

Customer Churn Analytics

Data Science Academy

Vinicius Biazon

Customer churn (cancellation rates) occurs when customers or subscribers stop doing business with a company or service.

One sector in which knowing and predicting cancellation rates is particularly useful is the telecommunications, because most customers have several options to choose from within a geographic location.

In this project, I predicted the customer churn using a telecommunications data set offered on IBM Sample Data Sets portal. I used logistic regression, decision tree and random forest as Machine Learning models.

<https://community.ibm.com/community/user/businessanalytics/blogs/steven-macko/2019/07/11/telco-customer-churn-1113>

Step 1 - Defining the working directory and Loading the packages

```
# Defining the working directory
setwd("D:/Documentos/FCD/BigDataRAzure/Cap06")
getwd()
```

```
## [1] "D:/Documentos/FCD/BigDataRAzure/Cap06"
```

```
# Loading the packages
library(plyr)
library(corrplot)
library(ggplot2)
library(gridExtra)
library(ggthemes)
library(caret)
library(MASS)
library(randomForest)
library(party)
```

Step 2 - Loading and Cleaning Data

```
churn <- read.csv('Telco-Customer-Churn.csv')

# Removing all lines with missing values.
sapply(churn, function(x) sum(is.na(x)))
```

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

```
churn <- churn[complete.cases(churn), ]
```

```
# Changing "No internet service" to "No" in six columns:
# "OnlineSecurity", "OnlineBackup", "DeviceProtection", "TechSupport", "streamingTV", "streamingMovies"
cols_recode1 <- c(10:15)
for(i in 1:ncol(churn[,cols_recode1])) {
  churn[,cols_recode1][,i] <- as.factor(mapvalues
                                         (churn[,cols_recode1][,i], from=c("No internet service"),to=c(

```

```
# Changing "No phone service" to "No" on "MultipleLines" column
churn$MultipleLines <- as.factor(mapvalues(churn$MultipleLines,
                                           from=c("No phone service"),
                                           to=c("No")))

```

```
# Grouping "tenure" into five categories:
# "0-12 Month", "12-24 Month", "24-48 Month", "48-60 Month", "> 60 Month"
min(churn$tenure); max(churn$tenure)
```

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

```
## [1] 72
```

```
group_tenure <- function(tenure){
  if (tenure >= 0 & tenure <= 12){
    return('0-12 Month')
  }else if(tenure > 12 & tenure <= 24){
    return('12-24 Month')
  }else if (tenure > 24 & tenure <= 48){
    return('24-48 Month')
  }else if (tenure > 48 & tenure <=60){
    return('48-60 Month')
  }else if (tenure > 60){
    return('> 60 Month')
  }
}
```

```

churn$tenure_group <- sapply(churn$tenure,group_tenure)
churn$tenure_group <- as.factor(churn$tenure_group)

# Changing "SeniorCitizen" column values from 0 or 1 to "No" or "Yes".
churn$SeniorCitizen <- as.factor(mapvalues(churn$SeniorCitizen,
                                           from=c("0","1"),
                                           to=c("No", "Yes")))

# Removing unnecessary columns for analysis.
churn$customerID <- NULL
churn$tenure <- NULL

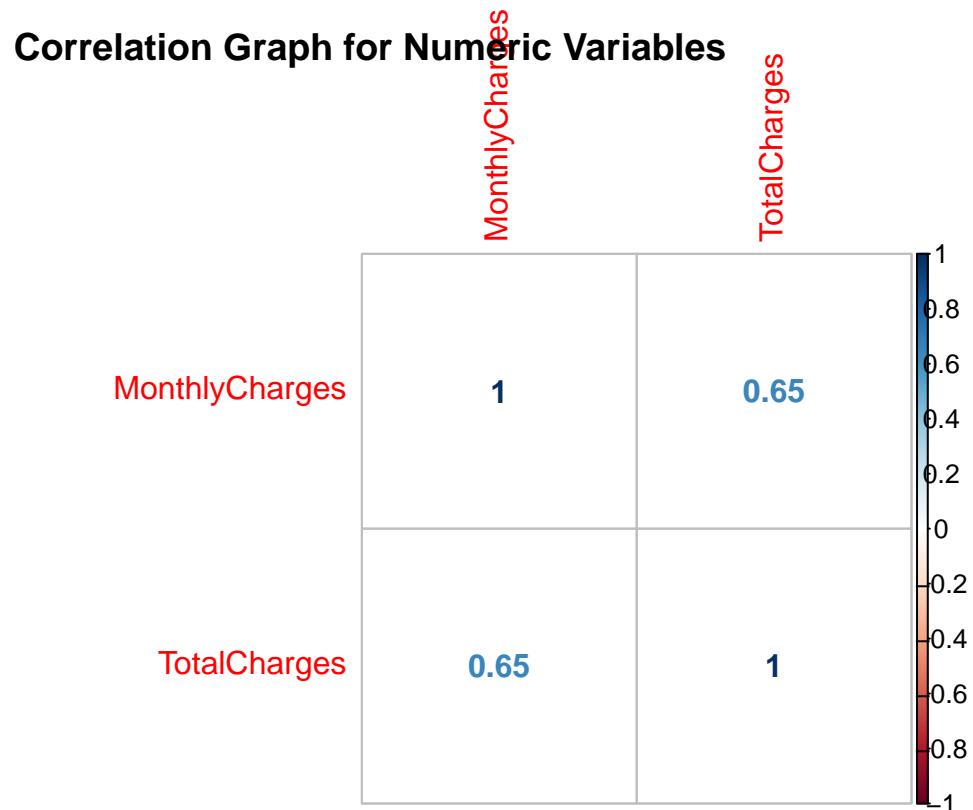
```

Step 3 - Exploratory data analysis

```

# Numeric variables correlation
numeric.var <- sapply(churn, is.numeric)
corr.matrix <- cor(churn[,numeric.var])
corrplot(corr.matrix, main="\n\nCorrelation Graph for Numeric Variables", method="number")

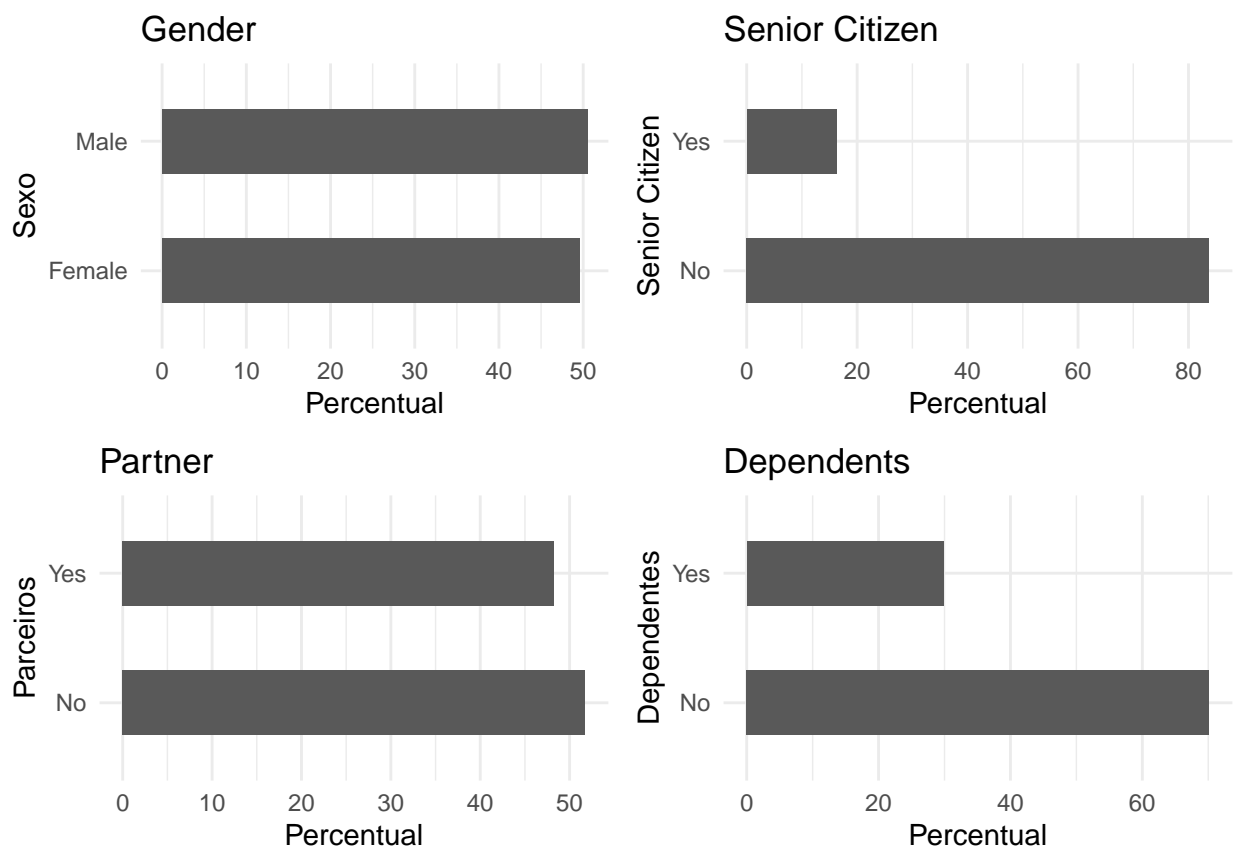
```



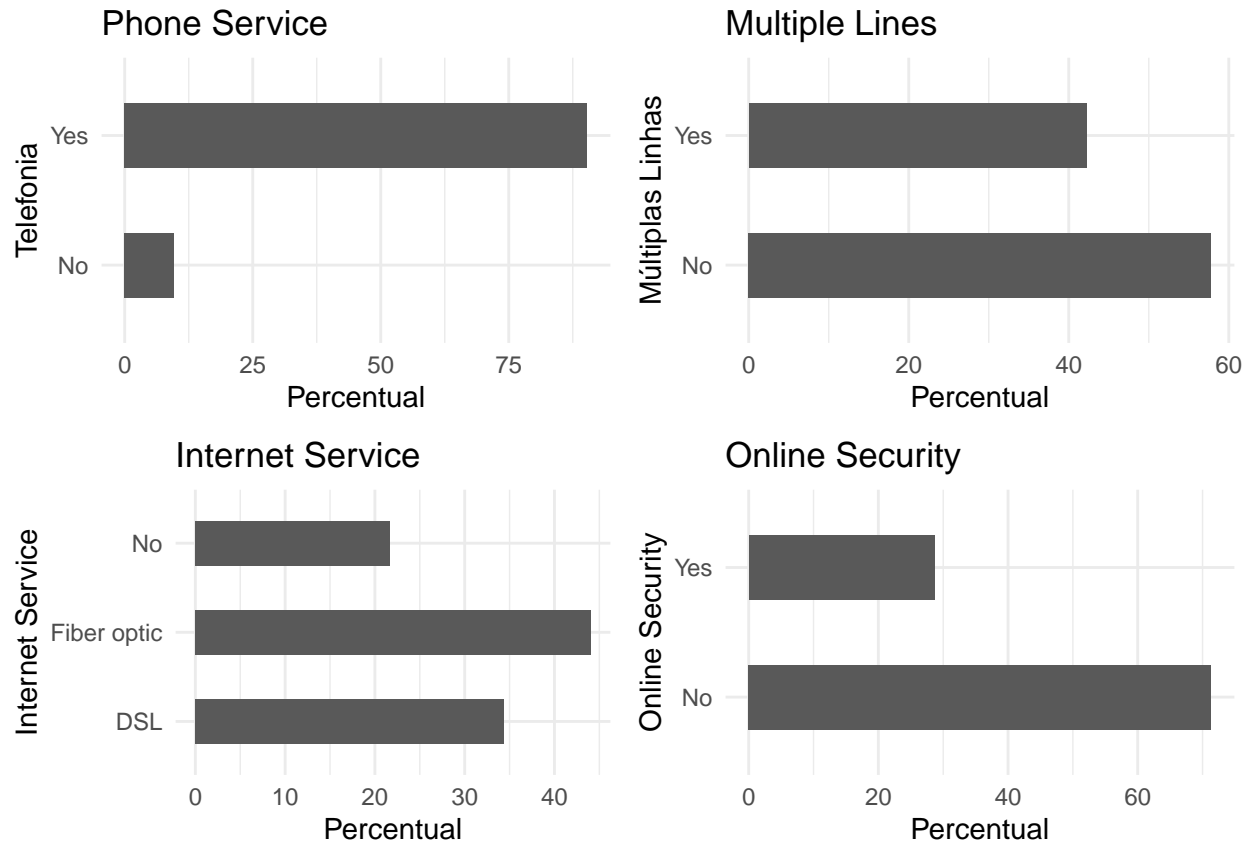
```
# Removing Total Charges to avoid overfitting
churn$TotalCharges <- NULL
```

```
# Categorical variable bar graphs
```

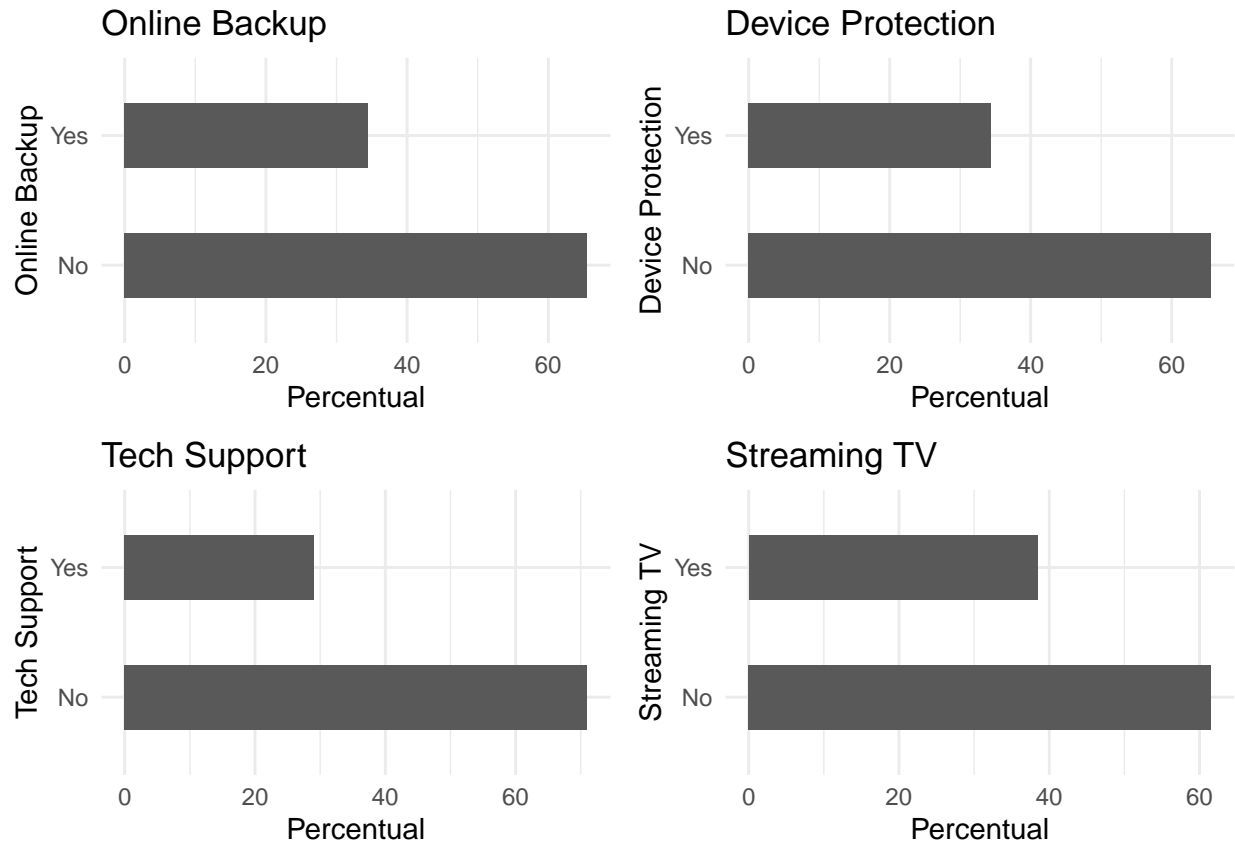
```
p1 <- ggplot(churn, aes(x=gender)) + ggtitle("Gender") + xlab("Sexo") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p2 <- ggplot(churn, aes(x=SeniorCitizen)) + ggtitle("Senior Citizen") + xlab("Senior Citizen") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p3 <- ggplot(churn, aes(x=Partner)) + ggtitle("Partner") + xlab("Parceiros") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p4 <- ggplot(churn, aes(x=Dependents)) + ggtitle("Dependents") + xlab("Dependentes") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
grid.arrange(p1, p2, p3, p4, ncol=2)
```



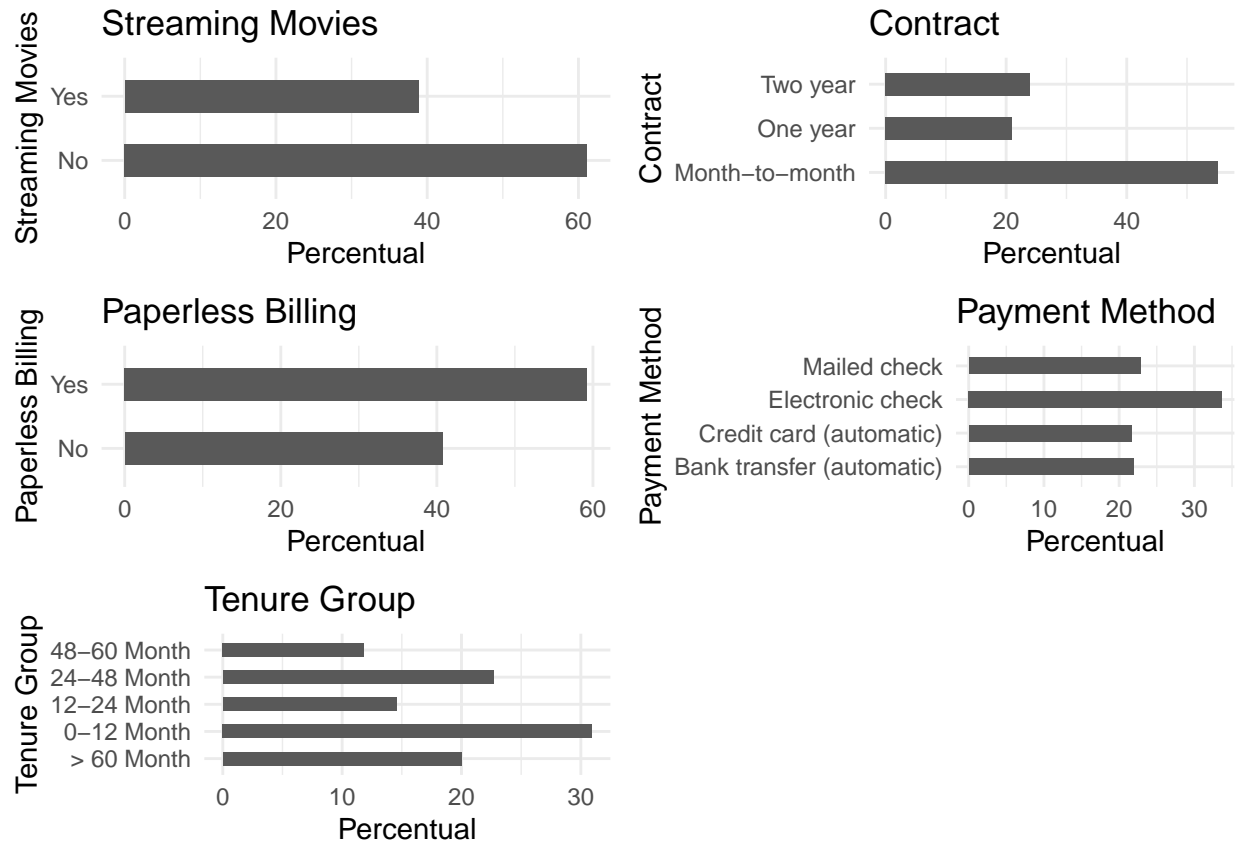
```
p5 <- ggplot(churn, aes(x=PhoneService)) + ggtitle("Phone Service") + xlab("Telefonia") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p6 <- ggplot(churn, aes(x=MultipleLines)) + ggtitle("Multiple Lines") + xlab("Múltiplas Linhas") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p7 <- ggplot(churn, aes(x=InternetService)) + ggtitle("Internet Service") + xlab("Internet Service") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p8 <- ggplot(churn, aes(x=OnlineSecurity)) + ggtitle("Online Security") + xlab("Online Security") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
grid.arrange(p5, p6, p7, p8, ncol=2)
```



```
p9 <- ggplot(churn, aes(x=OnlineBackup)) + ggtitle("Online Backup") + xlab("Online Backup") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p10 <- ggplot(churn, aes(x=DeviceProtection)) + ggtitle("Device Protection") + xlab("Device Protection") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p11 <- ggplot(churn, aes(x=TechSupport)) + ggtitle("Tech Support") + xlab("Tech Support") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p12 <- ggplot(churn, aes(x=StreamingTV)) + ggtitle("Streaming TV") + xlab("Streaming TV") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
grid.arrange(p9, p10, p11, p12, ncol=2)
```



```
p13 <- ggplot(churn, aes(x=StreamingMovies)) + ggtitle("Streaming Movies") + xlab("Streaming Movies") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p14 <- ggplot(churn, aes(x=Contract)) + ggtitle("Contract") + xlab("Contract") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p15 <- ggplot(churn, aes(x=PaperlessBilling)) + ggtitle("Paperless Billing") + xlab("Paperless Billing") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p16 <- ggplot(churn, aes(x=PaymentMethod)) + ggtitle("Payment Method") + xlab("Payment Method") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
p17 <- ggplot(churn, aes(x=tenure_group)) + ggtitle("Tenure Group") + xlab("Tenure Group") +
  geom_bar(aes(y = 100*(..count..)/sum(..count..)), width = 0.5) + ylab("Percentual") + coord_flip() +
grid.arrange(p13, p14, p15, p16, p17, ncol=2)
```



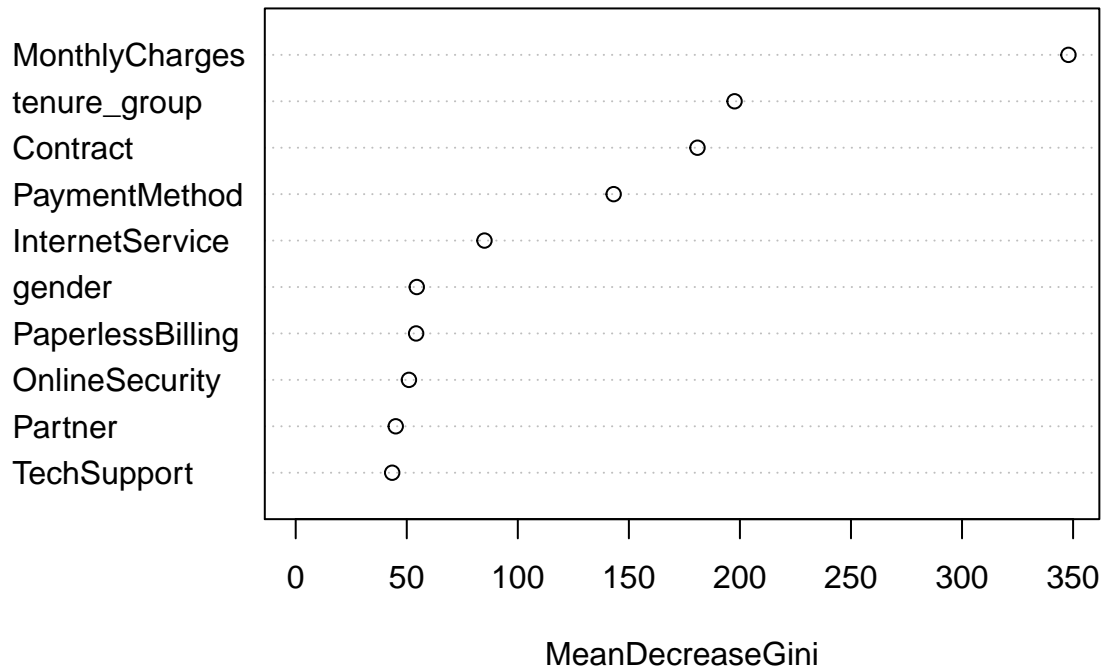
All categorical variables will be maintained because appear to have a reasonably wide distribution.

Step 4 - Feature Selection with Random Forest

```
# Dividing data into training and test - 70:30 ratio
intrain <- createDataPartition(churn$Churn,p=0.7,list=FALSE)
training <- churn[intrain,]
testing <- churn[-intrain,]

# Feature Selection
rfModel <- randomForest(Churn ~., data = training)
pred_rf <- predict(rfModel, testing)
varImpPlot(rfModel, sort=T, n.var = 10, main = 'Top 10 Feature Importance')
```

Top 10 Feature Importance



Step 5 - Predictive Modeling: Logistic Regression

```
# Logistic regression model fitting
LogModel <- glm(Churn ~ MonthlyCharges+tenure_group+Contract+PaymentMethod+InternetService, family=binomial)
print(summary(LogModel))
```

```
##
## Call:
## glm(formula = Churn ~ MonthlyCharges + tenure_group + Contract +
##      PaymentMethod + InternetService, family = binomial(link = "logit"),
##      data = training)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.6966  -0.6677  -0.3301   0.7607   3.0480
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -2.241722   0.300292  -7.465 8.32e-14 ***
## MonthlyCharges    0.003641   0.003603   1.010 0.312263
## tenure_group0-12 Month  1.551455   0.185930   8.344 < 2e-16 ***
## tenure_group12-24 Month  0.661943   0.188238   3.517 0.000437 ***
## tenure_group24-48 Month  0.347888   0.173568   2.004 0.045035 *
## tenure_group48-60 Month  0.228620   0.191299   1.195 0.232050
```



```
## ContractOne year          -0.914781    0.123360   -7.416 1.21e-13 ***
## ContractTwo year         -1.899855    0.206253   -9.211 < 2e-16 ***
## PaymentMethodCredit card (automatic) -0.003688    0.132202   -0.028 0.977745
## PaymentMethodElectronic check      0.557206    0.110710    5.033 4.83e-07 ***
## PaymentMethodMailed check    0.011122    0.135003    0.082 0.934341
## InternetServiceFiber optic    0.895961    0.154298    5.807 6.37e-09 ***
## InternetServiceNo          -0.794355    0.181777   -4.370 1.24e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 5702.8  on 4923  degrees of freedom
## Residual deviance: 4208.6  on 4911  degrees of freedom
## AIC: 4234.6
##
## Number of Fisher Scoring iterations: 6
```

```
# Accuracy
testing$Churn <- as.character(testing$Churn)
testing$Churn[testing$Churn=="No"] <- "0"
testing$Churn[testing$Churn=="Yes"] <- "1"
fitted.results <- predict(LogModel,newdata=testing,type='response')
fitted.results <- ifelse(fitted.results > 0.5,1,0)
misClasificError <- mean(fitted.results != testing$Churn)
print(paste('Logistic Regression Accuracy',1-misClasificError))
```

```
## [1] "Logistic Regression Accuracy 0.790796963946869"
```

```
# Logistic Regression Confusion Matrix
print("Logistic Regression - Confusion Matrix"); table(testing$Churn, fitted.results > 0.5)
```

```
## [1] "Logistic Regression - Confusion Matrix"
```

```
##
## FALSE TRUE
## 0 1393 155
## 1 286 274
```

```
# Odds Ratio
exp(cbind(OR=coef(LogModel), confint(LogModel)))
```

```
## Waiting for profiling to be done...
```

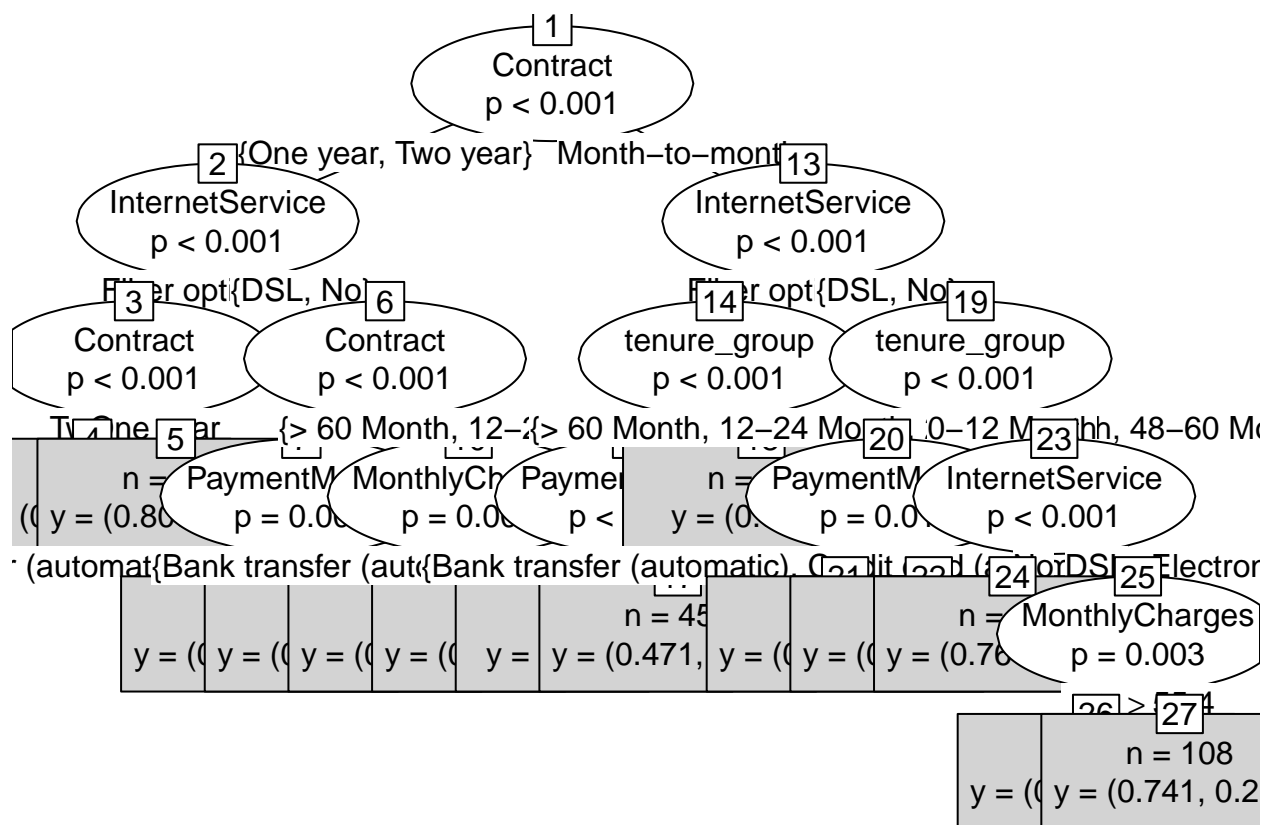
```
##
## OR      2.5 %    97.5 %
## (Intercept) 0.1062754 0.05872199 0.1906195
## MonthlyCharges 1.0036476 0.99659493 1.0107751
## tenure_group0-12 Month 4.7183316 3.28895774 6.8207422
## tenure_group12-24 Month 1.9385556 1.34386242 2.8123553
## tenure_group24-48 Month 1.4160735 1.01059388 1.9969128
## tenure_group48-60 Month 1.2568642 0.86382548 1.8300625
```

```
## ContractOne year          0.4006045 0.31363087 0.5088190
## ContractTwo year         0.1495904 0.09858153 0.2216115
## PaymentMethodCredit card (automatic) 0.9963189 0.76875267 1.2911200
## PaymentMethodElectronic check 1.7457876 1.40677379 2.1716125
## PaymentMethodMailed check 1.0111842 0.77636886 1.3182203
## InternetServiceFiber optic 2.4496880 1.81253220 3.3192432
## InternetServiceNo        0.4518728 0.31575969 0.6441814
```

For each unit increase in the Monthly Charge, there is a 2.5% reduction in the probability of the customer canceling the subscription.

Step 6 - Predictive Modeling: Decision Tree

```
tree <- ctree(Churn ~ MonthlyCharges+tenure_group+Contract+PaymentMethod+InternetService, training)
plot(tree, type='simple')
```



The Contract is the most important variable for predicting customer churn. If a customer is on a monthly contract, and in the 0 to 12 month tenure group, and using PaperlessBilling, he is more likely to cancel the subscription.

```
# Decision Tree Confusion Matrix
pred_tree <- predict(tree, testing)
print("Decision Tree Confusion Matrix"); table(Predicted = pred_tree, Actual = testing$Churn)
```

```
## [1] "Decision Tree Confusion Matrix"
```

```
##           Actual
## Predicted    0    1
##      No  1344  262
##      Yes   204  298
```

```
# Accuracy
p1 <- predict(tree, training)
tab1 <- table(Predicted = p1, Actual = training$Churn)
tab2 <- table(Predicted = pred_tree, Actual = testing$Churn)
print(paste('Decision Tree Accuracy', sum(diag(tab2))/sum(tab2)))
```

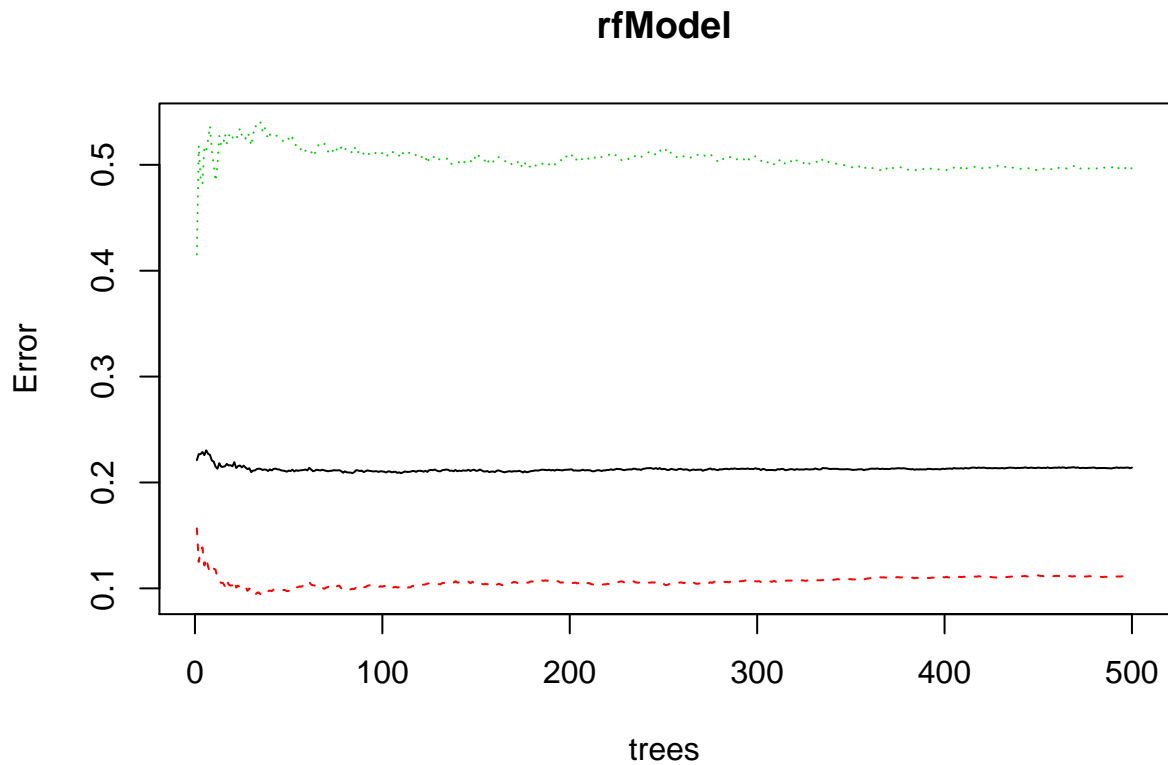
```
## [1] "Decision Tree Accuracy 0.778937381404175"
```

Step 7 - Predictive Modeling: Random Forest

```
rfModel <- randomForest(Churn ~ MonthlyCharges+tenure_group+Contract+PaymentMethod+InternetService, data = testing, ntree = 500)
print(rfModel)
```

```
##
## Call:
## randomForest(formula = Churn ~ MonthlyCharges + tenure_group + Contract + PaymentMethod + InternetService, data = testing, ntree = 500)
##           Type of random forest: classification
##           Number of trees: 500
## No. of variables tried at each split: 2
##
## OOB estimate of error rate: 21.41%
## Confusion matrix:
##      No Yes class.error
## No  3212 403  0.1114799
## Yes   651 658  0.4973262
```

```
plot(rfModel)
```



The prediction is very good when predicting “No”, but the error rate is much higher when predicting “yes”.

```
# Predicting values with test data
pred_rf <- predict(rfModel, testing)

# Confusion Matrix
print("Random Forest - Confusion Matrix"); table(testing$Churn, pred_rf)
```

```
## [1] "Random Forest - Confusion Matrix"
```

```
##      pred_rf
##      No  Yes
## 0 1364 184
## 1  253 307
```