

# Coursework (CM3111)

## Big Data Analytics

Alistair Quinn

December 16, 2017

---

## 1 Data Exploration

### 1.1 Dataset Choice

I have chosen a dataset that was used in the CoIL 2000 Challenge, which is called the Caravan Insurance Challenge dataset available here: <https://www.kaggle.com/uciml/caravan-insurance-challenge>

The data set is free to use for non-commercial use. Although sourced from the UCL Machine Learning organization's kaggle page, the dataset is owned by Sentient Machine Research.

### 1.2 Technology-Platform

The dataset is contained in a csv file that is 245KB in size. As the dataset is so small I not need to use Big Data technology such as Hadoop. Instead I will use RStudio on a windows PC. I chose the dataset based on my current ability, and I found the idea of the dataset interesting.

### 1.3 Problem Statement Data Exploration

Each row in the table corresponds to a post code. The task with this dataset is to identify potential purchasers of caravan insurance policies. The class label in the dataset is called CARAVAN and has two values, 0 or 1. CARAVAN is 1 when that row would potentially purchase a caravan insurance policy.

First I'll load the dataset I'm using.

```
#Set WD
setwd("D:/RGU/3rdYear/Semester1/Big Data Analytics/Coursework/wd")
#Load Data
df <- read.csv("Data/caravan-insurance-challenge.csv")
```

I will now explore my dataset to identify its features and learn more about it.

During my data exploration, I will be visualising my data using a package called ggplot2. This is a plotting system for R. I have chosen to use this package as I have experience with it from the labs in the course. It has some powerful plotting functions, and works in an intuitive way.

```
library(ggplot2)
## Warning: package 'ggplot2' was built under R version 3.4.2
```

Going to take a look at the number of rows and columns

```
#Rows and Cols
nrow(df)
## [1] 9822
ncol(df)
## [1] 87
```

There are currently 9822 rows, and 87 columns in the dataset

Going to take a look at the features of the dataset

```
#Names of columns
names(df)
```

## [1]	"ORIGIN"	"MOSTYPE"	"MAANTHUI"	"MGEMOMV"	"MGEMLEEF"	"MOSHOOFD"
## [7]	"MGODRK"	"MGODPR"	"MGODOV"	"MGODGE"	"MRELGE"	"MRELSA"
## [13]	"MRELOV"	"MFALLEEN"	"MFGEKIND"	"MFWEKIND"	"MOPLHOOG"	"MOPLMIDD"
## [19]	"MOPLLAAG"	"MBERHOOG"	"MBERZELF"	"MBERBOER"	"MBERMIDD"	"MBERARBG"
## [25]	"MBERARBO"	"MSKA"	"MSKB1"	"MSKB2"	"MSKC"	"MSKD"
## [31]	"MHHUUR"	"MHKOOP"	"MAUT1"	"MAUT2"	"MAUTO"	"MZFONDS"
## [37]	"MZPART"	"MINKM30"	"MINK3045"	"MINK4575"	"MINK7512"	"MINK123M"
## [43]	"MINKGEM"	"MKOOPKLA"	"PWAPART"	"PWABEDR"	"PWALAND"	"PPERSAUT"
## [49]	"PBESAUT"	"PMOTSCO"	"PVRAAUT"	"PAANHANG"	"PTRACTOR"	"PWERKT"
## [55]	"PBROM"	"PLEVEN"	"PPERSONG"	"PGEZONG"	"PWAOREG"	"PBRAND"
## [61]	"PZEILPL"	"PPLEZIER"	"PFIETS"	"PINBOED"	"PBYSTAND"	"AWAPART"
## [67]	"AWABEDR"	"AWALAND"	"APERSAUT"	"ABESAUT"	"AMOTSCO"	"AVRAAUT"
## [73]	"AAANHANG"	"ATTRACTOR"	"AWERKT"	"ABROM"	"ALEVEN"	"APERSONG"
## [79]	"AGEZONG"	"AWAOREG"	"ABRAND"	"AZEILPL"	"APLEZIER"	"AFIETS"
## [85]	"AINBOED"	"ABYSTAND"	"CARAVAN"			

Variables beginning with M are demographic statistics of the postal code.

- **ORIGIN:** *train* or *test*, as described above
- **MOSTYPE:** Customer Subtype; see L0
- **MAANTHUI:** Number of houses 1 - 10
- **MGEMOMV:** Avg size household 1 - 6
- **MGEMLEEF:** Avg age; see L1
- **MOSHOOFD:** Customer main type; see L2

Variables beginning with P and A refer to product ownership and insurance statistics of the postal code. Variables beginning with P refer to contribution policies.

- PWAPART: Contribution private third party insurance
- PWABEDR: Contribution third party insurance (firms) ...
- PWALAND: Contribution third party insurance (agriculture)
- PPERSAUT: Contribution car policies
- PBESAUT: Contribution delivery van policies
- PMOTSCO: Contribution motorcycle/scooter policies
- PVRAAUT: Contribution lorry policies
- PAANHANG: Contribution trailer policies
- PTRACTOR: Contribution tractor policies
- PWERKT: Contribution agricultural machines policies
- PBROM: Contribution moped policies
- PLEVEN: Contribution life insurances
- PPERSONG: Contribution private accident insurance policies
- PGEZONG: Contribution family accidents insurance policies
- PWAOREG: Contribution disability insurance policies
- PBRAND: Contribution fire policies
- PZEILPL: Contribution surfboard policies
- PPLEZIER: Contribution boat policies
- PFIETS: Contribution bicycle policies
- PINBOED: Contribution property insurance policies
- PBYSTAND: Contribution social security insurance policies

variables beginning with A refer to number of policies.

- **AWAPART:** Number of private third party insurance 1 - 12
- **AWABEDR:** Number of third party insurance (firms) ...
- **AWALAND:** Number of third party insurance (agriculture)
- **APERSAUT:** Number of car policies
- **ABESAUT:** Number of delivery van policies
- **AMOTSCO:** Number of motorcycle/scooter policies
- **AVRAAUT:** Number of lorry policies
- **AAANHANG:** Number of trailer policies
- **TRACTOR:** Number of tractor policies
- **AWERKT:** Number of agricultural machines policies
- **ABROM:** Number of moped policies
- **ALEVEN:** Number of life insurances
- **APERSONG:** Number of private accident insurance policies
- **AGEZONG:** Number of family accidents insurance policies
- **AWAOREG:** Number of disability insurance policies
- **ABRAND:** Number of fire policies
- **AZEILPL:** Number of surfboard policies
- **APLEZIER:** Number of boat policies
- **AFIETS:** Number of bicycle policies
- **AINBOED:** Number of property insurance policies
- **ABYSTAND:** Number of social security insurance policies

Going to check the factors of the dataset

```
#Names of columns
sapply(df, levels)
```

I have omitted the result of the above code, as it was far too large. Most of the columns in the current data are numeric values but are actually supposed to be factors. I will refactor these columns during pre-processing. The only variable that has been turned into a factor by R is the first one, ORIGIN. I will explore this factor later.

There are 4 keys that relate to this dataset. A key for customer subtype:

#### **L0: Customer subtype**

- 1: High Income, expensive child
- 2: Very Important Provincials
- 3: High status seniors
- 4: Affluent senior apartments
- 5: Mixed seniors
- 6: Career and childcare
- 7: Dinki's (double income no kids)
- 8: Middle class families
- 9: Modern, complete families
- 10: Stable family
- 11: Family starters
- 12: Affluent young families
- 13: Young all american family
- 14: Junior cosmopolitan
- 15: Senior cosmopolitans
- 16: Students in apartments
- 17: Fresh masters in the city
- 18: Single youth
- 19: Suburban youth
- 20: Ethnically diverse
- 21: Young urban have-nots
- 22: Mixed apartment dwellers
- 23: Young and rising
- 24: Young, low educated
- 25: Young seniors in the city
- 26: Own home elderly
- 27: Seniors in apartments
- 28: Residential elderly
- 29: Porchless seniors: no front yard
- 30: Religious elderly singles
- 31: Low income catholics
- 32: Mixed seniors
- 33: Lower class large families
- 34: Large family, employed child
- 35: Village families
- 36: Couples with teens 'Married with children'
- 37: Mixed small town dwellers
- 38: Traditional families
- 39: Large religious families
- 40: Large family farms
- 41: Mixed rurals

A key for average age:

**L1: average age keys:**

1: 20-30 years 2: 30-40 years 3: 40-50 years 4: 50-60 years 5: 60-70 years 6: 70-80 years

A key of customer main types:

**L2: customer main type keys:**

- 1: Successful hedonists
- 2: Driven Growers
- 3: Average Family
- 4: Career Loners
- 5: Living well
- 6: Cruising Seniors
- 7: Retired and Religious
- 8: Family with grown ups
- 9: Conservative families
- 10: Farmers

A key of percentage ranges:

**L3: percentage keys:**

- 0: 0%
- 1: 1 - 10%
- 2: 11 - 23%
- 3: 24 - 36%
- 4: 37 - 49%
- 5: 50 - 62%
- 6: 63 - 75%
- 7: 76 - 88%
- 8: 89 - 99%
- 9: 100%

and a key of total number ranges:

**L4: total number keys:**

- 0: 0
- 1: 1 - 49
- 2: 50 - 99
- 3: 100 - 199
- 4: 200 - 499
- 5: 500 - 999
- 6: 1000 - 4999
- 7: 5000 - 9999
- 8: 10,000 - 19,999
- 9:  $\geq 20,000$

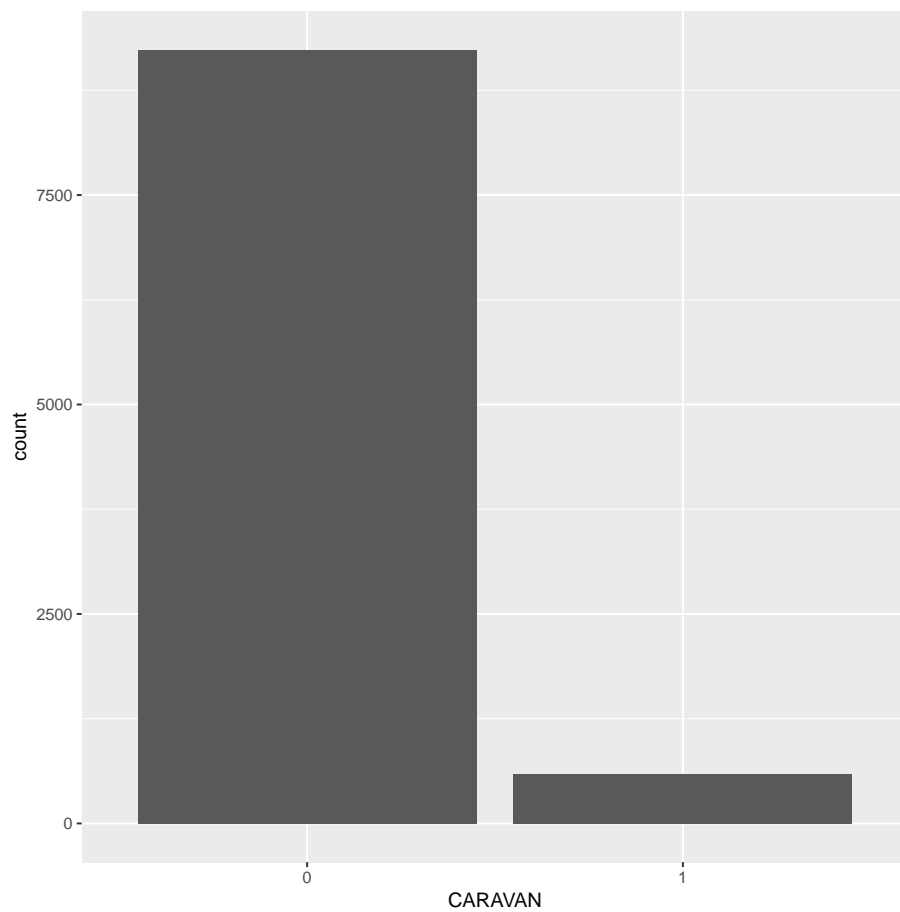


I will use these keys to turn the appropriate columns into factors later, by mapping the numeric values in the dataset to the appropriate value from the keys.

I will now take a look at the class label distribution

```
#Class label freq
classLabelFreq <- data.frame(df$CARAVAN)
classLabelFreq$df.CARAVAN <- as.factor(df$CARAVAN)

#Class label Distribution Plot
ggplot(classLabelFreq, aes(x=df.CARAVAN)) + geom_bar() + labs(x="CARAVAN")
```



```
#Size of each factor level
length(classLabelFreq[classLabelFreq$df.CARAVAN=="0",])
## [1] 9236
```

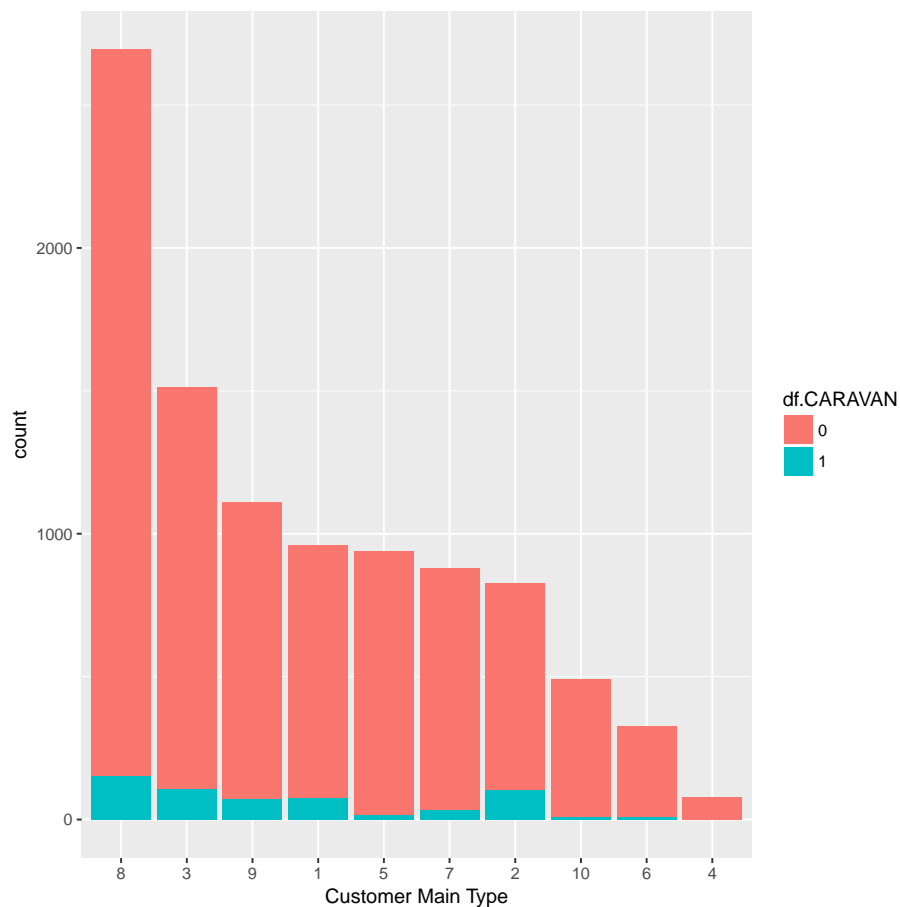
```
length(classLabelFreq[classLabelFreq$df.CARAVAN=="1",])
## [1] 586
```

There are 586 records that are likely to want caravan insurance. 9236 records that do not. Dataset has an imbalanced distribution in the class label. I will use a resampling technique during pre-processing to compensate for this.

Lets take a look at the distribution of Main Customer Type

```
#Cust main type
custMainType <- data.frame(df$MOSHOOFD,df$CARAVAN)
custMainType$df.MOSHOOFD <- as.factor(custMainType$df.MOSHOOFD)
custMainType$df.CARAVAN <- as.factor(custMainType$df.CARAVAN)

#Plot of Customer Main Type
ggplot(custMainType,aes(x=reorder(df.MOSHOOFD,df.MOSHOOFD,function(x)-length(x)),fill=df.CARAVAN
```



Most frequent Main Customer Type is 8:Family with Grown Ups, 2nd Most frequent is 3:Average Family and the 3rd most frequent is 9:conservative families.

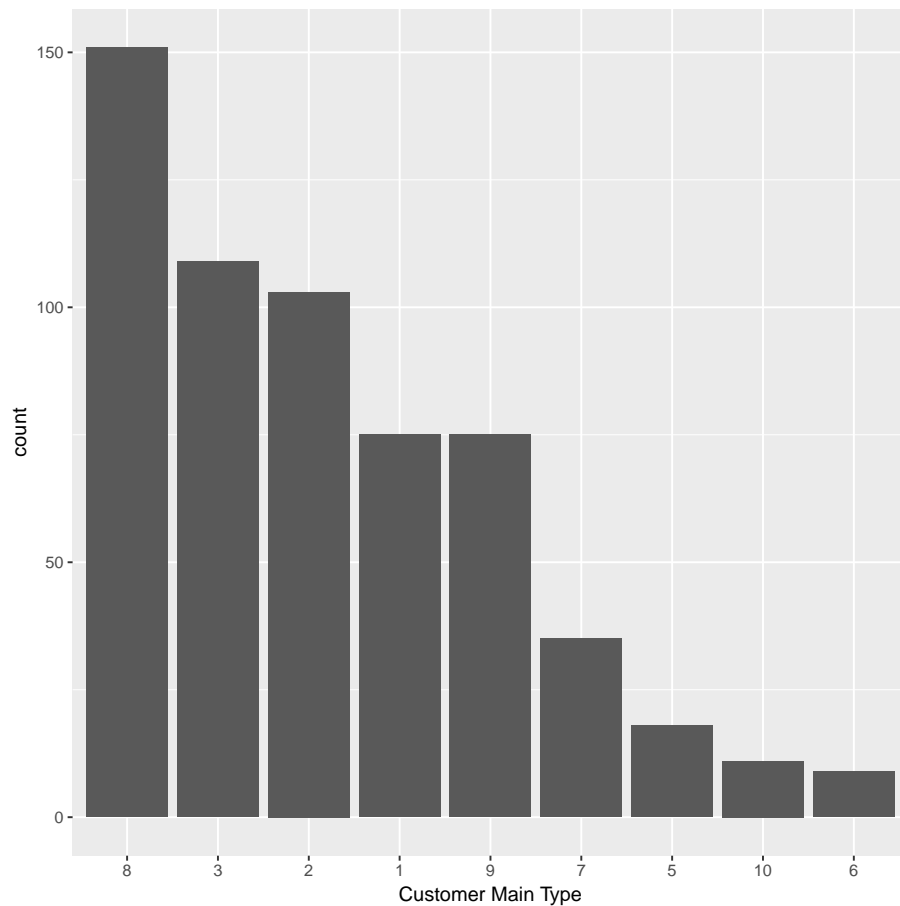
The least frequent type is 4:Career Loners ,2nd least frequent is 6:Cruising Seniors and the 3rd least frequent is 10:Farmers.

Comparing to where CARAVAN is TRUE, you can see that the two most frequent main types are the same as the whole dataset. Customer Main Type 2: Driven Growers seems to be a bit more prominent where CARAVAN is true. You can also see that there are no instances of group 4: Career Loners in the rows where CARAVAN is true.

Lets take a closer look at the rows where CARAVAN is TRUE:

```
#Wants caravan
wantsCaravan <- df[df$CARAVAN==1,]
wantsCaravan$MOSHOOFD <- as.factor(wantsCaravan$MOSHOOFD)
wantsCaravan$MOSTYPE <- as.factor(wantsCaravan$MOSTYPE)

#Plot of Customer Main Type where wants caravan
ggplot(wantsCaravan,aes(x=reorder(MOSHOOFD,MOSHOOFD,function(x)-length(x)))) + geom_bar() +
```



```
#Max and Min
mainCustType = table(wantsCaravan$MOSH00FD)
names(which.max(mainCustType))
## [1] "8"

names(which.min(mainCustType))
## [1] "6"
```

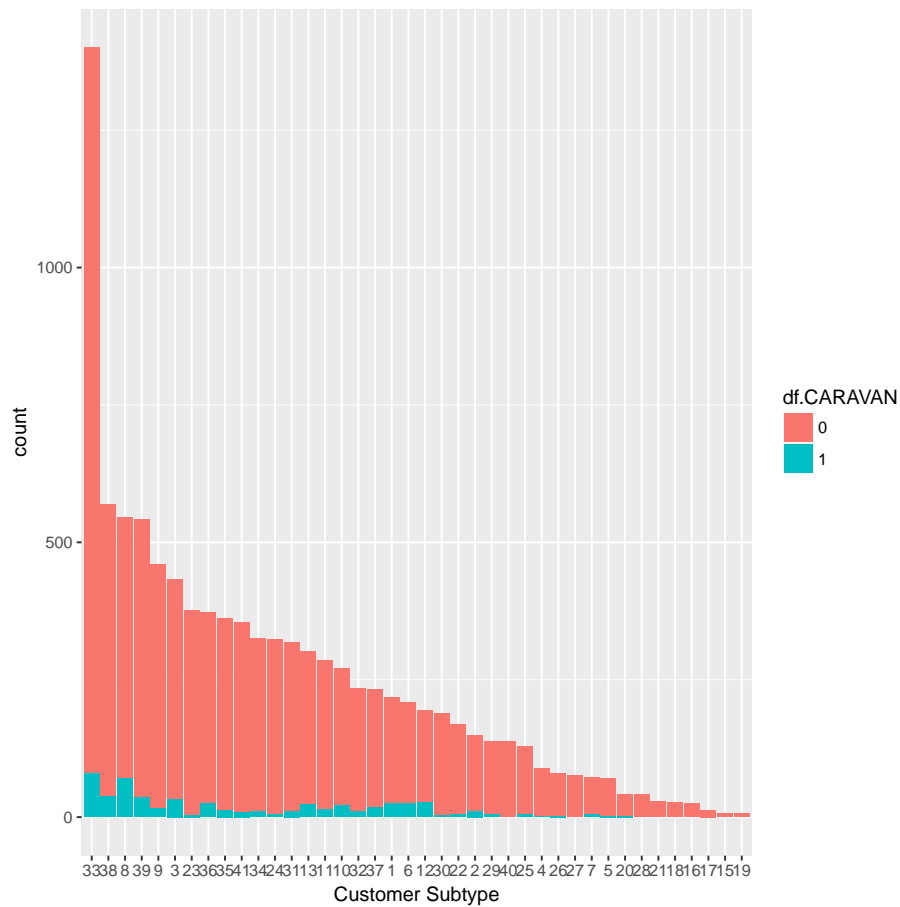
The top 3 Main Customer Types are 8:Family with Grown Ups, 3: Average Family and 2:Driven Growers. As stated before the first 2 types are the same as the dataset as a whole. Type 2:Driven Growers has now overtaken 9:conservative families. 9:Conservative Families is joint 4th. Interesting that category 1: Successful Hedonists has the same number of occurrences as 9:Conservative Families. Hedonists are people who devote their lives to the pursuit of pleasure. There is perhaps a connection there between the idea of traveling by caravan and being a hedonist.

The 3 least frequent Main Custom Types 6:Cruising Seniors ,10:Farmers and 5:Living Well. As stated before there are no instances of 4:Career Loners when CARAVAN is TRUE. There is too small a number of instances of 4:Career Loners in the whole dataset to really say there is correlation but it is possible.

Now going to take a look at Customer Subtype

```
#Sub cust type
subCustType <- data.frame(df$MOSTYPE,df$CARAVAN)
subCustType$df.MOSTYPE <- as.factor(subCustType$df.MOSTYPE)
subCustType$df.CARAVAN <- as.factor(subCustType$df.CARAVAN)

#Plot of Customer subtype
ggplot(subCustType,aes(x=reorder(df.MOSTYPE,df.MOSTYPE,function(x)-length(x)),fill=df.CARAVAN))
```



The top 3 subtypes are 33:Lower class large families ,38:Traditional families and 8:Middle class families.

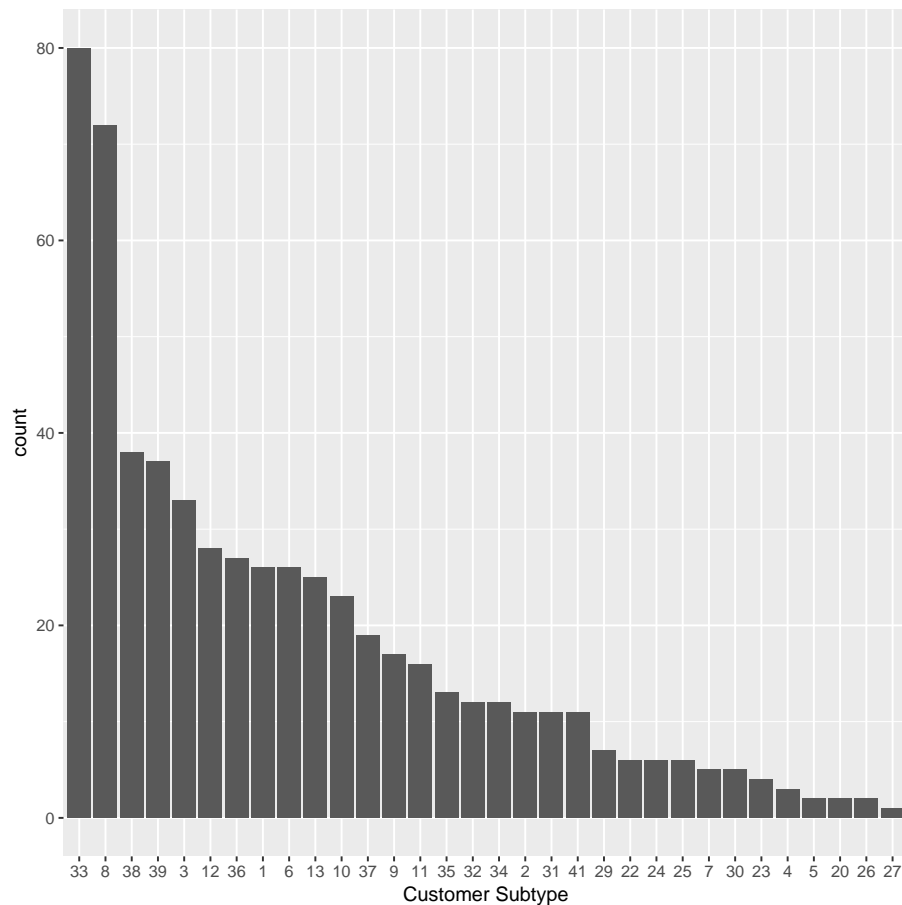
The bottom 3 are 19:Suburban Youth ,15:Senior cosmopolitans and 17:Fresh masters in the city. The dataset contains mostly data about families.

Taking a look at when CARAVAN is TRUE, you can see that It doesn't follow the same trend in frequency as the dataset as a whole This might mean that It could be used as a useful predictor. You can also see that there are no rows where CARAVAN is true that have subtypes of 40:Large family farms, 21:Young urban have-nots, 18:Single youth, 16:Students in apartments, 17:Fresh masters in the city, 15:Senior cosmopolitans and 19:Suburban youth.

This could suggest that postcode areas most comprised with families are more likely to purchase caravan insurance.

Going to take a look at customer subtype when CARAVAN is TRUE

```
#Plot of Customer Subtype where wants caravan
ggplot(wantsCaravan,aes(x=reorder(MOSTYPE,MOSTYPE,function(x)-length(x)))) + geom_bar() + la
```



```

#Max and Min
subCustType = table(wantsCaravan$MOSTYPE)
names(which.max(subCustType))

## [1] "33"

names(which.min(subCustType))

## [1] "27"

```

Top 3 subtypes are the same as the dataset as a whole but subtype 8:Middle class families is now more frequent than 38:Traditional families. This might mean that areas consisting of lower to middle class families are more likely to purchase caravan insurance than areas with mostly traditional families.

The bottom 3 subtypes are 27:Seniors in apartments ,26:Own home elderly and 20:Ethnically diverse. This supports the theory that areas with higher amounts of families are the most likely to purchase caravan insurance.

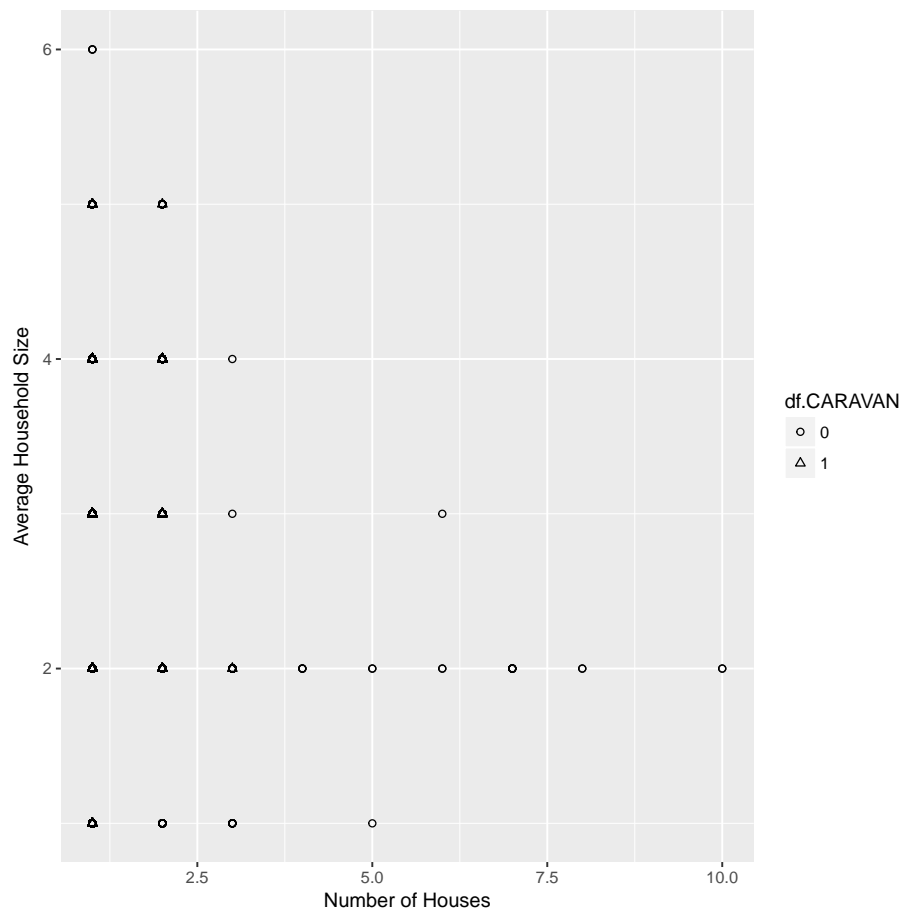
I am not going to take a look at two columns. MAANTHUI(Number of houses) and MGEMOMV(Avg size of household). MAATHUI is in the range of 1-10 and MGEMOMV is in the range of 1-6. They are the only two numeric values in the dataset the rest are factors. I will look at them together to see if there is any correlation. I will use ggplot2 again to make a scatter plot. These are integer values there will be overlap.

```

#Number of Houses and Avg size of household
houseData<-data.frame(df$MAANTHUI,df$MGEMOMV,df$CARAVAN)
houseData$df.CARAVAN<-as.factor(houseData$df.CARAVAN)

#ScatterPlot of both
ggplot(houseData,aes(x=df.MAANTHUI,y=df.MGEMOMV)) + geom_point(aes(shape=df.CARAVAN)) + scale_shape_size()

```



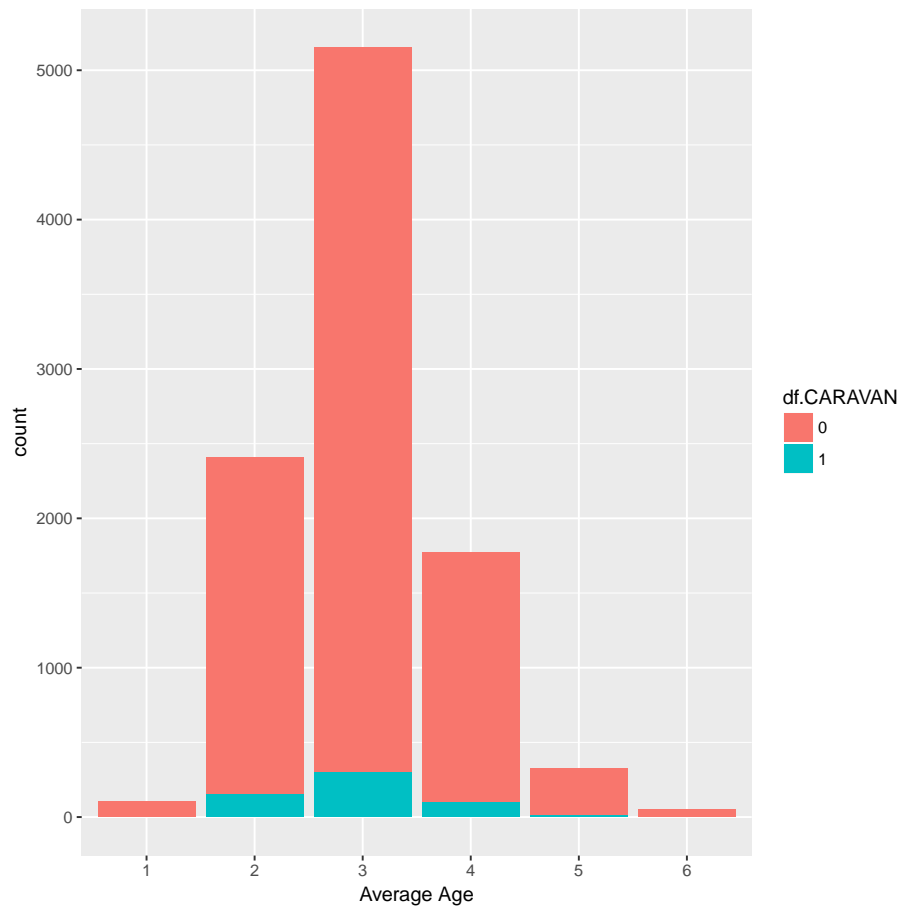
Taking a look at the results, average size of household decreases as the number of houses increases which makes sense. Looking at the points where CARAVAN is TRUE, There are no points when the number of houses is greater than 3. There are also no points when the average house size is greater than 5. This shows a potential connection between number of houses and CARAVAN. If I had to remove one of the two variables I would remove Average house size as I think number of houses has a greater correlation to CARAVAN equaling TRUE

I will now take a look at the average age variable, MGEMLEEF.

```
#Average Age
averageAge <- data.frame(df$MGEMLEEF,df$CARAVAN)
averageAge$df.MGEMLEEF <- as.factor(averageAge$df.MGEMLEEF)
averageAge$df.CARAVAN <- as.factor(averageAge$df.CARAVAN)

#Plot of Average Age
ggplot(averageAge,aes(x=df.MGEMLEEF,fill=df.CARAVAN)) + geom_bar() + labs(x="Average Age")
```





Average age is a factor, where each level is a range of ages. Lowest age range is 1:20-30 years, highest is 6:70-80 years. Looking at the graph, levels 3:40-50, 2:30-40 and 4:50-60 are the top 3 most frequent values. The extremes of 1 and 6 are the two lowest and do not have very many occurrences. Looking at instances where CARAVAN is TRUE, They are in the same order as the whole dataset. Except there are no instances where CARAVAN is TRUE that contain 1 and 6 for average age. There might be a trend that areas where the average age is 1:20-30 or 6:70-80 would not buy caravan insurance. There isn't enough data to be sure and as the data where CARAVAN is true follows a similar trend to the data as a whole it's unlikely there is a correlation between average age and CARAVAN.

The first column in the data set is ORIGIN. It has two values

```
#Levels of ORIGIN
levels(df$ORIGIN)
## [1] "test" "train"
```

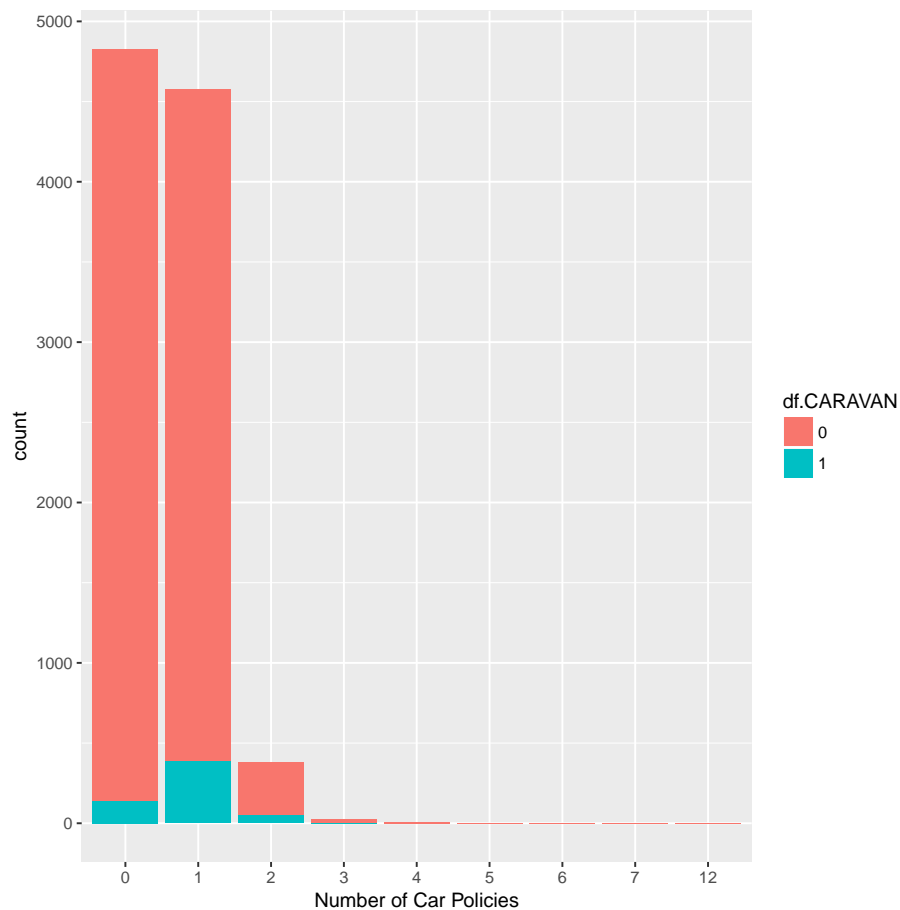
This is the original source of the row from the challenge. The rows are already split into a train set and test set. I will remove this column later during pre-processing, as I plan to resample the data and split the data into train and test sets myself.

I have now looked at all the main variables.

I am not going to take a look at the APERSAUT column, which is the Number of Car Policies column. This is a factor, where each level equates to a range of values. The ranges are defined in the total number key. I have a feeling that it might be an important predictor, so want to try and confirm my theory.

```
#Number of Car Policies
numberOfCarPolicies <- data.frame(df$APERSAUT,df$CARAVAN)
numberOfCarPolicies$df.APERSAUT <- as.factor(numberOfCarPolicies$df.APERSAUT)
numberOfCarPolicies$df.CARAVAN <- as.factor(numberOfCarPolicies$df.CARAVAN)

#Plot of APERSAUT
ggplot(numberOfCarPolicies,aes(x=reorder(df.APERSAUT,df.APERSAUT,function(x)-length(x)),fill=
```

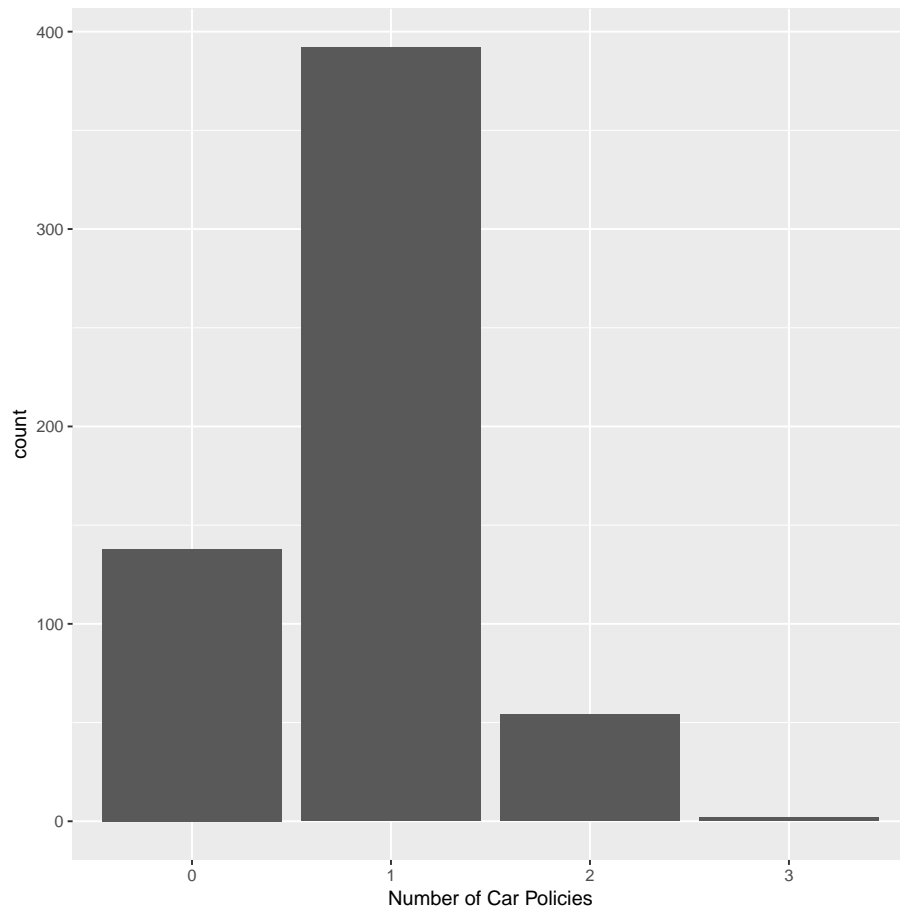


There is a factor level of 12 here on the graph. There is only supposed to be levels 0-9. This is likely a mistake in the dataset. When I factor this column later during pre processing and remove n/a values this wrong value should get removed.

Looking at when CARAVAN is TRUE, there is a possible relation between when the number of car policies are the level 1:1-49 and when CARAVAN is TRUE.

```
#Number of Car Policies When Caravan is TRUE
wantsCaravan$APERSAUT <- as.factor(wantsCaravan$APERSAUT)

#Plot of number of car policies (caravan is TRUE)
ggplot(wantsCaravan, aes(x=APERSAUT)) + geom_bar() + labs(x="Number of Car Policies")
```



This factor is supposed to have 9 levels in total, but this column factored only contains the lowest 4 possible factors, 0,1,2 and 3. 1 is the most frequent. In this case 1 represents the range of values 1-49. 0 represents 0, 2 represents the range of values 50-99 and 3 represents the range of values 100-199.

This would suggest that post codes that contain a number of cars in the range of 1-49 are most likely to purchase caravan insurance. There are only 2 rows in the dataset that contain 3 for this column so this could be considered an outlier. Based on the graph the range of values for number of cars could be 0-99, or 0-199 if you include the two rows that have the value 3. There isn't enough data here to be sure.

I was originally confused by this result. I had assumed that I would find the opposite, that areas with large numbers of cars would likely need caravan insurance as I assumed you would need a car to tow a caravan. After further study of the information about the dataset, in this case CARAVAN actually refers to mobile homes. This would suggest the data is based on american post codes (or

area codes) and that 'caravan' refers to a self driving vehicle that you can sleep in rather than a traditional british caravan.

Based on the graph its possible there is a correlation between this variable and CARAVAN being TRUE. This variable might be an important predictor.

## 1.4 Pre-Processing

I will now pre-process my data before I begin building models

First I will rename all the columns. This will make them a little easier to understand, and means I do not have to keep referring to the keys. I will use the `dp` package to do this.

```
library(dplyr)

## Warning: package 'dplyr' was built under R version 3.4.2

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
## filter, lag

## The following objects are masked from 'package:base':
##
## intersect, setdiff, setequal, union
```

This is a package used for manipulating frame objects in `r`. I will now use the `rename` function to rename the columns

```
#Rename columns
df<-rename(df, Customer_Subtype=MOSTYPE, Number_of_Houses=MAANTHUI, Avg_Size_Household=MGEMOMV)
#Show new names
names(df)
```

I will now refactor all the appropriate columns

```
#Refactoring
#Customer Subtype Refactor
df$Customer_Subtype <- factor(df$Customer_Subtype,levels=c(1:41),labels=c("High Income, expensive child", "High Income, moderate child", "High Income, low income child", "Moderate Income, expensive child", "Moderate Income, moderate child", "Moderate Income, low income child", "Low Income, expensive child", "Low Income, moderate child", "Low Income, low income child"))

## Error in '$<-'.data.frame'('*tmp*', Customer_Subtype, value = structure(integer(0),
. Label = c("High Income, expensive child", : replacement has 0 rows,
data has 9822

#Average Age Refactor
df$Avg_Age <- factor(df$Avg_Age,levels=c(1:6),labels=c("20-30 years", "30-40 years", "40-50 years", "50-60 years", "60-70 years", "70+ years"))
```



The ROSE package (Random Over-Sampling Examples) is a package that helps deal with binary classification with imbalanced classes, making it perfect for what I'm trying to do.

I will now use the `ovun.sample` function from the package to oversample my dataset so that there is roughly even distribution of the class label.

```
#Resample Train(Oversampling)
df<-ovun.sample(CARAVAN~.,data=df,method="over")$data
```

## 2 Modelling/Classification

I have decided to create a random forest model for classifying my dataset. I have chosen this model because it is supposed to be a high performing classifier. It is supposed to be a robust model that can handle unbalanced data such as the dataset I have chosen. I have also already had experience with this type of model from the labs in the course. To create my models I will be using the `caret` and `randomForest` R libraries.

```
## [1] 1e+05

library(caret)
## Warning: package 'caret' was built under R version 3.4.2
## Loading required package: lattice
library(randomForest)
## Warning: package 'randomForest' was built under R version 3.4.2
## randomForest 4.6-12
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:dplyr':
##
## combine
## The following object is masked from 'package:ggplot2':
##
## margin
```

The `caret` package (Classification and Regression Training) contains useful for splitting data that I will use to split my data into train and test sets. Specifically I will use the function `createDataPartition` to create a single train and test set from all of my data that I will use to build my initial model. I will also use the

function createFolds function to generate folds for 10 fold cross validation.

The randomForest package contains functions for creating randomForest models, and for evaluating variable importance in models as well as functions to calculate and plot various accuracies and other useful information about random forest models. I will use the function randomForest to create my model, and I will use the importance function to try and improve the accuracy of my model. I will also use the plot function from the package to plot error rates.

First I will create a train and test set.

```
#Partition dataset using caret
part<-createDataPartition(y=df$CARAVAN,p=0.7,list=FALSE)
train<-df[part,]
test<-df[-part,]
```

I will write a function to build a random forest model, using the randomForest function from the randomForest package. I will pass the training set, test set and allow the passing of ntrees and nodeSize as I plan to varie these values later.

```
#Function to build random forest model
buildModel<-function(trainData,testData,ntrees=100,nodeSize=1){
  #build random forest model
  model<-randomForest(trainData[,ncol(trainData)],trainData[,ncol(trainData)],xtest=testData[,ncol(testData)],ntrees=ntrees,nodeSize=nodeSize)
  #Return model
  return(model)
}
```

I will create a function to display error rates and accuracies from a model.

```
#Print Error rates and accuracies
displayResultsFromModel<-function(model,trainRows,testRows){
  #PLOT
  print("TRAIN")
  #Train OOB Error
  print(paste("Train OOB Error: ",model$err.rate[nrow(model$test$err.rate),1,drop=FALSE],sep=" "))
  #Train Factor Level 0 Error
  print(paste("Train CARAVAN=0 Error: ",model$err.rate[nrow(model$test$err.rate),2,drop=FALSE],sep=" "))
  #Train Factor Level 1 Error
  print(paste("Train CARAVAN=1 Error: ",model$err.rate[nrow(model$test$err.rate),3,drop=FALSE],sep=" "))
  #Train Accuracy
  trainAuc<-sum(diag(model$confusion))/trainRows
  print(paste("Train Accuracy: ",trainAuc,"%",sep=" "))

  #Print blank line between train and test results
  print(" ")
}
```



```

print("TEST")
#Test Error
print(paste("Test Error: ",model$test$err.rate[nrow(model$test$err.rate),1,drop=FALSE],sep=" "))
#Train Factor Level 0 Error
print(paste("Test CARAVAN=0 Error: ",model$test$err.rate[nrow(model$test$err.rate),2,drop=FALSE],sep=" "))
#Train Factor Level 1 Error
print(paste("Test CARAVAN=1 Error: ",model$test$err.rate[nrow(model$test$err.rate),3,drop=FALSE],sep=" "))
#Test Accuracy
testAuc<-sum(diag(model$test$confusion))/testRows
print(paste("Test Accuracy: ",testAuc,sep=" "))
}

```

I will now use these functions to build and test my model.

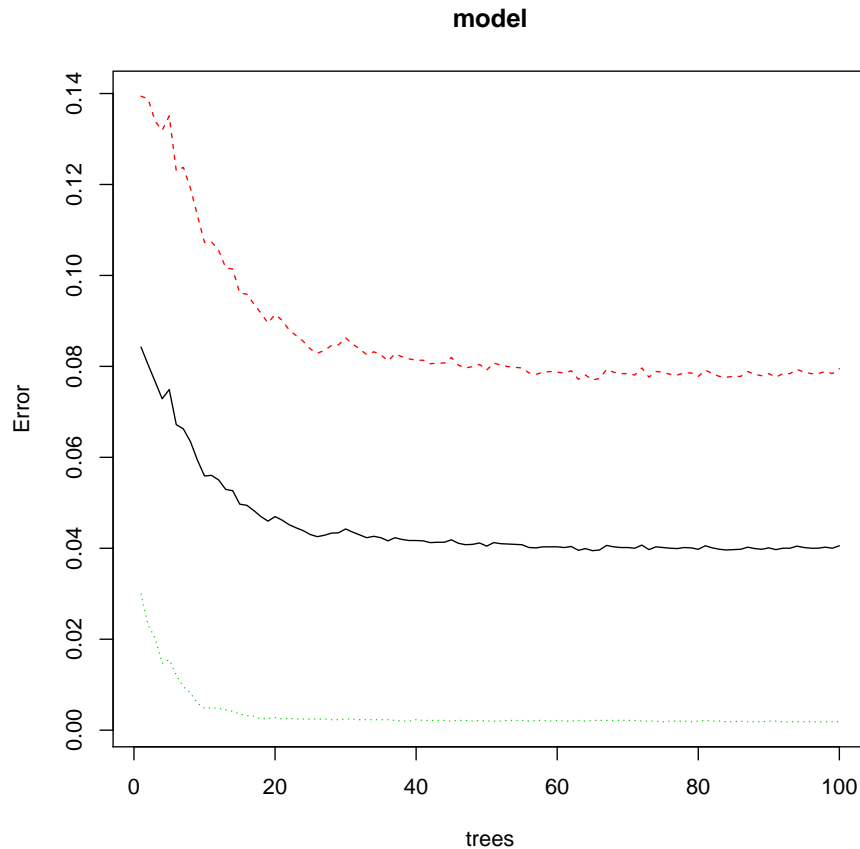
```

#Build model and display accuracies
model<-buildModel(train,test)
#Display Values
displayResultsFromModel(model,nrow(train),nrow(test))

## [1] "TRAIN"
## [1] "Train OOB Error: 0.0405520006167605"
## [1] "Train CARAVAN=0 Error: 0.079492731209403"
## [1] "Train CARAVAN=1 Error: 0.00184473481936972"
## [1] "Train Accuracy: 0.95944799938324%"
## [1] " "
## [1] "TEST"
## [1] "Test Error: 0.040849379161418"
## [1] "Test CARAVAN=0 Error: 0.0787003610108303"
## [1] "Test CARAVAN=1 Error: 0.00322927879440258"
## [1] "Test Accuracy: 0.959150620838582"

#Plot Error Rates
plot(model)

```



Now I will take a look at the error rates and accuracies of my model. During initial testing, accuracies were around 55-57 range. The train error rate tended to be lower than the test error rate but this was expected. The error rates for the CARAVAN=1 were extremely high, near 80 percent. Error rates for CARAVAN=0 were extremely low less than 1 percent. This isn't great and could be improved. More data where CARAVAN=TRUE is really needed. Perhaps different resampling methods like bootstrapping might generate better results.

I will write a function to validate the model. I will use a 10 fold cross validation method.

```
#Function to perform 10 fold cross validation
validateModel <- function(data,ntrees=100,nodeSize=1){
  #Frame to hold results
  results<-data.frame(OOB=as.numeric(),trainFalseError=as.numeric(),trainTrueError=as.numeric())
  #Folds generated using Caret packages createFolds
```

```

folds<-createFolds(data$CARAVAN,k=10,list=TRUE,returnTrain=FALSE)
for (i in 1:10){
  #Keep one set for testing, rest training
  trainData<-data[-c(folds[[i]]),]
  testData<-data[c(folds[[i]]),]
  model<-randomForest(trainData[, -ncol(trainData)],trainData[,ncol(trainData)],xtest=testData[,ncol(testData)])
  #TRAIN
  oob<-model$serr.rate[nrow(model$test$serr.rate),1,drop=FALSE]
  trainFalse<-model$serr.rate[nrow(model$test$serr.rate),2,drop=FALSE]
  trainTrue<-model$serr.rate[nrow(model$test$serr.rate),3,drop=FALSE]
  trainAccuracy<-sum(diag(model$confusion))/nrow(trainData)
  #TEST
  testError<-model$test$serr.rate[nrow(model$test$serr.rate),1,drop=FALSE]
  testFalse<-model$test$serr.rate[nrow(model$test$serr.rate),2,drop=FALSE]
  testTrue<-model$test$serr.rate[nrow(model$test$serr.rate),3,drop=FALSE]
  testAccuracy<-sum(diag(model$test$confusion))/nrow(testData)
  #Create new Row in results with values
  results[nrow(results)+1,<-c(oob,trainFalse,trainTrue,testError,testFalse,testTrue,trainAccuracy,testAccuracy)
}
#Return results
return(results)
}

```

I will also write a function to display the results. I will display the data frame as well as averages.

```

#Takes results and displays them as a whole and with averages
displayResults<-function(results){
  #Display results
  print(results)
  #TRAIN
  #OOB
  print(paste("Average OOB: ",sum(results$OOB)/nrow(results),sep=""))
  #Train CARAVAN=0 Error
  print(paste("Average CARAVAN=0 Error: ",sum(results$trainFalseError)/nrow(results),sep=""))
  #Train Caravan=1 Error
  print(paste("Average CARAVAN=1 Error: ",sum(results$trainTrueError)/nrow(results),sep=""))
  #Train Accuracy
  print(paste("Average Train Accuracy: ",sum(results$trainAccuracy)/nrow(results),sep=""))

  #Print blank line between train and test results
  print(" ")

  #Test Error
  print(paste("Average Test Error: ",sum(results$testError)/nrow(results),sep=""))
  #Test CARAVAN=0 Error
}

```

```

print(paste("Average CARAVAN=0 Error: ",sum(results$testFalseError)/nrow(results),sep=""),
      #Test CARAVAN=1 Error
print(paste("Average CARAVAN=1 Error: ",sum(results$testTrueError)/nrow(results),sep=""),
      #Test Accuracy
print(paste("Average Test Accuracy: ",sum(results$testAccuracy)/nrow(results),sep=""))
}

```

I will now use these functions to validate my model

```

#Validate Model
validateResult <- validateModel(df)
displayResults(validateResult)

##           OOB trainFalseError trainTrueError testError testFalseError
## 1  0.03484049      0.06772525    0.0021523377 0.03185745    0.06175515
## 2  0.03352324      0.06592878    0.0013153175 0.03507825    0.07034632
## 3  0.03382106      0.06640202    0.0014348918 0.03779698    0.07367281
## 4  0.03538443      0.06953802    0.0014350634 0.03559871    0.06926407
## 5  0.03370113      0.06580055    0.0017936147 0.03077754    0.06067172
## 6  0.03370315      0.06713186    0.0004782973 0.03669725    0.07359307
## 7  0.03490255      0.06809432    0.0019131890 0.02482461    0.04870130
## 8  0.03532234      0.06905679    0.0017936147 0.03453859    0.06818182
## 9  0.03544015      0.06952965    0.0015544661 0.03239741    0.06283857
## 10 0.03454480      0.06797401    0.0013154748 0.03182309    0.06168831
##      testTrueError trainAccuracy testAccuracy
## 1  0.002152853      0.9651595    0.9681425
## 2  0.000000000      0.9664768    0.9649217
## 3  0.002152853      0.9661789    0.9622030
## 4  0.002150538      0.9646156    0.9644013
## 5  0.001076426      0.9662989    0.9692225
## 6  0.000000000      0.9662969    0.9633028
## 7  0.001076426      0.9650975    0.9751754
## 8  0.001076426      0.9646777    0.9654614
## 9  0.002152853      0.9645598    0.9676026
## 10 0.002150538      0.9654552    0.9681769
## [1] "Average OOB: 0.0345183335916085"
## [1] "Average CARAVAN=0 Error: 0.0677181249100548"
## [1] "Average CARAVAN=1 Error: 0.00151862669542208"
## [1] "Average Train Accuracy: 0.965481666408392"
## [1] " "
## [1] "Average Test Error: 0.0331389858382584"
## [1] "Average CARAVAN=0 Error: 0.065071313662863"
## [1] "Average CARAVAN=1 Error: 0.0013988911652025"
## [1] "Average Test Accuracy: 0.966861014161742"

```

Similar results to training the initial model, accuracies in the 55-57 range roughly. Although there isn't much variance in any of the values. Right now

the model is almost just saying that all the rows are FALSE. Train accuracy was higher than test accuracy again. Train error rate for when CARAVAN=1 was extremely high, near 90 percent. Compared to when CARAVAN=0 which was 1 percent.

### 3 Improving Performance

I will now try to improve the performance of my model

I will start by trying to fine tune the ntree attribute of my model by testing my model with values between 1 and 100 for ntree. I was originally going to use the range 100 to 1000 and increment by 100 trees, but during testing I found that the lowest accuracies were within the 1 to 100 range and that after 100 they only increased. I have set the increment to 10 as I don't want to overfit the model.

I will write a function to do this, that will return the optimal number of trees based on test accuracy.

```
#Using same train and test set as before
#Tweak number of trees
testNTrees <- function(trainData,testData){
  ntrees<-10
  results<-NULL
  results<-data.frame(NTrees=as.numeric(),OOB=as.numeric(),trainFalseError=as.numeric(),trainTrueError=as.numeric())
  for (i in 1:10){
    trainData=train
    testData=test
    model<-randomForest(trainData[,,-ncol(trainData)],trainData[,ncol(trainData)],xtest=testData[,,-ncol(testData)],
      #TRAIN
    oob<-model$err.rate[nrow(model$test$err.rate),1,drop=FALSE]
    trainFalse<-model$err.rate[nrow(model$test$err.rate),2,drop=FALSE]
    trainTrue<-model$err.rate[nrow(model$test$err.rate),3,drop=FALSE]
    trainAccuracy<-sum(diag(model$confusion))/nrow(trainData)
    #TEST
    testError<-model$test$err.rate[nrow(model$test$err.rate),1,drop=FALSE]
    testFalse<-model$test$err.rate[nrow(model$test$err.rate),2,drop=FALSE]
    testTrue<-model$test$err.rate[nrow(model$test$err.rate),3,drop=FALSE]
    testAccuracy<-sum(diag(model$test$confusion))/nrow(testData)
    #Create new row in results with new data
    results[nrow(results)+1,]<-c(ntrees,oob,trainFalse,trainTrue,testError,testFalse,testTrue,testAccuracy)
    results
    ntrees <-ntrees + 10
  }
  #return max row
```

```

    ntrees<-results$NTrees[which.max(results$testAccuracy)]
    return(ntrees)
}

```

I will now use this function to find ntrees.

```

#Get ntrees
ntrees<-testNTrees(train,test)
ntrees
## [1] 70

```

During testing, ntree values were not high. Usually 100 or 200 was returned

I will not build a second model with the new values of ntrees to compare it to the original model

```

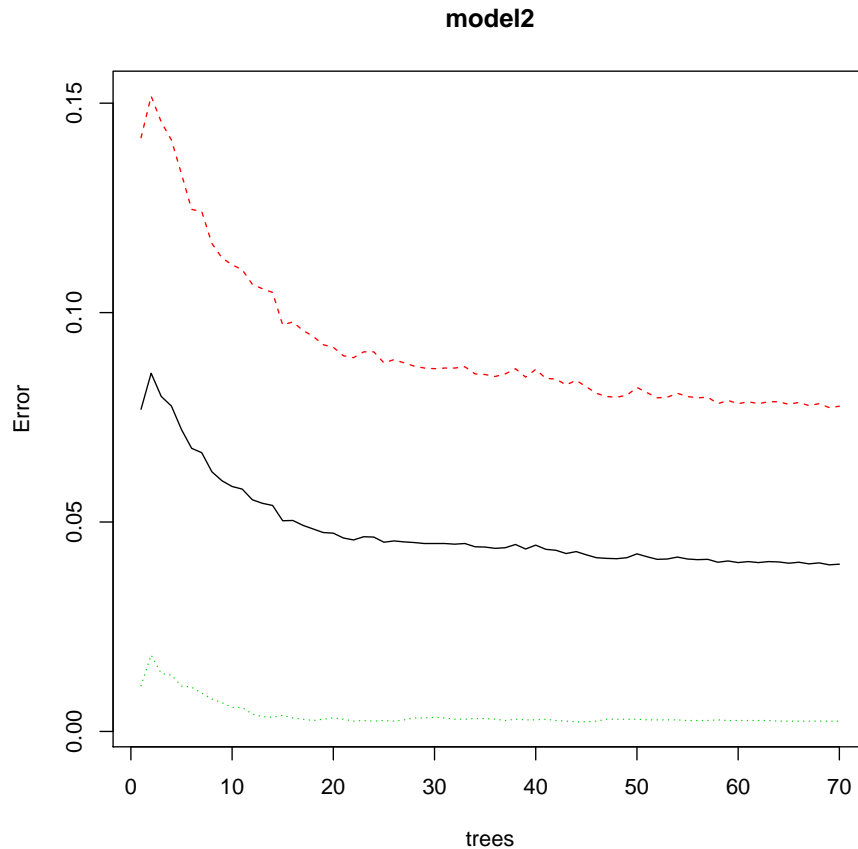
#Build second model with new ntree and nodesize
model2<-buildModel(train,test,ntrees=ntrees)
#Display Values
displayResultsFromModel(model2)

## [1] "TRAIN"
## [1] "Train OOB Error: 0.0399352401511063"
## [1] "Train CARAVAN=0 Error: 0.0776368697803897"
## [1] "Train CARAVAN=1 Error: 0.00245964642582629"

## Error in displayResultsFromModel(model2): argument "trainRows"
## is missing, with no default

#Plot Errors
plot(model2)

```



I will now validate the new model

```
#Validate second model, 10 fold cross validation
validateResults2<-validateModel(df,ntrees)

## Error: cannot allocate vector of size 2.1 Gb

displayResults(validateResults2)

## Error in print(results): object 'validateResults2' not found
```

During testing, the changes caused a greater range in test errors and accuracies. Values now varied greatly. Train errors were increased slightly but test errors were improved slightly also. Change was roughly 1 percent in each case. Error rates for CARAVAN=1 increased also and were near 90 percent again.

I will not try and fine tune the randomForest function variable nodesize. I will write a function to do this that works in a similar way to the function I wrote to test ntrees.

```

#Tweek Nodesize
testNodeSize <- function(trainData,testData,ntrees){
  nsize<-0
  results<-data.frame(Nodesize=as.numeric(),OOB=as.numeric(),trainFalseError=as.numeric(),trainTrueError=as.numeric())
  for (i in 1:floor(nrow(trainData)/100)){
    model<-randomForest(trainData[,,-ncol(trainData)],trainData[,ncol(trainData)],xtest=testData[,ncol(testData)])
    #TRAIN
    oob<-model$err.rate[nrow(model$test$err.rate),1,drop=FALSE]
    trainFalse<-model$err.rate[nrow(model$test$err.rate),2,drop=FALSE]
    trainTrue<-model$err.rate[nrow(model$test$err.rate),3,drop=FALSE]
    trainAccuracy<-sum(diag(model$confusion))/nrow(trainData)
    #TEST
    testError<-model$test$err.rate[nrow(model$test$err.rate),1,drop=FALSE]
    testFalse<-model$test$err.rate[nrow(model$test$err.rate),2,drop=FALSE]
    testTrue<-model$test$err.rate[nrow(model$test$err.rate),3,drop=FALSE]
    testAccuracy<-sum(diag(model$test$confusion))/nrow(testData)
    results[nrow(results)+1,]<-c(Nodesize=nsize,OOB=oob,trainFalseError=trainFalse,trainTrueError=trainTrue)
    nsize<-nsize+1
  }
  #Return node size
  nodeSize<-results$Nodesize[which.max(results$testAccuracy)]
  return(nodeSize)
}

```

I will now use the function to find nodesize.

```

#Get node size
nodeSize<-testNodeSize(train,test,ntrees)
nodeSize
## [1] 31

```

During testing, the values for nodeSize where in the low 100s. Train and test error rates have improved.

I will now test the accuracy of my new model, with the new values for ntree and nodesize.

```

#Build second model with new ntree and nodesize
model3<-buildModel(train,test,ntrees=ntrees,nodeSize=nodeSize)
#Display Values
displayResultsFromModel(model3)
## [1] "TRAIN"
## [1] "Train OOB Error: 0.0826459023976563"
## [1] "Train CARAVAN=0 Error: 0.138107021342406"
## [1] "Train CARAVAN=1 Error: 0.0275172943889316"
## Error in displayResultsFromModel(model3): argument "trainRows"

```



is missing, with no default

*#Plot Results*

I will now validate the new model

*#Validate second model, 10 fold cross validation*

```
validateResults3<-validateModel(df,ntrees,nodeSize)
displayResults(validateResults3)
```

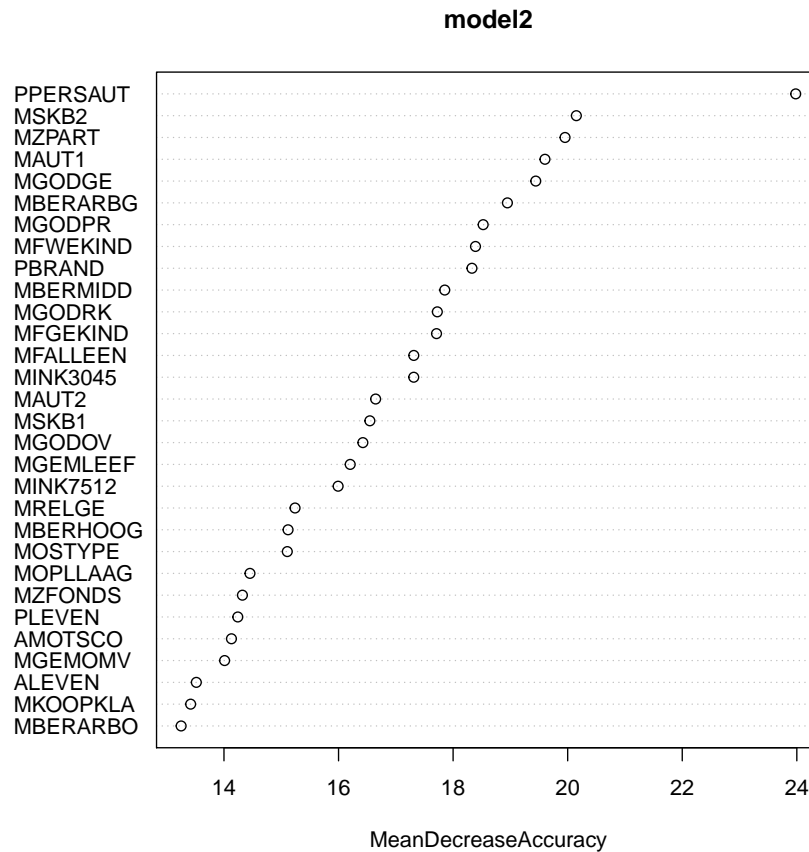
```
##          OOB trainFalseError trainTrueError  testError testFalseError
## 1  0.06326837      0.1115255      0.01530551 0.05882353      0.1125541
## 2  0.06404414      0.1146397      0.01375105 0.05561555      0.1007584
## 3  0.06194903      0.1115255      0.01267488 0.06422018      0.1147186
## 4  0.06524346      0.1172862      0.01351190 0.05453564      0.1007584
## 5  0.06117308      0.1120067      0.01064339 0.06310680      0.1222944
## 6  0.06314464      0.1145194      0.01207701 0.05723542      0.1018418
## 7  0.06212894      0.1123677      0.01219658 0.06206152      0.1103896
## 8  0.06128936      0.1107903      0.01207845 0.05234754      0.1018418
## 9  0.06416792      0.1165784      0.01207701 0.05828386      0.1114719
## 10 0.06332834      0.1106833      0.01626211 0.06260119      0.1147186
##      testTrueError trainAccuracy testAccuracy
## 1  0.005382131      0.9367316      0.9411765
## 2  0.010764263      0.9359559      0.9443844
## 3  0.013993541      0.9380510      0.9357798
## 4  0.008611410      0.9347565      0.9454644
## 5  0.004301075      0.9388269      0.9368932
## 6  0.012917115      0.9368554      0.9427646
## 7  0.013993541      0.9378711      0.9379385
## 8  0.003225806      0.9387106      0.9476525
## 9  0.005382131      0.9358321      0.9417161
## 10 0.010764263      0.9366717      0.9373988
## [1] "Average OOB: 0.0629737255822842"
## [1] "Average CARAVAN=0 Error: 0.113192285422822"
## [1] "Average CARAVAN=1 Error: 0.0130577869405343"
## [1] "Average Train Accuracy: 0.937026274417716"
## [1] " "
## [1] "Average Test Error: 0.0588831235732841"
## [1] "Average CARAVAN=0 Error: 0.109134761951663"
## [1] "Average CARAVAN=1 Error: 0.00893352778452956"
## [1] "Average Test Accuracy: 0.941116876426716"
```

Error rates and accuracies have improved. Train accuracy is not in the high 50s low 60 range. Same for test accuracy. Train accuracy is still higher than test accuracy. The biggest improvement is in the error rate for when CARAVAN = 1. Error rates are now in mid 60s low 70s range. This still isn't great but is a good improvement over the original model. Error rates for when CARAVAN=0

however have also increased by around 1-2 percent. This isn't a huge increase and was worth the drop in the error rates for CARAVAN=1.

Now I will use the importance function from the randomForest package, to create plots of mean decrease in accuracy and mean decrease in Gini. I will also create a list of each ordered highest to lowest which I will use to try and remove variables which are making my model less accurate.

```
#Mean Decrease in accuracy
meanDecreaseAccuracy<-importance(model2,type=1)
#Order highest to lowest
meanDecreaseAccuracy<-meanDecreaseAccuracy[order(-meanDecreaseAccuracy),,drop=FALSE]
#Plot
varImpPlot(model2,type=1)
```

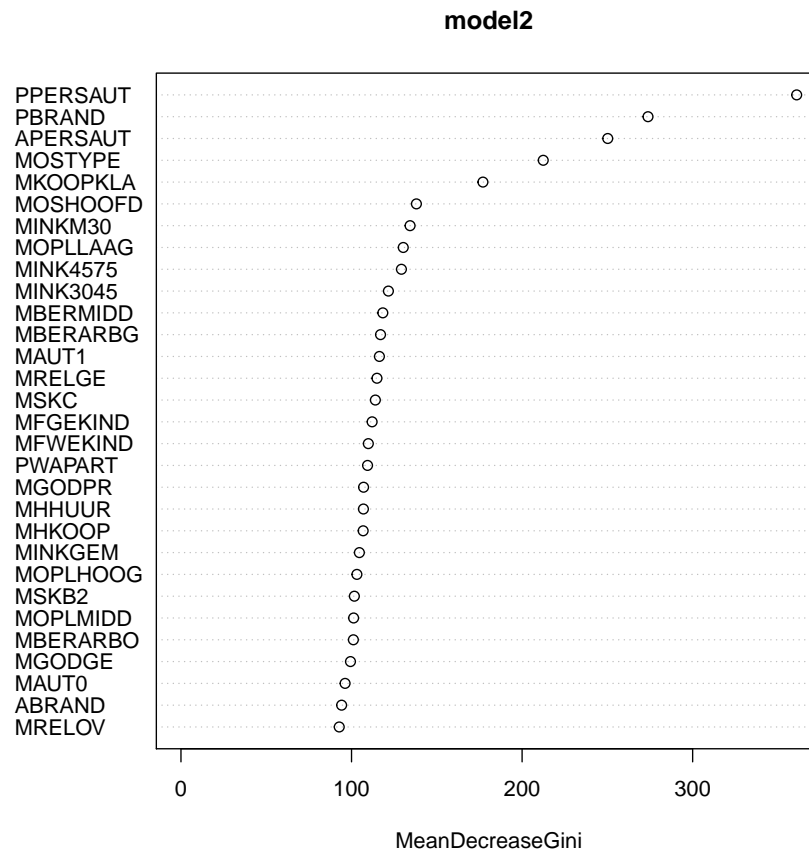


```
#Mean decrease in node impurity
meanDecreaseGini<-importance(model2,type=2)
```

```

#Order highest to lowest
meanDecreaseGini<-meanDecreaseGini[order(-meanDecreaseGini),,drop=FALSE]
#Plot
varImpPlot(model2,type=2)

```



I will remove any columns that have a negative mean decrease in accuracy value(if any).

```

#Get negative or 0 MDA
cols<-rownames(meanDecreaseAccuracy[meanDecreaseAccuracy<0,,drop=FALSE])
#Show cols being removed
cols

## NULL

#Remove cols
if (!is.null(cols)){

```

```

train<-train[,!(colnames(train) %in% cols)]
test<-test[,!(colnames(test) %in% cols)]
}

```

I will now test the accuracy of the final model

```

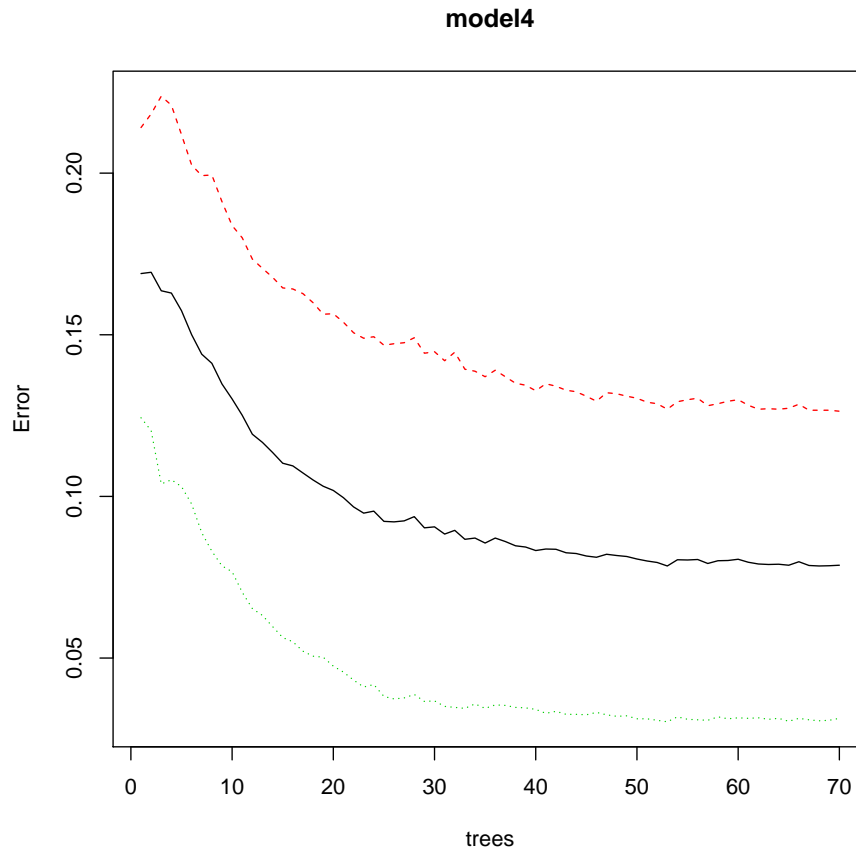
#Build final model, with removed columns based on mean decrease in accuracy
model4<-buildModel(train,test,ntrees=ntrees,nodeSize=nodeSize)
#Display values
displayResultsFromModel(model4)

## [1] "TRAIN"
## [1] "Train OOB Error: 0.0787140544291111"
## [1] "Train CARAVAN=0 Error: 0.126353232291989"
## [1] "Train CARAVAN=1 Error: 0.0313604919292852"

## Error in displayResultsFromModel(model4): argument "trainRows"
## is missing, with no default

#Plot Error
plot(model4)

```



I will now validate the final model

```
#Validate final model, 10 fold cross validation
validateResults4<-validateModel(df,ntrees,nodeSize)

## Error: cannot allocate vector of size 261.9 Mb

displayResults(validateResults4)

## Error in print(results): object 'validateResults4' not found
```

During testing, only one column was removed but this did cause an increase in accuracy of around 1-2 percent in both the training and testing accuracies. Error rate averages did slightly decrease, but there was a change in the range of values. There were now error rates in the 80s range for when CARAVAN=1. But the overall accuracy of the model did increase with this change.

## 4 Conclusions

My model was not very accurate, the real issue being the error rate in predicting when CARAVAN=1. More data would be helpful and I think a better resampling technique such as bootstrap might have helped. A different model such as a neural network might have performed better. By the final model, my model was around 60-62 percent accuracy and error rates for CARAVAN=1 were in the 60-80 percent range.

I did originally try and remove highly correlated variables, but didn't have much success so I omitted this from the coursework. If I had more time I would retry this to see if it could improve the accuracy.

If I had more time, I would have tried to combine some of the variables into new ones. I would have tried to combine the 5 social class variables for example into one variable, and classify each row as one social class. Due to time constraints I was not able to do this.

I feel like I have learned much completing this coursework despite the fact the accuracy of my model wasn't great.

## 5 References

- <http://ggplot2.org/>
- <https://cran.r-project.org/web/packages/dplyr/dplyr.pdf>
- <https://cran.r-project.org/web/packages/ROSE/ROSE.pdf>
- <https://www.analyticsvidhya.com/blog/2016/03/practical-guide-deal-imbalanced-classification-problems/>
- <https://cran.r-project.org/web/packages/randomForest/randomForest.pdf>
- <http://topepo.github.io/caret/index.html>