

table(test.y, result_m2>0.2)

##
## test.y FALSE TRUE
##      0 43559 1133
##      1 3321    95

plot(performance(prediction_pca,"tpr","fpr"), colorize=T)
abline(0,1)

```



```
pca_auc_ROCR2 <- performance(prediction_pca, measure = "auc")
pca_auc_ROCR2@y.values[[1]]
```

```
## [1] 0.5073262
```

Maybe don't need session

(0) Ensemble Learning Used by Paper

- Lasso Ensemble Algorithm
- Aggregating base learner: Weighted base=learner

Not Sure if we use this!!

```
# Check how the model fits the data
# From model.final, we can calculate the difference in the two deviances from the summary(model): 987.2
#Number of regressors in the model: 813-802=11
pchisq(473.75,12)

## [1] 1
```

```

#The area below 473.75 is one which means the area above it is almost zero. This means that our model has a good fit.
print(paste("Pearson's X^2 =",round(sum(residuals(model.final,type="pearson")^2),3)))

## [1] "Pearson's X^2 = 1270781.333"

qchisq(0.95,802)

## [1] 868.9936

#781.61<868.99, so we fail to reject the null hypothesis and conclude that the logistic model fits the data well.

sum(cs_train[,2])

## [1] 4941

sum(cs_test[,2])

## [1] 3416

```

Cluster Attempt

```

# Cluster Grouping Majority Data (Try to cluster data by age/monthly income)
set.seed(1) # for reproducibility

head(cs_train_maj)

```

```

##          X SeriousDlqin2yrs RevolvingUtilizationOfUnsecuredLines age
## 30484 30484                 0                         0.09474044 47
## 74216 74216                 0                         0.05277307 68
## 54077 54077                 0                         0.73665267 36
## 86784 86784                 0                         0.02365700 36
## 14388 14388                 0                         0.00581900 44
## 31442 31442                 0                         0.11523783 40
##      NumberOfTime30.59DaysPastDueNotWorse DebtRatio MonthlyIncome
## 30484                               0 0.36842105           6250
## 74216                               0 0.22759781          11884
## 54077                               0 0.09445277           2000
## 86784                               0 0.21710409           5600
## 14388                               0 0.18525779           5100
## 31442                               1 2.49168646          2946
##      NumberOfOpenCreditLinesAndLoans NumberOfTimes90DaysLate
## 30484                           8                         0
## 74216                          13                         0
## 54077                           6                         0
## 86784                           9                         0
## 14388                          24                         0
## 31442                          17                         0
##      NumberRealEstateLoansOrLines NumberOfTime60.89DaysPastDueNotWorse

```

```

## 30484           1           0
## 74216           2           0
## 54077           0           0
## 86784           1           0
## 14388           1           0
## 31442           3           0
##   NumberOfDependents
## 30484           3
## 74216           1
## 54077           1
## 86784           0
## 14388           2
## 31442           0

# Test with smaller group based on monthly income range
cs_train_maj_g1 <- cs_train_maj[cs_train_maj$NA_Indicator==0,] # Filter out NA monthly income first
cs_train_maj_g1 <- cs_train_maj_g1[cs_train_maj_g1$MonthlyIncome<20000,] #38035 obs
head(cs_train_maj_g1)

## [1] X           SeriousDlqin2yrs
## [3] RevolvingUtilizationOfUnsecuredLines age
## [5] NumberOfTime30.59DaysPastDueNotWorse DebtRatio
## [7] MonthlyIncome           NumberOfOpenCreditLinesAndLoans
## [9] NumberOfTimes90DaysLate NumberRealEstateLoansOrLines
## [11] NumberOfTime60.89DaysPastDueNotWorse NumberOfDependents
## <0 rows> (or 0-length row.names)

# Create a dat with the two predictors of interest
dat <- cs_train_maj_g1[,c(4,7)] # Age and MonthlyIncome
head(dat)

## [1] age      MonthlyIncome
## <0 rows> (or 0-length row.names)

n_maj <- nrow(dat) # get number of rows

# Initial assignments to three groups that will need to update
assignments <- factor(sample(c(1,2,3), n_maj, replace = TRUE))
#plot(dat, col=assignments, xlim = c(0,110), asp=1)
#plot(dat, col=assignments)

```

Maybe don't need session (END)