

Author: Fangjung (Kristy) Lin

Purpose: Application of Tree methods and Random forests in Telecom Industry (Case 4)

Introduction and overview

- Using Tree methods and Random forests to predict customer churn in telecom industry
Goal: To know more about what make the customer churn and predicting if he/she will churn or not.

Data set: The data was downloaded from IBM Sample Data Sets for customer retention programs. You can download it from here: <https://www.kaggle.com/blastchar/telco-customer-churn/download>

The data set includes information about:

- Customers who left within the last month – the column is called Churn
- Services that each customer has signed up for – phone, multiple lines, internet, online security, online backup, device protection, tech support, and streaming TV and movies
- Customer account information – how long they've been a customer, contract, payment method, paperless billing, monthly charges, and total charges
- Demographic info about customers – gender, age range, and if they have partners and dependents

Content of the analysis

Import necessary packages

Part1. Data Acquisition and Data Cleaning

Part2. Exploratory Data Analysis

Part3. Decision Tree

Part4. Random Forest

Part5. Conclusion

Import necessary packages

```
1 install.packages("tidyverse")
2 library(tidyverse)
3 install.packages("recipes")
4 library(recipes)
5 library(readr)
6 library(tidyverse)
7 library(modelr)
8 library(ggplot2)
9 install.packages("GGally")
10 library(GGally)
11 install.packages("corrplot")
12 library(corrplot)
13 library(gridExtra)
14 library(caret)
15 install.packages("e1071")
16 library("e1071")
17 install.packages("gplots")
18 install.packages("ROCR")
19 library(ROCR)
20 library(rpart)
21 install.packages('rpart.plot')
22 library(rpart.plot)
23 library(randomForest)
24 library(gplots)
25 install.packages('pROC')
26 library(pROC)
```

Part1. Data Acquisition and Data Cleaning

Data Acquisition

Walk through telecom data set

```
In [2]: 1 tele_churn <- read_csv("Desktop/churn.csv")
2 head(tele_churn)
```

customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	... DeviceProtection	TechSupport
7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	No	...	No
5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Yes	...	Yes
3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Yes	...	No
7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Yes	...	Yes
9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	No	...	No
9305-CDSKC	Female	0	No	No	8	Yes	Yes	Fiber optic	No	...	Yes

Explore the structure of data

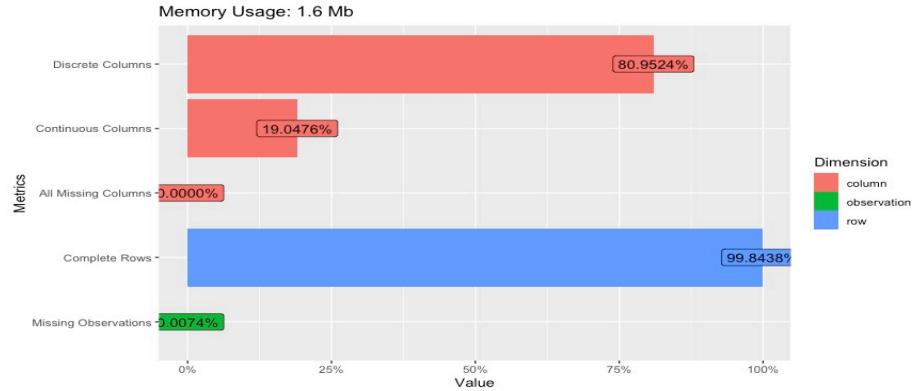
We can find that the data set contains 7043 observations and 21 variables

```
In [73]: 1 str(tele_churn)
tibble [7,043 x 21] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
$ customerID : chr [1:7043] "7590-VHVEG" "5575-GNVDE" "3668-QPYBK" "7795-CFOCW" ...
$ gender      : chr [1:7043] "Female" "Male" "Male" "Male" ...
$ SeniorCitizen: num [1:7043] 0 0 0 0 0 0 0 0 ...
$ Partner     : chr [1:7043] "Yes" "No" "No" "No" ...
$ Dependents  : chr [1:7043] "No" "No" "No" "No" ...
$ tenure      : num [1:7043] 1 34 2 45 2 8 22 10 28 62 ...
$ PhoneService: chr [1:7043] "No" "Yes" "Yes" "No" ...
$ MultipleLines: chr [1:7043] "No phone service" "No" "No" "No phone service" ...
$ InternetService: chr [1:7043] "DSL" "DSL" "DSL" "DSL" ...
$ OnlineSecurity: chr [1:7043] "No" "Yes" "Yes" "Yes" ...
$ OnlineBackup:  chr [1:7043] "Yes" "No" "Yes" "No" ...
$ DeviceProtection: chr [1:7043] "No" "Yes" "No" "Yes" ...
$ TechSupport:  chr [1:7043] "No" "No" "No" "Yes" ...
$ StreamingTV:  chr [1:7043] "No" "No" "No" "No" ...
$ StreamingMovies: chr [1:7043] "No" "No" "No" "No" ...
$ Contract:    chr [1:7043] "Month-to-month" "One year" "Month-to-month" "One year" ...
$ PaperlessBilling: chr [1:7043] "Yes" "No" "Yes" "No" ...
$ PaymentMethod:  chr [1:7043] "Electronic check" "Mailed check" "Mailed check" "Bank transfer (automatic)" ...
$ MonthlyCharges: num [1:7043] 29.9 57 53.9 42.3 70.7 ...
$ TotalCharges:  num [1:7043] 29.9 1889.5 108.2 1840.8 151.7 ...
$ Churn       : chr [1:7043] "No" "No" "Yes" "No" ...
```

Overview of the data

Let's take an overview on some aspects about the data. We can find that most of the data is discrete, there is not a lot of missing values, but we need to deal with them by removing or imputing and no columns missing completely.

```
In [6]: 1 install.packages("DataExplorer")
2 library(DataExplorer)
3 plot_intro(tele_churn)
```



Check the missing value

We are checking if this dataset has any missing values.

```
In [74]: 1 unlist(lapply(tele_churn,function(x) sum(is.na(x))))
```

```
customerID 0
gender 0
SeniorCitizen 0
Partner 0
Dependents 0
tenure 0
PhoneService 0
MultipleLines 0
InternetService 0
OnlineSecurity 0
OnlineBackup 0
DeviceProtection 0
TechSupport 0
StreamingTV 0
StreamingMovies 0
Contract 0
PaperlessBilling 0
PaymentMethod 0
MonthlyCharges 0
TotalCharges 11
Churn 0
```

I use lapply to check the number of missing values in each column. We found that there are 11 missing values in "TotalCharges" column. So, let's remove all rows with missing values.

Missing value deletion

In this step, I deleted the rows which have missing values.

```
In [75]: 1 tele_churn <- na.omit(tele_churn)
```

Part2. Exploratory Data Analysis

Let's see the overall summarize of our target variable

Summary of the data frame

```
In [9]: 1 summary(tele_churn)

  customerID      gender      SeniorCitizen      Partner
Length:7032    Length:7032    Min.   :0.0000    Length:7032
Class :character Class :character  1st Qu.:0.0000    Class :character
Mode  :character Mode  :character   Median :0.0000    Mode  :character
                           Mean   :0.1624
                           3rd Qu.:0.0000
                           Max.  :1.0000

  Dependents      tenure      PhoneService      MultipleLines
Length:7032    Min.   : 1.00    Length:7032    Length:7032
Class :character 1st Qu.: 9.00    Class :character  Class :character
Mode  :character Median :29.00    Mode  :character  Mode  :character
                           Mean   :32.42
                           3rd Qu.:55.00
                           Max.  :72.00

  InternetService  OnlineSecurity  OnlineBackup  DeviceProtection
Length:7032    Length:7032    Length:7032    Length:7032
Class :character Class :character  Class :character  Class :character
Mode  :character Mode  :character  Mode  :character  Mode  :character

  TechSupport      StreamingTV      StreamingMovies      Contract
Length:7032    Length:7032    Length:7032    Length:7032
Class :character Class :character  Class :character  Class :character
Mode  :character Mode  :character  Mode  :character  Mode  :character

  PaperlessBilling PaymentMethod      MonthlyCharges      TotalCharges
Length:7032    Length:7032    Min.   : 18.25    Min.   : 18.8
Class :character Class :character  1st Qu.: 35.59    1st Qu.: 401.4
Mode  :character Mode  :character  Median : 70.35    Median :1397.5
                           Mean   : 64.80    Mean   :2283.3
                           3rd Qu.: 89.86    3rd Qu.:3794.7
                           Max.  :118.75    Max.  :8684.8

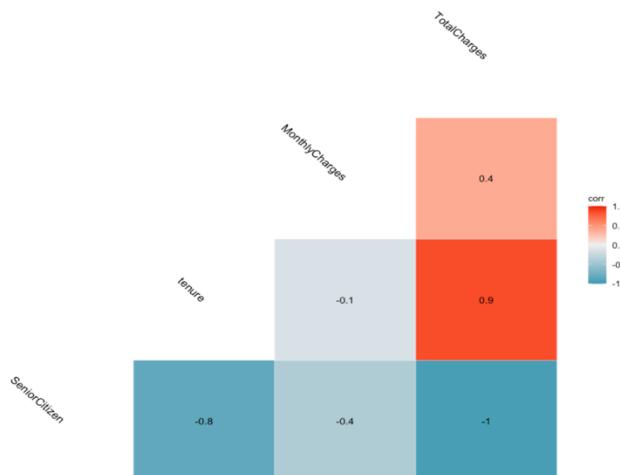
  Churn
Length:7032
Class :character
Mode  :character
```

From the summarize above, I find that the min/max difference between tenure is quite huge. The difference between MonthlyCharges and TotalCharges are quite huge. I need to have a deeper analysis in these variables to see if they can influence the customer churn.

Correlation among numeric variables

First, let's check the correlation between the numerical variables.

```
In [76]: 1 numeric.var <- sapply(tele_churn, is.numeric)
2 corr.matrix <- cor(tele_churn[,numeric.var])
3 options(repr.plot.width = 15, repr.plot.height = 8)
4 # corrplot(corr.matrix, main="\n\nCorrelation Plot for Numerical Variables", method = "number")
5
6 ggcormat(corr.matrix, name = "corr", label = TRUE, hjust = 1, label_size = 3.5,
7           angle = -45, size = 4)
```



- SeniorCitizen has a completely negative relation with TotalCharges. I assume that SeniorCitizen can be a factor that influence customer churn.
- tenure has a strong positive relation with totalcharges, which we need to take a deeper analysis to find if tenure has a positive relation with customer churn.

Creating new column tenure_bin and converting the tenure in years

```
In [77]: 1 #senior citizen is in integer form
2 tele_churn$SeniorCitizen <- as.factor(tele_churn$SeniorCitizen)
3
4 tele_churn <- mutate(tele_churn,tenure_bin=tenure)
5 tele_churn$tenure_bin[tele_churn$tenure_bin >= 0 & tele_churn$tenure_bin <= 12] <- "0 - 1 years"
6 tele_churn$tenure_bin[tele_churn$tenure_bin >= 13 & tele_churn$tenure_bin <= 24] <- "1 - 2 years"
7 tele_churn$tenure_bin[tele_churn$tenure_bin >= 25 & tele_churn$tenure_bin <= 36] <- "2 - 3 years"
8 tele_churn$tenure_bin[tele_churn$tenure_bin >= 37 & tele_churn$tenure_bin <= 48] <- "3 - 4 years"
9 tele_churn$tenure_bin[tele_churn$tenure_bin >= 49 & tele_churn$tenure_bin <= 60] <- "4 - 5 years"
10 tele_churn$tenure_bin[tele_churn$tenure_bin >= 61 & tele_churn$tenure_bin <= 72] <- "5 - 6 years"
11
12 #converting the newly created column into factor tenure_bin.
13 tele_churn$tenure_bin = as.factor(tele_churn$tenure_bin)
```

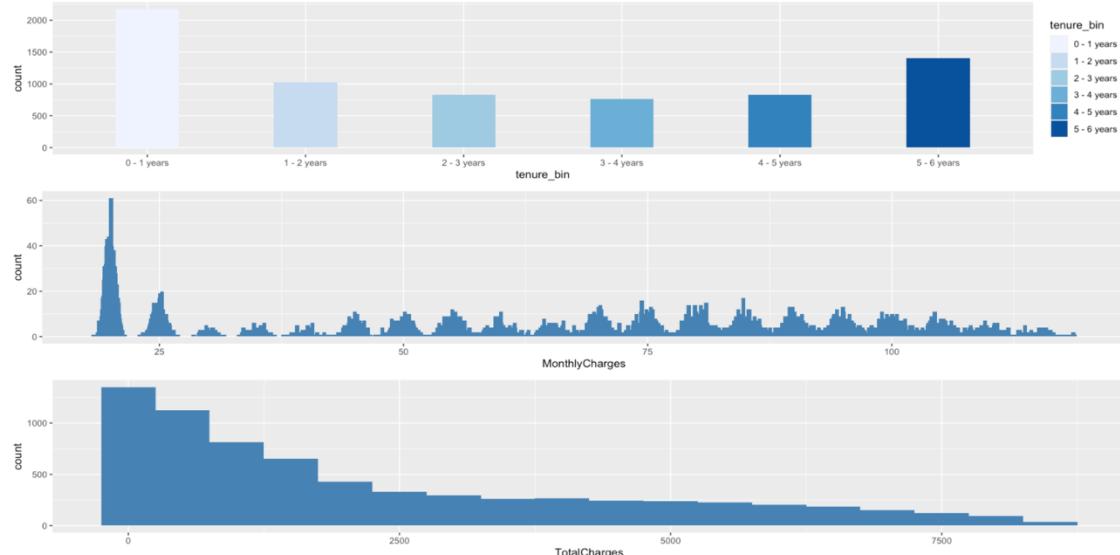
In the column MultipleLines modifying the No phone service -> No

```
In [78]: 1 tele_churn$MultipleLines[which(tele_churn$MultipleLines=="No phone service")]<- "No"
```

Graphical Representation of the Numerical Variables histograms

```
In [79]: 1 plot_A <- ggplot(tele_churn,aes(x = tenure_bin,fill=tenure_bin))+ geom_bar(width=0.4)+ scale_fill_brewer(palette="Blues")
2 plot_B <- ggplot(tele_churn,aes(x = MonthlyCharges))+ geom_bar(fill="steelblue",width=0.4)
3 plot_C <- ggplot(tele_churn,aes(x = TotalCharges))+ geom_histogram(fill="steelblue",binwidth =500)
4
5 options(repr.plot.width = 15, repr.plot.height = 8)
6 grid.arrange(plot_A,plot_B,plot_C)
```

Warning message:
 "position_stack requires non-overlapping x intervals"



From the tenure bin, we can find:

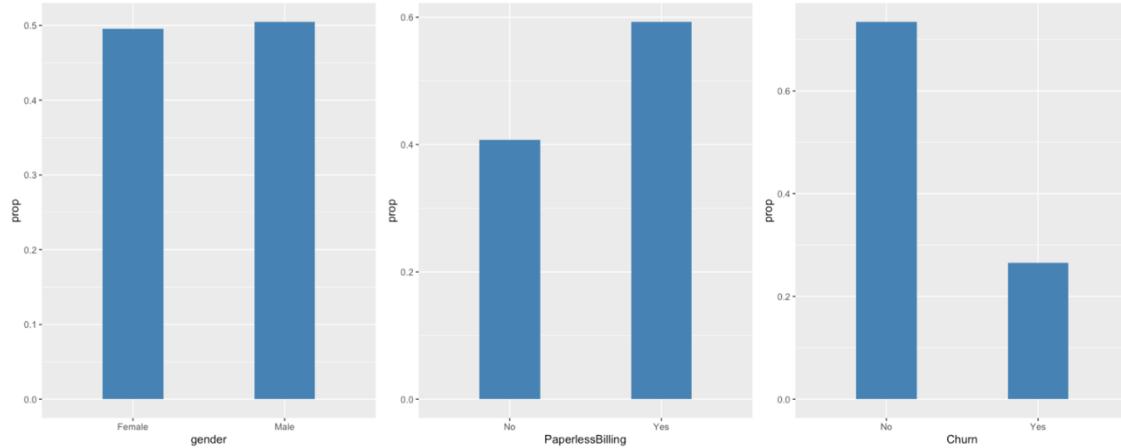
- Most people tenure year is between 0-1 years
- The trend of monthly charges and total charges is the similar.

Visualize the distribution of each variable

```
In [81]: 1 plot4 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = PaymentMethod,y=..prop..,group=2)
2                                         ,fill="Steelblue",stat='count',width=0.4)
3 plot5 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = gender,y=..prop..,group=2)
4                                         ,fill="Steelblue",stat='count',width=0.4)
5 plot6 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = SeniorCitizen,y=..prop..,group=2)
6                                         ,fill="Steelblue",stat='count',width=0.4)
7 plot7 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = Partner,y=..prop..,group=2)
8                                         ,fill="Steelblue",width=0.4,stat='count')
9 plot8 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = Dependents,y=..prop..,group=2)
10                                         ,fill="Steelblue",width=0.4,stat='count')
11 plot9 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = PhoneService,y=..prop..,group=2),
12                                         fill="Steelblue",width=0.4,stat='count')
13 plot10 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = MultipleLines,y=..prop..,group=2),
14                                         fill="Steelblue",width=0.4,stat='count')
15 plot11 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = InternetService,y=..prop..,group=2),
16                                         fill="Steelblue",width=0.4,stat='count')
17 plot12 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = OnlineSecurity,y=..prop..,group=2),
18                                         fill="Steelblue",width=0.4,stat='count')
19 plot13 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = OnlineBackup,y=..prop..,group=2),
20                                         fill="Steelblue",width=0.4,stat='count')
21 plot14 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = DeviceProtection,y=..prop..,group=2),
22                                         fill="Steelblue",width=0.4,stat='count')
23 plot15 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = TechSupport,y=..prop..,group=2),
24                                         fill="Steelblue",width=0.4,stat='count')
25 plot16 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = StreamingTV,y=..prop..,group=2),
26                                         fill="Steelblue",width=0.4,stat='count')
27 plot17 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = StreamingMovies,y=..prop..,group=2),
28                                         fill="Steelblue",width=0.4,stat='count')
29 plot18 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = Contract,y=..prop..,group=2),
30                                         fill="Steelblue",width=0.4,stat='count')
31 plot19 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = PaperlessBilling,y=..prop..,group=2),
32                                         fill="Steelblue",width=0.4,stat='count')
33 plot20 <- ggplot(data=tele_churn) + geom_bar(mapping = aes(x = Churn,y=..prop..,group=2),
34                                         fill="Steelblue",width=0.4,stat='count')
```

Gender/PaperlessBilling/Churn

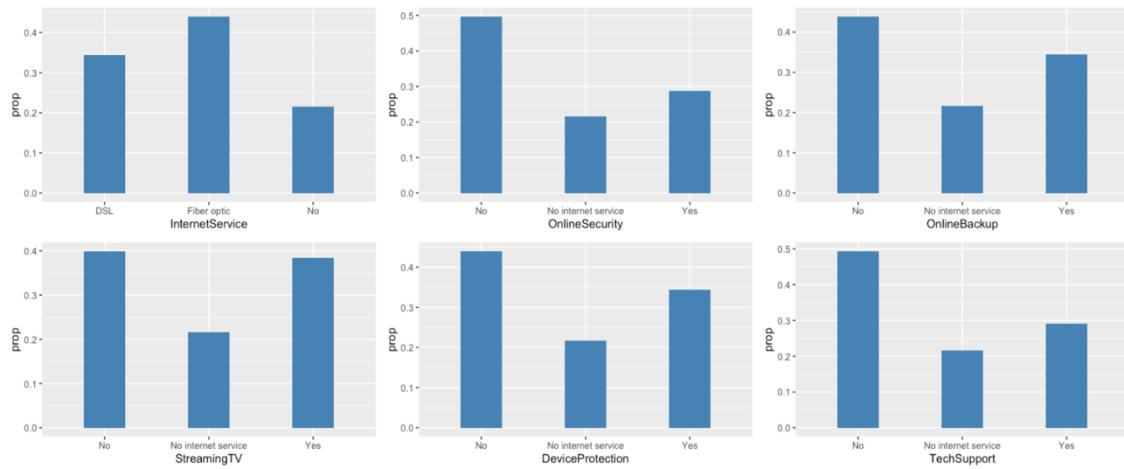
```
In [82]: 1 options(repr.plot.width = 15, repr.plot.height = 6)
2 grid.arrange(plot5,plot19,plot20,ncol=3)
```



- The data includes almost equal proportion of males and females.
- Almost 58% customers are on paperless billing.
- 26% of the customers have churned from the platform. From the distribution of customer churn, we can find that it's an imbalance data set.

InternetService/ OnlineSecurity/ OnlineBackup/ StreamingTV/ DeviceProtection/ TechSupport

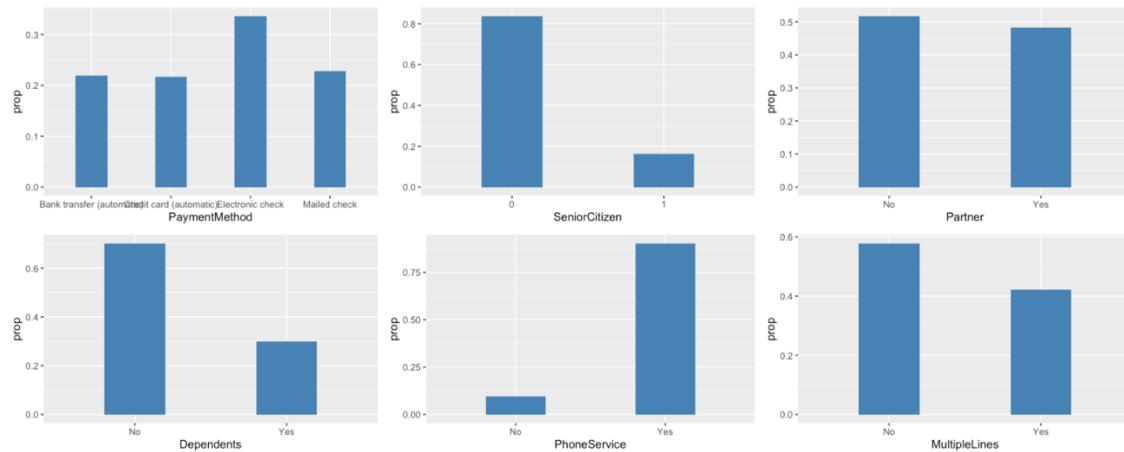
```
In [16]: 1 grid.arrange(plot11,plot12,plot13,plot16,plot14,plot15,nrow=2, ncol=3)
```



- Almost 40% of the customers have subscribed for the Fibre optic internet service.
- Almost 50% of the customers have no online security and almost 45% customers have no online backup.
- Almost 50% customers have no techsupport access and 40% have no streamingtv as a service.
- 45% of the customers have no service of device protection.

PaymentMethod/ SeniorCitizen/ Partner/ Dependents/ PhoneService/ MultipleLines

```
In [83]: 1 grid.arrange(plot4,plot6,plot7,plot8,plot9,plot10,nrow=2,ncol=3)
```



- Maximum number of customers have subscribed for electronic check for their payments.
- Very less i.e approx 20% of the customers are senior citizens.
- Equal number of customers with and without partners.
- 65% of the customers have no dependents.
- Almost 87% of the customers are with the phoneservice.

Summaries of Variables

In [84]:	<pre> 1 #Gender 2 cat("Gender") 3 type_counts1 <- table(tele_churn\$Gender) 4 type_counts1 / sum(type_counts1) * 100 5 6 #SeniorCitizen 7 cat("SeniorCitizen") 8 type_counts2 <- table(tele_churn\$SeniorCitizen) 9 type_counts2 / sum(type_counts2) * 100 10 11 #Partner 12 cat("Partner") 13 type_counts3 <- table(tele_churn\$Partner) 14 type_counts3 / sum(type_counts3) * 100 15 16 #Dependents 17 cat("SeniorCitizen") 18 type_counts4 <- table(tele_churn\$Dependents) 19 type_counts4 / sum(type_counts4) * 100 20 21 #PhoneService 22 cat("Dependents") 23 type_counts5 <- table(tele_churn\$PhoneService) 24 type_counts5 / sum(type_counts5) * 100 </pre>	<p>Gender</p> <table border="1"> <thead> <tr> <th>Female</th> <th>Male</th> </tr> </thead> <tbody> <tr> <td>49.53072</td> <td>50.46928</td> </tr> </tbody> </table> <p>SeniorCitizen</p> <table border="1"> <thead> <tr> <th>0</th> <th>1</th> </tr> </thead> <tbody> <tr> <td>83.75995</td> <td>16.24005</td> </tr> </tbody> </table> <p>Partner</p> <table border="1"> <thead> <tr> <th>No</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>51.74915</td> <td>48.25085</td> </tr> </tbody> </table> <p>SeniorCitizen</p> <table border="1"> <thead> <tr> <th>No</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>70.15074</td> <td>29.84926</td> </tr> </tbody> </table> <p>Dependents</p> <table border="1"> <thead> <tr> <th>No</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>9.67008</td> <td>90.32992</td> </tr> </tbody> </table>	Female	Male	49.53072	50.46928	0	1	83.75995	16.24005	No	Yes	51.74915	48.25085	No	Yes	70.15074	29.84926	No	Yes	9.67008	90.32992																						
Female	Male																																											
49.53072	50.46928																																											
0	1																																											
83.75995	16.24005																																											
No	Yes																																											
51.74915	48.25085																																											
No	Yes																																											
70.15074	29.84926																																											
No	Yes																																											
9.67008	90.32992																																											
In [85]:	<pre> 1 #MultipleLines 2 cat("MultipleLines") 3 type_counts6 <- table(tele_churn\$MultipleLines) 4 type_counts6 / sum(type_counts6) * 100 5 6 #InternetService 7 cat("InternetService") 8 type_counts7 <- table(tele_churn\$InternetService) 9 type_counts7 / sum(type_counts7) * 100 10 11 #OnlineSecurity 12 cat("OnlineSecurity") 13 type_counts8 <- table(tele_churn\$OnlineSecurity) 14 type_counts8 / sum(type_counts8) * 100 15 16 #OnlineBackup 17 cat("OnlineBackup") 18 type_counts9 <- table(tele_churn\$OnlineBackup) 19 type_counts9 / sum(type_counts9) * 100 20 21 #DeviceProtection 22 cat("DeviceProtection") 23 type_counts10 <- table(tele_churn\$DeviceProtection) 24 type_counts10 / sum(type_counts10) * 100 25 26 #TechSupport 27 cat("TechSupport") 28 type_counts11 <- table(tele_churn\$TechSupport) 29 type_counts11 / sum(type_counts11) * 100 </pre>	<p>MultipleLines</p> <table border="1"> <thead> <tr> <th>No</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>57.80717</td> <td>42.19283</td> </tr> </tbody> </table> <p>InternetService</p> <table border="1"> <thead> <tr> <th>DSL</th> <th>Fiber optic</th> <th>No</th> </tr> </thead> <tbody> <tr> <td>34.35722</td> <td>44.02730</td> <td>21.61547</td> </tr> </tbody> </table> <p>OnlineSecurity</p> <table border="1"> <thead> <tr> <th>No</th> <th>No</th> <th>internet service</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>49.72981</td> <td></td> <td>21.61547</td> <td>28.65472</td> </tr> </tbody> </table> <p>OnlineBackup</p> <table border="1"> <thead> <tr> <th>No</th> <th>No</th> <th>internet service</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>43.89932</td> <td></td> <td>21.61547</td> <td>34.48521</td> </tr> </tbody> </table> <p>DeviceProtection</p> <table border="1"> <thead> <tr> <th>No</th> <th>No</th> <th>internet service</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>43.99886</td> <td></td> <td>21.61547</td> <td>34.38567</td> </tr> </tbody> </table> <p>TechSupport</p> <table border="1"> <thead> <tr> <th>No</th> <th>No</th> <th>internet service</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>49.37429</td> <td></td> <td>21.61547</td> <td>29.01024</td> </tr> </tbody> </table>	No	Yes	57.80717	42.19283	DSL	Fiber optic	No	34.35722	44.02730	21.61547	No	No	internet service	Yes	49.72981		21.61547	28.65472	No	No	internet service	Yes	43.89932		21.61547	34.48521	No	No	internet service	Yes	43.99886		21.61547	34.38567	No	No	internet service	Yes	49.37429		21.61547	29.01024
No	Yes																																											
57.80717	42.19283																																											
DSL	Fiber optic	No																																										
34.35722	44.02730	21.61547																																										
No	No	internet service	Yes																																									
49.72981		21.61547	28.65472																																									
No	No	internet service	Yes																																									
43.89932		21.61547	34.48521																																									
No	No	internet service	Yes																																									
43.99886		21.61547	34.38567																																									
No	No	internet service	Yes																																									
49.37429		21.61547	29.01024																																									
In [86]:	<pre> 1 #StreamingTV 2 cat("StreamingTV") 3 type_counts12 <- table(tele_churn\$StreamingTV) 4 type_counts12 / sum(type_counts12) * 100 5 6 #Streaming Movies 7 cat("Streaming Movies") 8 type_counts13 <- table(tele_churn\$StreamingMovies) 9 type_counts13 / sum(type_counts13) * 100 10 11 #Contract 12 cat("Contract") 13 type_counts14 <- table(tele_churn\$Contract) 14 type_counts14 / sum(type_counts14) * 100 15 16 #PaperlessBilling 17 cat("PaperlessBilling") 18 type_counts15 <- table(tele_churn\$PaperlessBilling) 19 type_counts15 / sum(type_counts15) * 100 20 21 #PaymentMethod 22 cat("PaymentMethod") 23 type_counts16 <- table(tele_churn\$PaymentMethod) 24 type_counts16 / sum(type_counts16) * 100 25 26 #Churn 27 cat("Churn") 28 type_counts17 <- table(tele_churn\$Churn) 29 type_counts17 / sum(type_counts17) * 100 </pre>	<p>StreamingTV</p> <table border="1"> <thead> <tr> <th>No</th> <th>No</th> <th>internet service</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>39.94596</td> <td></td> <td>21.61547</td> <td>38.43857</td> </tr> </tbody> </table> <p>Streaming Movies</p> <table border="1"> <thead> <tr> <th>No</th> <th>No</th> <th>internet service</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>39.54778</td> <td></td> <td>21.61547</td> <td>38.83675</td> </tr> </tbody> </table> <p>Contract</p> <table border="1"> <thead> <tr> <th>Month-to-month</th> <th>One year</th> <th>Two year</th> </tr> </thead> <tbody> <tr> <td>55.10523</td> <td>20.93288</td> <td>23.96189</td> </tr> </tbody> </table> <p>PaperlessBilling</p> <table border="1"> <thead> <tr> <th>No</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>40.7281</td> <td>59.2719</td> </tr> </tbody> </table> <p>PaymentMethod</p> <table border="1"> <thead> <tr> <th>Bank transfer (automatic)</th> <th>Credit card (automatic)</th> <th>Electronic check</th> </tr> </thead> <tbody> <tr> <td>21.92833</td> <td>21.62969</td> <td>33.63197</td> </tr> <tr> <td>Mailed check</td> <td></td> <td></td> </tr> <tr> <td>22.81001</td> <td></td> <td></td> </tr> </tbody> </table> <p>Churn</p> <table border="1"> <thead> <tr> <th>No</th> <th>Yes</th> </tr> </thead> <tbody> <tr> <td>73.4215</td> <td>26.5785</td> </tr> </tbody> </table>	No	No	internet service	Yes	39.94596		21.61547	38.43857	No	No	internet service	Yes	39.54778		21.61547	38.83675	Month-to-month	One year	Two year	55.10523	20.93288	23.96189	No	Yes	40.7281	59.2719	Bank transfer (automatic)	Credit card (automatic)	Electronic check	21.92833	21.62969	33.63197	Mailed check			22.81001			No	Yes	73.4215	26.5785
No	No	internet service	Yes																																									
39.94596		21.61547	38.43857																																									
No	No	internet service	Yes																																									
39.54778		21.61547	38.83675																																									
Month-to-month	One year	Two year																																										
55.10523	20.93288	23.96189																																										
No	Yes																																											
40.7281	59.2719																																											
Bank transfer (automatic)	Credit card (automatic)	Electronic check																																										
21.92833	21.62969	33.63197																																										
Mailed check																																												
22.81001																																												
No	Yes																																											
73.4215	26.5785																																											

Dependent variable Churn vs Tenure_bin

```
In [87]: 1 options(repr.plot.width = 15, repr.plot.height = 6)
2 ggplot(tele_churn, aes(x = tenure_bin , y = MonthlyCharges)) + geom_point(aes(colour=factor(Churn)))
3 ggplot(tele_churn, aes(x = MonthlyCharges , y = TotalCharges)) + geom_point(aes(colour=factor(Churn)))
```



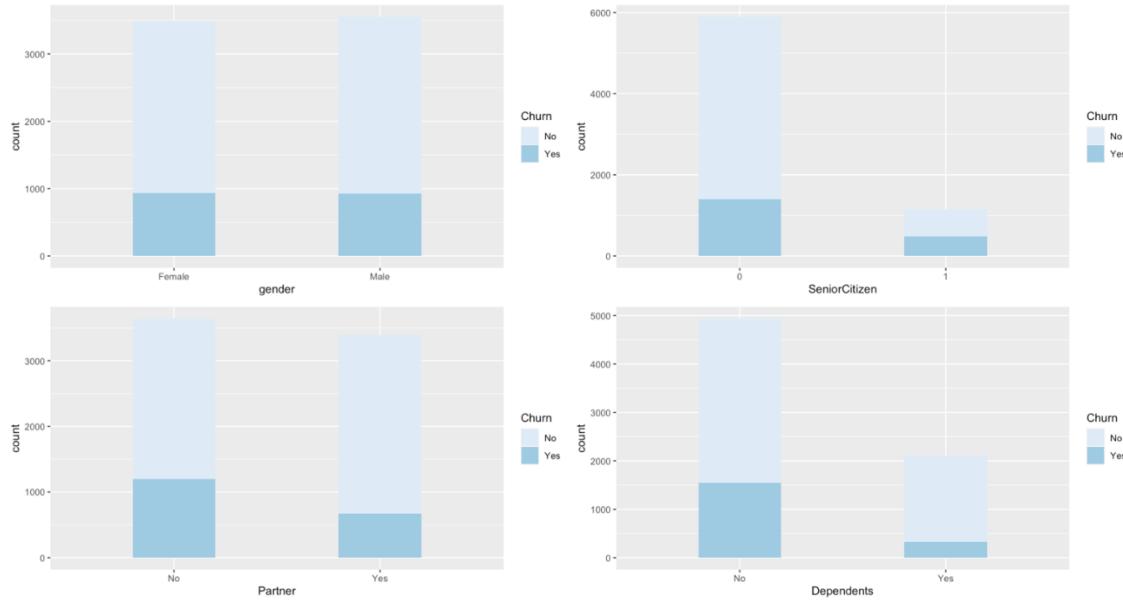
- Maximum Customers churned from the platform are the one having a tenure of 0-1 years.
- Maximum Churned customers have a Monthly charge more than \$65.

Relationship between the variables (dependent variable Churn vs Categorical Variables)

```
In [88]: 1 #gender vs churn
2 plot_relation_1 "ggplot(tele_churn, aes(x = gender,fill=Churn)) + geom_bar(width = 0.4)
3 *scale_fill_brewer(palette="Blues")
4 #SeniorCitizen vs Churn
5 plot_relation_2 "ggplot(tele_churn, aes(x = SeniorCitizen,fill=Churn)) + geom_bar(width = 0.4)
6 *scale_fill_brewer(palette="Blues")
7 #partner vs Churn
8 plot_relation_3 "ggplot(tele_churn, aes(x = Partner,fill=Churn)) + geom_bar(width = 0.4)
9 *scale_fill_brewer(palette="Blues")
10 #Dependents vs Churn
11 plot_relation_4 "ggplot(tele_churn, aes(x = Dependents,fill=Churn)) + geom_bar(width = 0.4)
12 *scale_fill_brewer(palette="Blues")
13 #PhoneService vs Churn
14 plot_relation_5 "ggplot(tele_churn, aes(x = PhoneService,fill=Churn)) + geom_bar(width = 0.4)
15 *scale_fill_brewer(palette="Blues")
16 #MultipleLines vs Churn
17 plot_relation_6 "ggplot(tele_churn, aes(x = MultipleLines,fill=Churn)) + geom_bar(width = 0.4)
18 *scale_fill_brewer(palette="Blues")
19 #InternetServices vs Churn
20 plot_relation_7 "ggplot(tele_churn, aes(x = InternetService,fill=Churn)) + geom_bar(width = 0.4)
21 *scale_fill_brewer(palette="Blues")
22 #OnlineSecurity vs Churn
23 plot_relation_8 "ggplot(tele_churn, aes(x = OnlineSecurity,fill=Churn)) + geom_bar(width = 0.4)
24 *scale_fill_brewer(palette="Blues")
25 #OnlineBackup vs Churn
26 plot_relation_9 "ggplot(tele_churn, aes(x = OnlineBackup,fill=Churn)) + geom_bar(width = 0.4)
27 *scale_fill_brewer(palette="Blues")
28 #DeviceProtection vs Churn
29 plot_relation_10 "ggplot(tele_churn, aes(x = DeviceProtection,fill=Churn)) + geom_bar(width = 0.4)
30 *scale_fill_brewer(palette="Blues")
31 #TechSupport vs Churn
32 plot_relation_11 "ggplot(tele_churn, aes(x = TechSupport,fill=Churn)) + geom_bar(width = 0.4)
33 *scale_fill_brewer(palette="Blues")
34 #StreamingTV vs Churn
35 plot_relation_12 "ggplot(tele_churn, aes(x = StreamingTV,fill=Churn)) + geom_bar(width = 0.4)
36 *scale_fill_brewer(palette="Blues")
37 #StreamingMovies vs Churn
38 plot_relation_13 "ggplot(tele_churn, aes(x = StreamingMovies,fill=Churn)) + geom_bar(width = 0.4)
39 *scale_fill_brewer(palette="Blues")
40 #Contract vs Churn
41 plot_relation_14 "ggplot(tele_churn, aes(x = Contract,fill=Churn)) + geom_bar(width = 0.4)
42 *scale_fill_brewer(palette="Blues")
43 #PaperlessBilling vs Churn
44 plot_relation_15 "ggplot(tele_churn, aes(x = PaperlessBilling,fill=Churn)) + geom_bar(width = 0.4)
45 *scale_fill_brewer(palette="Blues")
46 #PaymentMethod vs Churn
47 plot_relation_16 "ggplot(tele_churn, aes(x = PaymentMethod,fill=Churn)) + geom_bar(width = 0.4)
48 *scale_fill_brewer(palette="Blues")
```

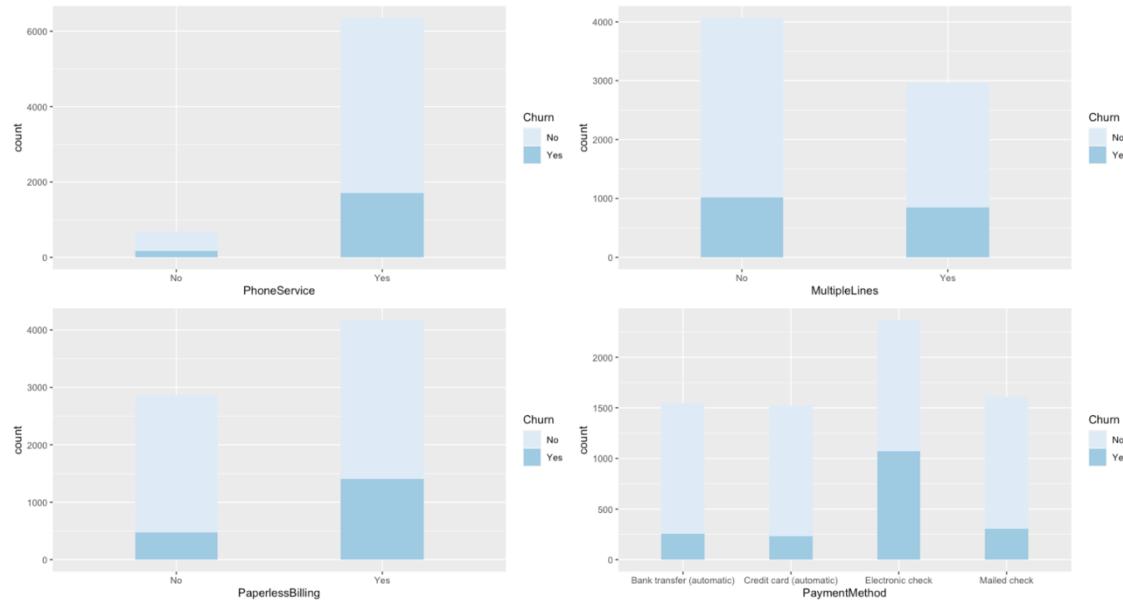
Visualize the result

```
In [89]: 1 options(repr.plot.width = 15, repr.plot.height = 8)
2 grid.arrange(plot_relation_1,plot_relation_2,plot_relation_3,plot_relation_4)
```

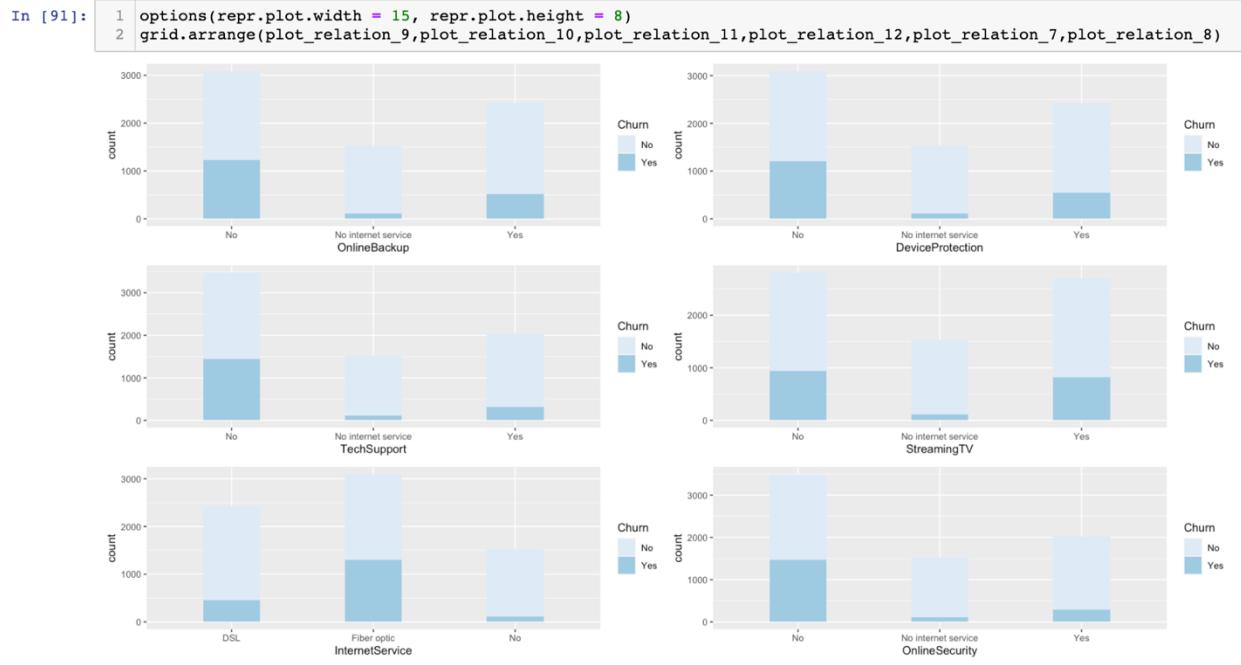


- We can see clearly that there is no specific gender correlated with Churned or not churned both genders have the same distribution for both 0 and 1
- The fraction of churned out of all SeniorCitizen is much bigger than the fraction of churned out of all non SeniorCitizen so being SeniorCitizen is indeed an important factor of churning
- Half the people have Partner and they have less probability of Churning than people who don't have as we can see the fraction of churned in Partner group is less than no Partner group
- Churn rate is more among the customers with no dependents. The ratio of churned people who don't have Dependents is 31%, in the people who have dependent is 15% only which is half the ratio which means that not having dependent increase the probability of churn

```
In [90]: 1 options(repr.plot.width = 15, repr.plot.height = 8)
2 grid.arrange(plot_relation_5,plot_relation_6,plot_relation_15,plot_relation_16)
```



- Churn rate is more with customers having phone service.
- Churn rate is more with customers having paperless billing.
- Churn rate is more with customers having electronic check as the payment mode.



- Churn rate is more with customers having month to month contract.
- Churn rate is more with customers having no online security and techsupport.
- Churn rate is almost equal among the subscribers with or without the streamingtv.

Part3. Decision Tree

A decision tree can give a nice and easy visualization of the prediction rules. More specifically, rpart classification tree will be grown. No tuning, just a simple tree.

Data frame manipulation

```
In [92]: 1 tele_churn$churn_number <- 0
          2 tele_churn$churn_number[tele_churn$Churn == 'Yes'] <- 1
          3 tele_churn$Churn <- NULL
```

Convert categorical variables to factor

```
In [93]: 1 tele_churn$gender <- as.factor(tele_churn$gender)
          2 tele_churn$SeniorCitizen <- as.factor(tele_churn$SeniorCitizen)
          3 tele_churn$Partner <- as.factor(tele_churn$Partner)
          4 tele_churn$Dependents <- as.factor(tele_churn$Dependents)
          5 tele_churn$PhoneService <- as.factor(tele_churn$PhoneService)
          6 tele_churn$MultipleLines <- as.factor(tele_churn$MultipleLines)
          7 tele_churn$InternetService <- as.factor(tele_churn$InternetService)
          8 tele_churn$OnlineSecurity <- as.factor(tele_churn$OnlineSecurity)
          9 tele_churn$OnlineBackup <- as.factor(tele_churn$OnlineBackup)
         10 tele_churn$DeviceProtection <- as.factor(tele_churn$DeviceProtection)
         11 tele_churn$TechSupport <- as.factor(tele_churn$TechSupport)
         12 tele_churn$StreamingTV <- as.factor(tele_churn$StreamingTV)
         13 tele_churn$StreamingMovies <- as.factor(tele_churn$StreamingMovies)
         14 tele_churn$Contract <- as.factor(tele_churn$Contract)
         15 tele_churn$PaperlessBilling <- as.factor(tele_churn$PaperlessBilling)
         16 tele_churn$PaymentMethod <- as.factor(tele_churn$PaymentMethod)
```

Remove unnecessary columns for modeling

We don't need the customerID and tenure (tenure from the dataset has already transformed the tenure values in tenure_bin).

```
In [94]: 1 tele_churn$customerID <- NULL
          2 tele_churn$tenure <- NULL
```

Creating dummy variables

```
In [95]: 1 trainDummy <- dummyVars('~.', data = tele_churn, fullRank = F)
```

```
In [96]: 1 train <- as.data.frame(predict(trainDummy, tele_churn))
          2 colnames(train)
```

```
'gender.Female' 'gender.Male' 'SeniorCitizen.0' 'SeniorCitizen.1' 'Partner.No' 'Partner.Yes' 'Dependents.No' 'Dependents.Yes'
'PhoneService.No' 'PhoneService.Yes' 'MultipleLines.No' 'MultipleLines.Yes' 'InternetService.DSL' 'InternetService.Fiber optic' 'InternetService.No'
'OnlineSecurity.No' 'OnlineSecurity.No internet service' 'OnlineSecurity.Yes' 'OnlineBackup.No' 'OnlineBackup.No internet service'
'OnlineBackup.Yes' 'DeviceProtection.No' 'DeviceProtection.No internet service' 'DeviceProtection.Yes' 'TechSupport.No'
'TechSupport.No internet service' 'TechSupport.Yes' 'StreamingTV.No' 'StreamingTV.No internet service' 'StreamingTV.Yes' 'StreamingMovies.No'
'StreamingMovies.No internet service' 'StreamingMovies.Yes' 'Contract.Month-to-month' 'Contract.One year' 'Contract.Two year'
'PaperlessBilling.No' 'PaperlessBilling.Yes' 'PaymentMethod.Bank transfer (automatic)' 'PaymentMethod.Credit card (automatic)'
'PaymentMethod.Electronic check' 'PaymentMethod.Mailed check' 'MonthlyCharges' 'TotalCharges' 'tenure_bin.0 - 1 years' 'tenure_bin.1 - 2 years'
'tenure_bin.2 - 3 years' 'tenure_bin.3 - 4 years' 'tenure_bin.4 - 5 years' 'tenure_bin.5 - 6 years' 'churn_number'
```

```
In [97]: 1 train$churn_number <- as.factor(ifelse(train$churn_number == 1, 'yes', 'no'))
```

Make sure the structure of data

```
In [99]: 1 str(tele_churn)

tibble [7,032 x 20] (S3: tbl_df/tbl/data.frame)
$ gender      : Factor w/ 2 levels "Female", "Male": 1 2 2 2 1 1 2 1 1 2 ...
$ SeniorCitizen : Factor w/ 2 levels "0", "1": 1 1 1 1 1 1 1 1 1 1 ...
$ Partner      : Factor w/ 2 levels "No", "Yes": 2 1 1 1 1 1 1 1 2 1 ...
$ Dependents   : Factor w/ 2 levels "No", "Yes": 1 1 1 1 1 1 2 1 1 2 ...
$ PhoneService  : Factor w/ 2 levels "No", "Yes": 1 2 2 1 2 2 2 1 2 2 ...
$ MultipleLines : Factor w/ 2 levels "No", "Yes": 1 1 1 1 1 2 2 1 2 1 ...
$ InternetService: Factor w/ 3 levels "DSL", "Fiber optic", ...: 1 1 1 1 2 2 2 1 2 1 ...
$ OnlineSecurity: Factor w/ 3 levels "No", "No internet service", ...: 1 3 3 3 1 1 1 3 1 3 ...
$ OnlineBackup   : Factor w/ 3 levels "No", "No internet service", ...: 3 1 3 1 1 1 3 1 1 3 ...
$ DeviceProtection: Factor w/ 3 levels "No", "No internet service", ...: 1 3 1 3 1 3 1 1 3 1 ...
$ TechSupport    : Factor w/ 3 levels "No", "No internet service", ...: 1 1 1 3 1 1 1 1 3 1 ...
$ StreamingTV    : Factor w/ 3 levels "No", "No internet service", ...: 1 1 1 1 1 3 3 1 3 1 ...
$ StreamingMovies: Factor w/ 3 levels "No", "No internet service", ...: 1 1 1 1 1 3 1 3 1 ...
$ Contract      : Factor w/ 3 levels "Month-to-month", ...: 1 2 1 2 1 1 1 1 2 ...
$ PaperlessBilling: Factor w/ 2 levels "No", "Yes": 2 1 2 1 2 2 2 1 2 1 ...
$ PaymentMethod   : Factor w/ 4 levels "Bank transfer (automatic)", ...: 3 4 4 1 3 3 2 4 3 1 ...
$ MonthlyCharges : num [1:7032] 29.9 57 53.9 42.3 70.7 ...
$ TotalCharges   : num [1:7032] 29.9 1889.5 108.2 1840.8 151.7 ...
$ tenure_bin     : Factor w/ 6 levels "0 - 1 years", ...: 1 3 1 4 1 1 2 1 3 6 ...
$ churn_number   : num [1:7032] 0 0 1 0 1 1 0 0 1 0 ...
- attr(*, "na.action") = 'omit' Named int [1:11] 489 754 937 1083 1341 3332 3827 4381 5219 6671 ...
-- attr(*, "names") = chr [1:11] "489" "754" "937" "1083" ...
```

Splitting Training and Test Data

Data splitting into test and train into 70/30 ratio.

```
In [100]: 1 set.seed(123)
2 split <- sample(2, nrow(train), replace = T, prob = c(0.7, 0.3))
3 traindf <- train[split == 1,]
4 testdf <- train[split == 2,]
5
6 dim(traindf)
7 dim(testdf)
```

4943 51
2089 51

Check whether any NA value exists or not

```
In [102]: 1 anyNA(tele_churn)
2 anyNA(traindf)
3 anyNA(testdf)
```

FALSE
FALSE
FALSE

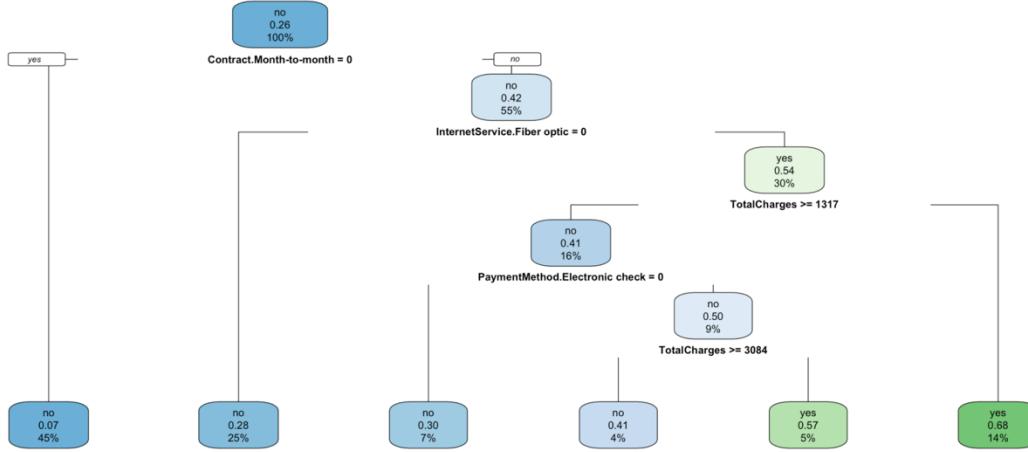
Train the model

I set up the method = class to make it a binary classification.

```
In [33]: 1 set.seed(123)
2 tree_df <- rpart(churn_number ~ ., data = traindf, method = "class", parms = list(split = "gini"))
```

Visualize the model

```
In [34]: 1 options(repr.plot.width = 18, repr.plot.height = 8)
2 rpart.plot(tree_df)
```



I noticed here that the model didn't use a lot of the features, so this is an indicator that feature engineering is important in the problem

Predict the value

```
In [35]: 1 predict(tree_df, data = traindf, type = "class") -> traintree_pred1
2 predict(tree_df, data = traindf, type = "prob") -> traintree_prob1
3 predict(tree_df, newdata = testdf, type = "class") -> testtree_pred1
4 predict(tree_df, newdata = testdf, type = "prob") -> testtree_prob1
```

Confusion Matrix of Training Set

```
In [36]: 1 confusionMatrix(data = traintree_pred1, reference = traindf$churn_number)

Confusion Matrix and Statistics

          Reference
Prediction   no  yes
      no    3319  680
      yes   324   620

Accuracy : 0.7969
95% CI : (0.7854, 0.808)
No Information Rate : 0.737
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.4255

McNemar's Test P-Value : < 2.2e-16

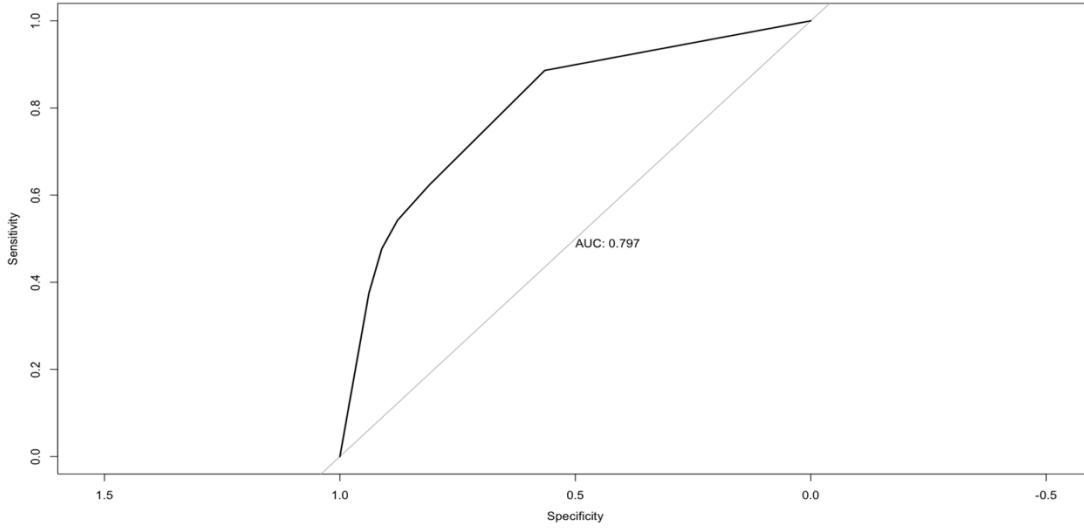
Sensitivity : 0.9111
Specificity : 0.4769
Pos Pred Value : 0.8300
Neg Pred Value : 0.6568
Prevalence : 0.7370
Detection Rate : 0.6715
Detection Prevalence : 0.8090
Balanced Accuracy : 0.6940

'Positive' Class : no
```

AUC-ROC Curve of Training Set

```
In [37]: 1 options(repr.plot.width = 15, repr.plot.height = 8)
2 traintree_actual <- ifelse(traindf$churn_number == "yes", 1,0)
3 roc <- roc(traintree_actual, traintree_prob[,2], plot= TRUE, print.auc=TRUE)
```

Setting levels: control = 0, case = 1
Setting direction: controls < cases



Confusing Matrix of Testing Set

```
In [38]: 1 confusionMatrix(data = testtree_pred1, reference = testdf$churn_number)
```

Confusion Matrix and Statistics

		Reference	
		no	yes
Prediction	no	1383	300
	yes	137	269

Accuracy : 0.7908
95% CI : (0.7727, 0.8081)
No Information Rate : 0.7276
P-Value [Acc > NIR] : 1.597e-11

Kappa : 0.4203

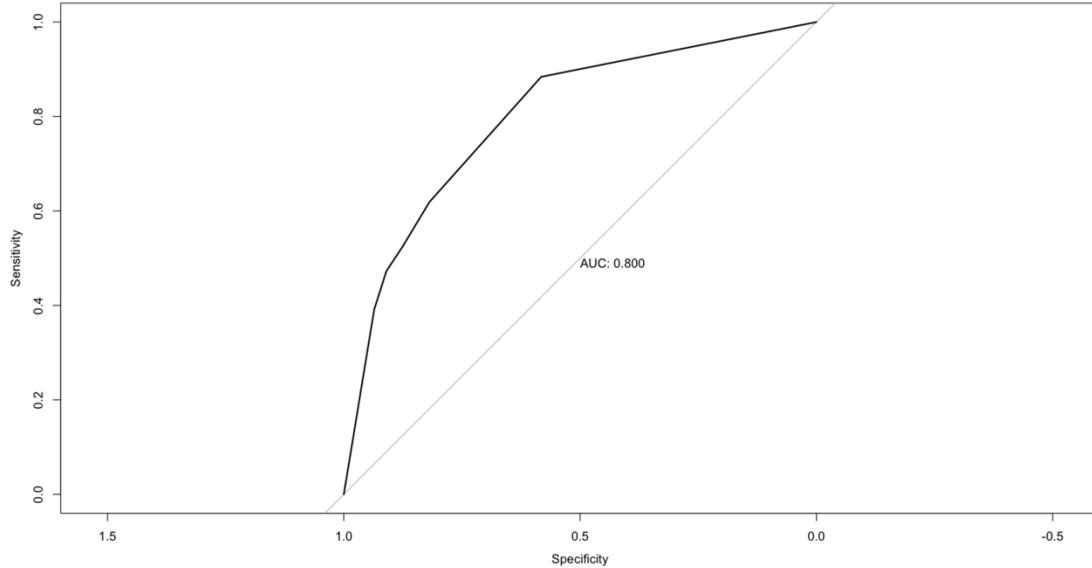
McNemar's Test P-Value : 9.225e-15

Sensitivity : 0.9099
Specificity : 0.4728
Pos Pred Value : 0.8217
Neg Pred Value : 0.6626
Prevalence : 0.7276
Detection Rate : 0.6620
Detection Prevalence : 0.8056
Balanced Accuracy : 0.6913

'Positive' Class : no

AUC-ROC Curve of Testing Set

```
In [39]: 1 testtree_actual <- ifelse(testdf$churn_number == "yes", 1,0)
2 roc <- roc(testtree_actual, testtree_prob[,2], plot = TRUE, print.auc = TRUE)
Setting levels: control = 0, case = 1
Setting direction: controls < cases
```



For the training set, the Accuracy is 0.7969 and the AUC is 0.797. For the testing set, the Accuracy is 0.7908 and the AUC is 0.80.

Part4. Random Forest

Data Preparation - In this part, I use the same data prepared for Classification Tree models to see if I can get a higher accuracy.

Train the model

```
In [40]: 1 names(traindf) <- make.names(names(traindf))
2 names(testdf) <- make.names(names(testdf))
3 rf_df <- randomForest(as.factor(churn_number) ~ ., data=traindf, importance = TRUE, ntree=500, do.trace=FALSE )
```

Predict the value

```
In [41]: 1 predict(rf_df, data = traindf, type = "class") -> trainrf_pred1
2 predict(rf_df, data = traindf, type = "prob") -> trainrf_prob1
3 predict(rf_df, newdata = testdf, type = "class") -> testrf_pred1
4 predict(rf_df, newdata = testdf, type = "prob") -> testrf_prob1
```

Confusion Matrix of training set

```
In [42]: 1 confusionMatrix(data = trainrf_pred1, reference = traindf$churn_number)
          Confusion Matrix and Statistics

          Reference
          Prediction   no   yes
                  no 3278  654
                  yes  365  646

          Accuracy : 0.7938
          95% CI : (0.7823, 0.8051)
          No Information Rate : 0.737
          P-Value [Acc > NIR] : < 2.2e-16

          Kappa : 0.4273

          Mcnemar's Test P-Value : < 2.2e-16

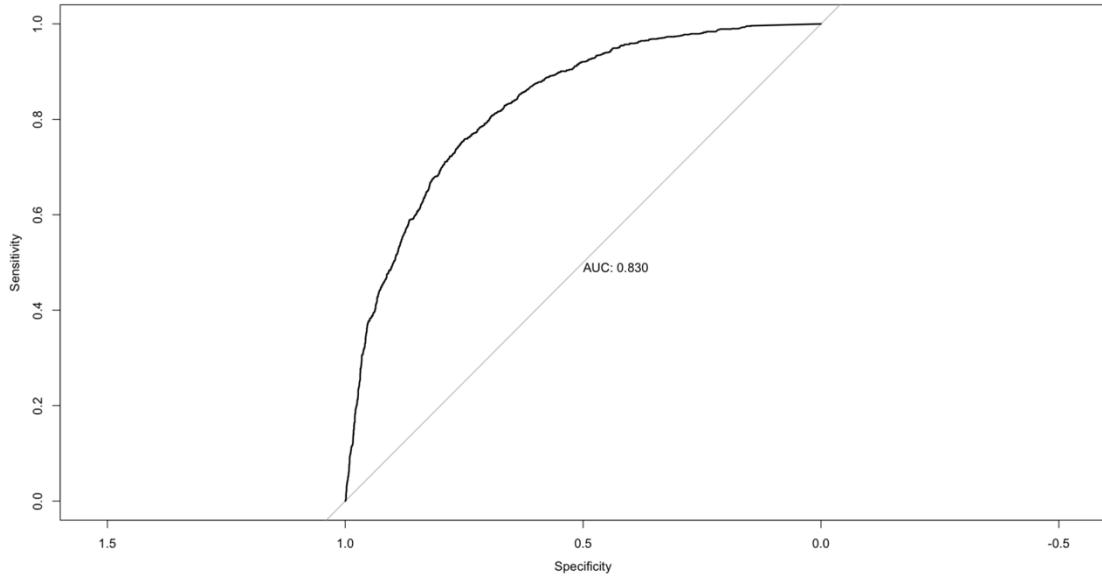
          Sensitivity : 0.8998
          Specificity : 0.4969
          Pos Pred Value : 0.8337
          Neg Pred Value : 0.6390
          Prevalence : 0.7370
          Detection Rate : 0.6632
          Detection Prevalence : 0.7955
          Balanced Accuracy : 0.6984

          'Positive' Class : no
```

AUC-ROC Curve of training set

```
In [43]: 1 options(repr.plot.width = 15, repr.plot.height = 8)
 2 rf_df_actual <- ifelse(traindf$churn_number == "yes", 1,0)
 3 roc <- roc(rf_df_actual, trainrf_prob[,2], plot= TRUE, print.auc=TRUE)

Setting levels: control = 0, case = 1
Setting direction: controls < cases
```



Confusion Matrix of testing set

```
In [44]: 1 confusionMatrix(data = testrf_pred1, reference = testdf$churn_number)

Confusion Matrix and Statistics

          Reference
Prediction    no  yes
      no   1370  291
      yes   150  278

      Accuracy : 0.7889
      95% CI : (0.7708, 0.8062)
      No Information Rate : 0.7276
      P-Value [Acc > NIR] : 6.410e-11

      Kappa : 0.4227

      Mcnemar's Test P-Value : 2.617e-11

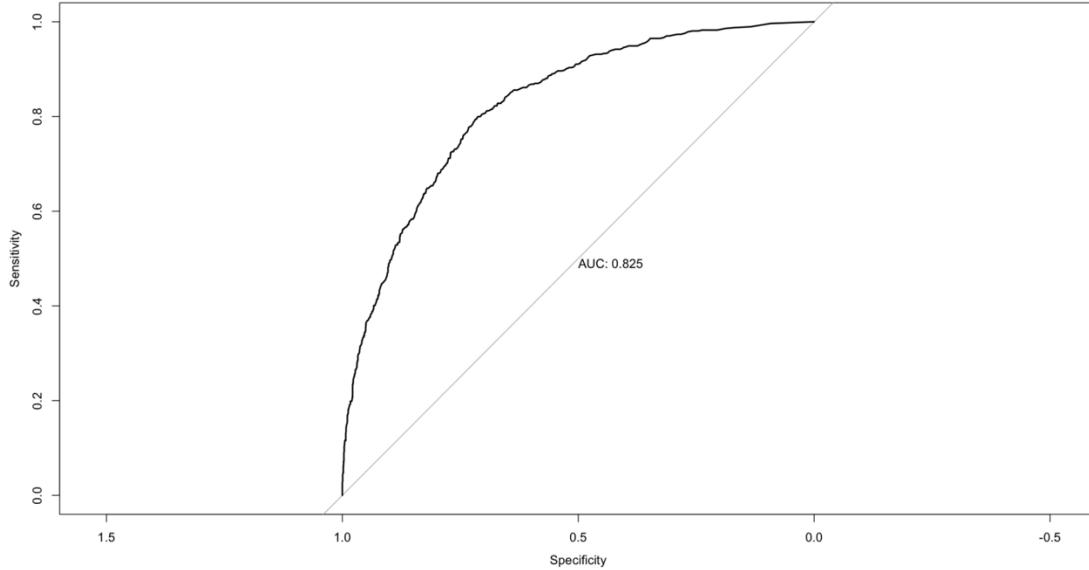
      Sensitivity : 0.9013
      Specificity : 0.4886
      Pos Pred Value : 0.8248
      Neg Pred Value : 0.6495
      Prevalence : 0.7276
      Detection Rate : 0.6558
      Detection Prevalence : 0.7951
      Balanced Accuracy : 0.6949

      'positive' Class : no
```

AUC-ROC Curve of testing set

```
In [45]: 1 options(repr.plot.width = 15, repr.plot.height = 8)
2 rf_df_actual <- ifelse(testdf$churn_number == "yes", 1,0)
3 roc <- roc(rf_df_actual, testrf_probl[,2], plot= TRUE, print.auc=TRUE)

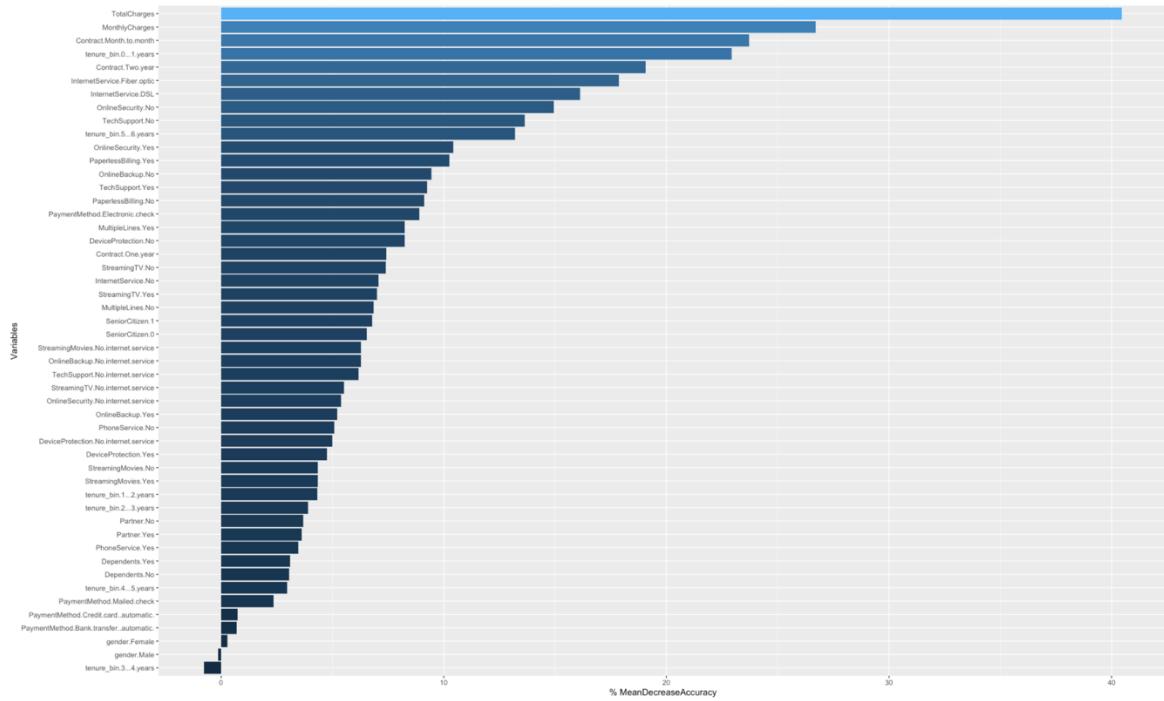
Setting levels: control = 0, case = 1
Setting direction: controls < cases
```



The training random forest model has the accuracy of 0.7957 and AUC of 0.829 for the training set. The testing random forest model has the Accuracy of 0.7927 and AUC of 0.828 for the testing set.

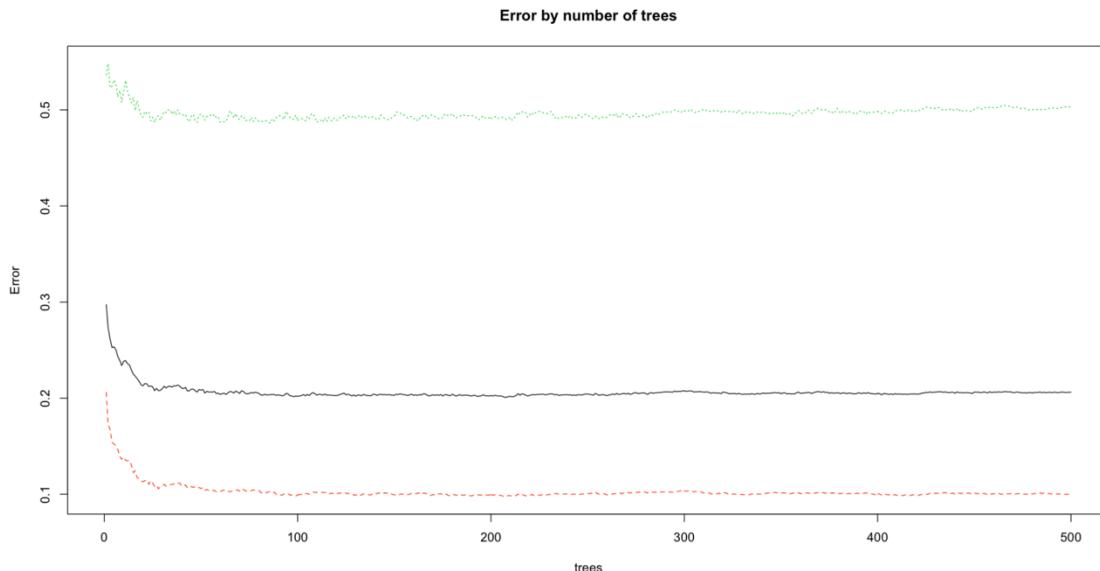
Plot the feature importance

```
In [46]:  
1 options(repr.plot.width = 20, repr.plot.height = 12)  
2 imp_rf_df <- importance(rf_df)  
3 imp_DF <- data.frame(Variables = row.names(imp_rf_df), MeanDecreaseAccuracy = imp_rf_df[,3])  
4 imp_DF <- imp_DF[order(imp_DF$MeanDecreaseAccuracy, decreasing = TRUE),]  
5  
6 ggplot(imp_DF, aes(x=reorder(Variables, MeanDecreaseAccuracy), y=MeanDecreaseAccuracy, fill=MeanDecreaseAccuracy))
```



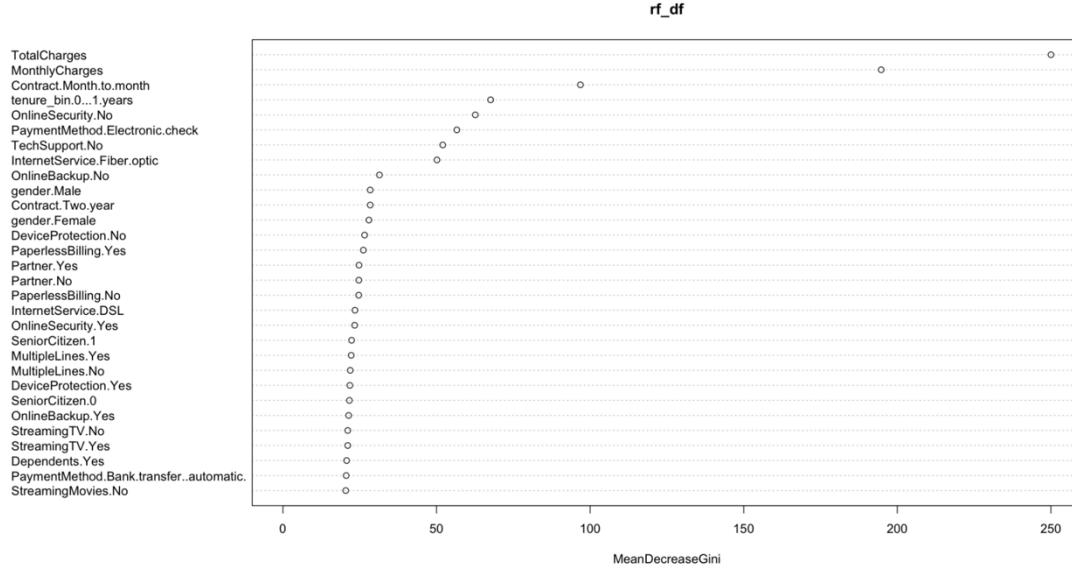
According to the Variable Importance plot, TotalCharges, MonthlyCharges, Tenure_year and Contract are the top 4 most important variables to predict churn. The PhoneService, Gender, SeniorCitizen, Dependents, Partner, MultipleLines, PaperlessBilling, StreamingTV, Movies, DeviceProtection and OnlineBackup have comparatively less effect on Churn.

```
In [47]:  
1 options(repr.plot.width = 15, repr.plot.height = 8)  
2 plot(rf_df, main = "Error by number of trees")
```



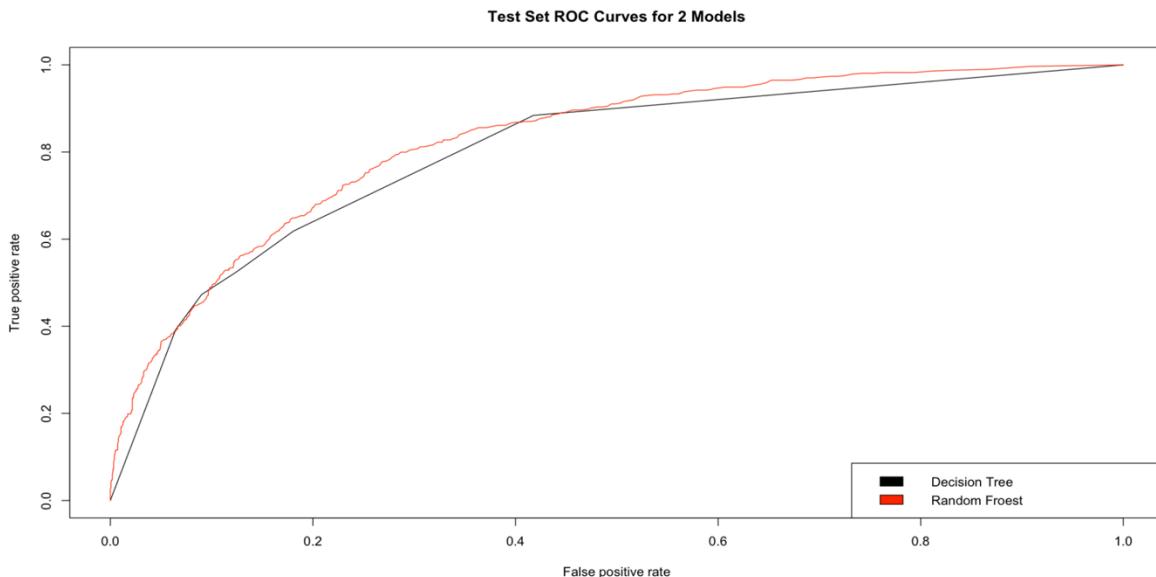
Variable Importance

```
In [48]: 1 options(repr.plot.width = 15, repr.plot.height = 8)
2 varImpPlot(rf_df,type=2)
```



Comparison of ROC and AUC for Decision Tree and Random Forest models

```
In [50]: 1 preds_list <- list(testtree_prob[,2], testrf_prob[,2])
2 m <- length(preds_list)
3 actuals_list <- rep(list(testdf$churn_number), m)
4
5 pred <- prediction(preds_list, actuals_list)
6 rocs <- performance(pred, "tpr", "fpr")
7 options(repr.plot.width = 15, repr.plot.height = 8)
8 plot(rocs, col = as.list(1:m), main = "Test Set ROC Curves for 2 Models")
9 legend(x = "bottomright",
10         legend = c("Decision Tree", "Random Froest"),
11         fill = 1:m)
```



The random forest model works better than the decision tree model. The accuracy is 0.7908 for Decision Tree and 0.7927 for Random Forest. The AUC score of decision tree is 0.79 and 0.828 for random forest.

Part5. Conclusion

- What makes the problem interesting from the viewpoint of analytics?

For me, I think the most interesting part is the difference between the decision tree and random forest. While random forest is a collection of decision trees, there are some differences. However, both of them can be used in different scenario and return the result we want!

If you input a training dataset with features and labels into a decision tree, it will formulate some set of rules, which will be used to make the predictions. If you put the features and labels into a decision tree, it will generate some rules that help predict customer churn. In comparison, the random forest algorithm randomly selects observations and features to build several decision trees and then averages the results. Another difference is "deep" decision trees might suffer from overfitting. Most of the time, random forest prevents this by creating random subsets of the features and building smaller trees using those subsets. Afterwards, it combines the subtrees.

- How did the chosen technique help to illuminate the, or solve the problem?

Tree-based models are very popular in machine learning. The decision tree, the foundation of tree-based models, is quite straightforward to interpret, but generally a weak predictor. Ensemble models can be used to generate stronger predictions from many trees, with random forest or gradient boosting as two of the most popular. All tree-based models can be used for regression or classification and can handle non-linear relationships quite well. It's awesome that we can use these models to predict the outcome and make the decision based on the result.

- What analysis do you think should be conducted next?

For customer churn prediction, the next analysis should be customer analysis. The company should focus on the customers who show non churn and find out the reason. Therefore, company can find out a better approach to retain the customers. The company should also focus on customers who show satisfaction. Company can propose some rewards for them to make sure that customers keep high satisfaction.

- Overall Summary

From the above analysis, we can get the result below.

Decision Tree

Confusion Matrix and Statistics

Reference

Prediction no yes no 1383 300 yes 137 269

Accuracy : 0.7908
95% CI : (0.7727, 0.8081)

No Information Rate : 0.7276
P-Value [Acc > NIR] : 1.597e-11

Kappa : 0.4203

McNemar's Test P-Value : 9.225e-15

Sensitivity : 0.9099
Specificity : 0.4728
Pos Pred Value : 0.8217
Neg Pred Value : 0.6626
Prevalence : 0.7276
Detection Rate : 0.6620

Detection Prevalence : 0.8056

Balanced Accuracy : 0.6913

Random Forest

Confusion Matrix and Statistics

Reference

Prediction no yes no 1375 288 yes 145 281

Accuracy : 0.7927
95% CI : (0.7747, 0.8099)
No Information Rate : 0.7276
P-Value [Acc > NIR] : 3.803e-12

Kappa : 0.4325

Mcnemar's Test P-Value : 8.849e-12

Sensitivity : 0.9046
Specificity : 0.4938
Pos Pred Value : 0.8268
Neg Pred Value : 0.6596
Prevalence : 0.7276
Detection Rate : 0.6582

Detection Prevalence : 0.7961

Balanced Accuracy : 0.6992

The AUC score of Decision Tree is 0.8, and the AUC score of Random Forest is 0.828. We can conclude that the performance of random forest is better than decision tree.