

Telco Churn

In this Project we shall predict customer churn for a Telecommunications company “Telco”. All credit for the dataset goes to Kaggle User “BlastChar” and can be accessed here: <https://www.kaggle.com/blatchar/telco-customer-churn>

Let's start off by loading our dataset and the required libraries.

```
#Let's import the dataset
```

```
library(readr)
WA_F <- read_csv("Telco Customer Churn/Telco.csv")
```

```
## Parsed with column specification:
## cols(
##   .default = col_character(),
##   SeniorCitizen = col_double(),
##   tenure = col_double(),
##   MonthlyCharges = col_double(),
##   TotalCharges = col_double()
## )
```

```
## See spec(...) for full column specifications.
```

```
mydata <-WA_F
```

```
#Let's load the Libraries
```

```
library(data.table)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:data.table':
##
##   between, first, last

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

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

```
library(ggplot2)
library(corrplot)
```

```
## corrplot 0.84 loaded
```

```
library(cowplot)
```

```
##
## *****

## Note: As of version 1.0.0, cowplot does not change the
## default ggplot2 theme anymore. To recover the previous
```

```
## behavior, execute:
## theme_set(theme_cowplot())

## *****

library(caret)

## Loading required package: lattice

library(class)
library(gridExtra)

##
## Attaching package: 'gridExtra'

## The following object is masked from 'package:dplyr':
##
## combine

library(tm)

## Loading required package: NLP

##
## Attaching package: 'NLP'

## The following object is masked from 'package:ggplot2':
##
## annotate

library(e1071)
```

Sections:

1. *Introduction and Dataset exploration*
2. *Variable Analysis and Data Exploration / Cleanup*
3. *Model Building*
4. *Prediction and Accuracy Evaluation*
5. *Conclusions and Future Work*

SECTION 1: INTRODUCTION

We will create a Logistic Regression model to predict customer churn.

We will be analysing variable by variable and modifying the data in a way that can be easily interpreted by a Logistic Regression model. Variables with text data will be split into “N” columns where “N” is the number of levels of the variable. These columns will be populated with 1s and 0s depending on the level of the original Variable column. Variables with numeric data can be fit into the regression model as they are.

All variables will then be standardized. This will give the variables zero-mean and unit-variance. This is done to mitigate the issue created by different ranges shown by different variables.

Summing this up, we will:

- Load our data
- Take a look at each variable to:
 - Check for NAs (and replace if necessary)
 - Manipulate Data to make it interpretable by a regression model
- Select which features we will include into our regression model
- Standardize all our data

- Split our data into Train and Test sets
- Fit our Logistic Regression Model
- Use the Train set to predict our Test Set
- Evaluate the accuracy of our model

Dataset Exploration Let us now take a look at our dataset.

```
#Let's take a look at our dataset
```

```
head(mydata)
```

```
## # A tibble: 6 x 21
##   customerID gender SeniorCitizen Partner Dependents tenure PhoneService
##   <chr>         <chr>         <dbl> <chr>    <chr>         <dbl> <chr>
## 1 7590-VHVEG Female             0 Yes     No             1 No
## 2 5575-GNVDE Male             0 No      No            34 Yes
## 3 3668-QPYBK Male             0 No      No             2 Yes
## 4 7795-CFOWC Male             0 No      No            45 No
## 5 9237-HQITU Female          0 No      No             2 Yes
## 6 9305-CDSKC Female          0 No      No             8 Yes
## # ... with 14 more variables: MultipleLines <chr>, InternetService <chr>,
## #   OnlineSecurity <chr>, OnlineBackup <chr>, DeviceProtection <chr>,
## #   TechSupport <chr>, StreamingTV <chr>, StreamingMovies <chr>,
## #   Contract <chr>, PaperlessBilling <chr>, PaymentMethod <chr>,
## #   MonthlyCharges <dbl>, TotalCharges <dbl>, Churn <chr>
```

Let's check for NAs. As per the below output, we have 11 NAs.

```
#Check for NAs
```

```
sum(is.na(mydata))
```

```
## [1] 11
```

```
#We have 11 Na's
```

Let us take a look at our **output Variable**; which is the column “Churn”. Let's check for NAs, and plot. There are much less people who churned than who didn't. As per the below output, we have no NAs.

```
#Let's look at our output
```

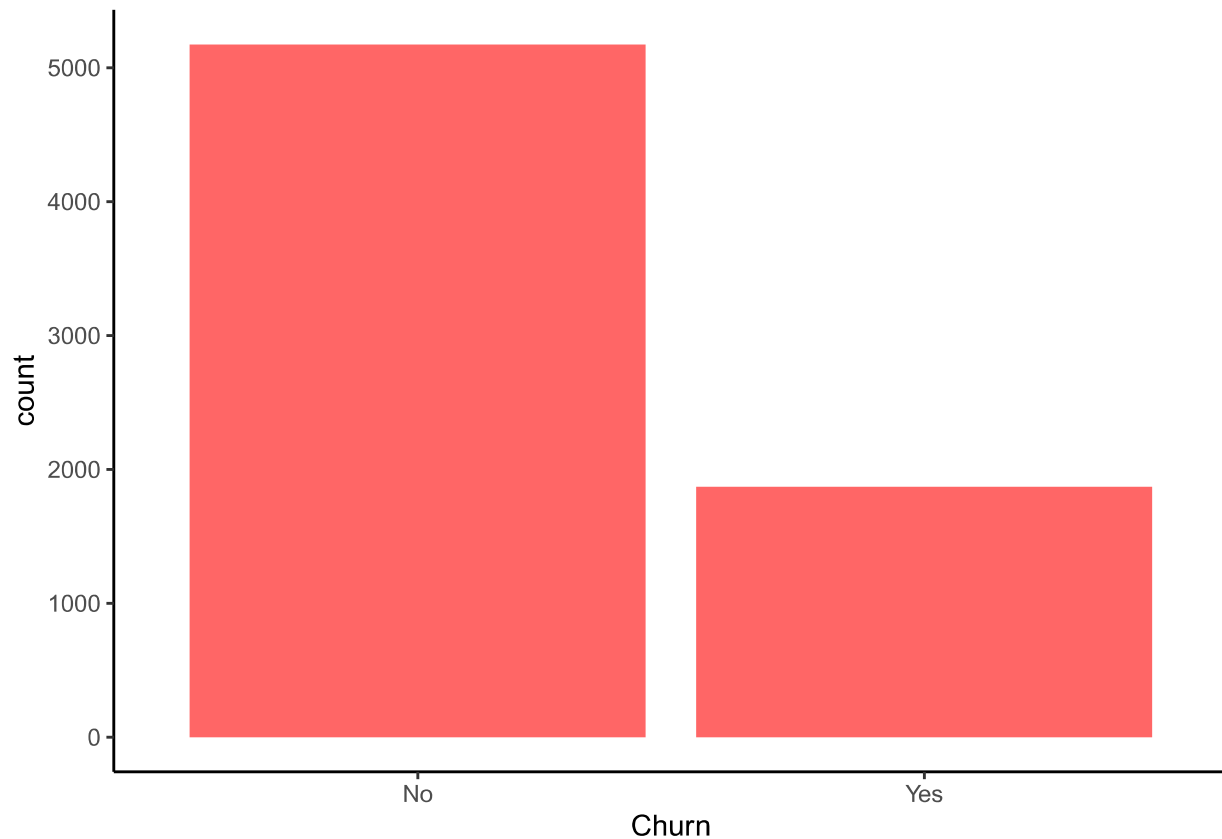
```
base::table(mydata$Churn)
```

```
##
##   No   Yes
## 5174 1869
```

```
sum(is.na(mydata$Churn))
```

```
## [1] 0
```

```
ggplot(mydata, aes(mydata$Churn)) + geom_bar(fill = "#FF6666") + theme_classic() + xlab("Churn")
```



Since we know that we will be using Logistic Regression - we need to switch our output Variable to 0s and 1s. The value of 1 represents customer churn.

I added a new column “ChurnLog” which contains binary data for Churn.

```
mydata$ChurnLog = 0
mydata$ChurnLog[mydata$Churn == "Yes"] = 1
mydata$ChurnLog[mydata$Churn == "No"] = 0

head(mydata)
```

```
## # A tibble: 6 x 22
##   customerID gender SeniorCitizen Partner Dependents tenure PhoneService
##   <chr>      <chr>      <dbl> <chr>   <chr>      <dbl> <chr>
## 1 7590-VHVEG Female          0 Yes    No          1 No
## 2 5575-GNVDE Male           0 No     No         34 Yes
## 3 3668-QPYBK Male           0 No     No          2 Yes
## 4 7795-CFOCW Male           0 No     No         45 No
## 5 9237-HQITU Female          0 No     No          2 Yes
## 6 9305-CDSKC Female          0 No     No          8 Yes
## # ... with 15 more variables: MultipleLines <chr>, InternetService <chr>,
## #   OnlineSecurity <chr>, OnlineBackup <chr>, DeviceProtection <chr>,
## #   TechSupport <chr>, StreamingTV <chr>, StreamingMovies <chr>,
## #   Contract <chr>, PaperlessBilling <chr>, PaymentMethod <chr>,
## #   MonthlyCharges <dbl>, TotalCharges <dbl>, Churn <chr>, ChurnLog <dbl>
```

Let's take a look at our column names. Creating a vector with the names of columns is useful if you want to call upon it to check the exact names of the columns.

```
ColVec <- colnames(mydata)
ColVec
```

```
## [1] "customerID"      "gender"           "SeniorCitizen"
## [4] "Partner"         "Dependents"       "tenure"
## [7] "PhoneService"    "MultipleLines"    "InternetService"
## [10] "OnlineSecurity"  "OnlineBackup"     "DeviceProtection"
## [13] "TechSupport"     "StreamingTV"      "StreamingMovies"
## [16] "Contract"        "PaperlessBilling" "PaymentMethod"
## [19] "MonthlyCharges"  "TotalCharges"     "Churn"
## [22] "ChurnLog"
```

SECTION 2: VARIABLE ANALYSIS AND DATA EXPLORATION / CLEANUP

->>> Variable Analysis

///GENDER///

First we check our levels. We have two, “Male” and “Female”. Then we check for NAs. We have no NAs.

```
base::table(mydata$gender)
```

```
##
## Female    Male
##   3488    3555
```

```
sum(is.na(mydata$gender))
```

```
## [1] 0
```

I created a dataframe using Dummy Variables. This dataframe contains two columns with binary values representing the Gender variable from mydata. This will be used to develop the logistic regression model.

```
Genderframe <- as.data.frame(mydata$gender)
dmy <- dummyVars(" ~ .", data = Genderframe)
Genderframe2 <- data.frame(predict(dmy, newdata = Genderframe))
head(Genderframe2)
```

```
##   X.mydata.gender.Female X.mydata.gender.Male
## 1                      1                      0
## 2                      0                      1
## 3                      0                      1
## 4                      0                      1
## 5                      1                      0
## 6                      1                      0
```

->>> Variable Analysis

///PARTNER///

First we check our levels. We have two, “Yes” and “No”. Then we check for NAs. I added a new column “PartnerLog” which contains binary data for Partner. As per the output - we have no NAs.

```
mydata$PartnerLog = 0
sum(is.na(mydata$Partner))
```

```
## [1] 0
```

```
mydata$PartnerLog[mydata$Partner == "Yes"] = 1
mydata$PartnerLog[mydata$Partner == "No"] = 0
head(mydata)
```

```
## # A tibble: 6 x 23
##   customerID gender SeniorCitizen Partner Dependents tenure PhoneService
##   <chr>      <chr>      <dbl> <chr>    <chr>      <dbl> <chr>
## 1 7590-VHVEG Female          0 Yes     No          1 No
## 2 5575-GNVDE Male           0 No      No         34 Yes
## 3 3668-QPYBK Male           0 No      No          2 Yes
## 4 7795-CFQCW Male           0 No      No         45 No
## 5 9237-HQITU Female         0 No      No          2 Yes
## 6 9305-CDSKC Female         0 No      No          8 Yes
## # ... with 16 more variables: MultipleLines <chr>, InternetService <chr>,
## #   OnlineSecurity <chr>, OnlineBackup <chr>, DeviceProtection <chr>,
## #   TechSupport <chr>, StreamingTV <chr>, StreamingMovies <chr>,
## #   Contract <chr>, PaperlessBilling <chr>, PaymentMethod <chr>,
## #   MonthlyCharges <dbl>, TotalCharges <dbl>, Churn <chr>, ChurnLog <dbl>,
## #   PartnerLog <dbl>
```

->>> Variable Analysis

```
///SENIOR CITIZEN///
```

Let's take a look at Senior Citizen. It is already in 1s and 0s, so we need not modify it. Also we have no NAs.

```
base::table(mydata$SeniorCitizen)
```

```
##
##    0    1
## 5901 1142
```

```
sum(is.na(mydata$SeniorCitizen))
```

```
## [1] 0
```

Variable Analysis

```
///DEPENDENTS///
```

Now let us take a look at dependents. As per the code output below, we have no NAs. We will create another column with Binary Data, as we did for Gender.

```
base::table(mydata$Dependents)
```

```
##
##   No  Yes
## 4933 2110
```

```
sum(is.na(mydata$Dependents))
```

```
## [1] 0
```

```
mydata$DependentsLog = 0
mydata$DependentsLog[mydata$Dependents == "Yes"] = 1
mydata$DependentsLog[mydata$Dependents == "No"] = 0
tail(mydata)
```

```
## # A tibble: 6 x 24
##   customerID gender SeniorCitizen Partner Dependents tenure PhoneService
##   <chr>      <chr>      <dbl> <chr>    <chr>      <dbl> <chr>
## 1 2569-WGERO Female          0 No      No          72 Yes
## 2 6840-RESVB Male           0 Yes     Yes          24 Yes
## 3 2234-XADUH Female         0 Yes     Yes          72 Yes
## 4 4801-JZAZL Female         0 Yes     Yes          11 No
```

```
## 5 8361-LTMKD Male          1 Yes      No          4 Yes
## 6 3186-AJIEK Male          0 No       No          66 Yes
## # ... with 17 more variables: MultipleLines <chr>, InternetService <chr>,
## #   OnlineSecurity <chr>, OnlineBackup <chr>, DeviceProtection <chr>,
## #   TechSupport <chr>, StreamingTV <chr>, StreamingMovies <chr>,
## #   Contract <chr>, PaperlessBilling <chr>, PaymentMethod <chr>,
## #   MonthlyCharges <dbl>, TotalCharges <dbl>, Churn <chr>, ChurnLog <dbl>,
## #   PartnerLog <dbl>, DependentsLog <dbl>
```

Variable Analysis

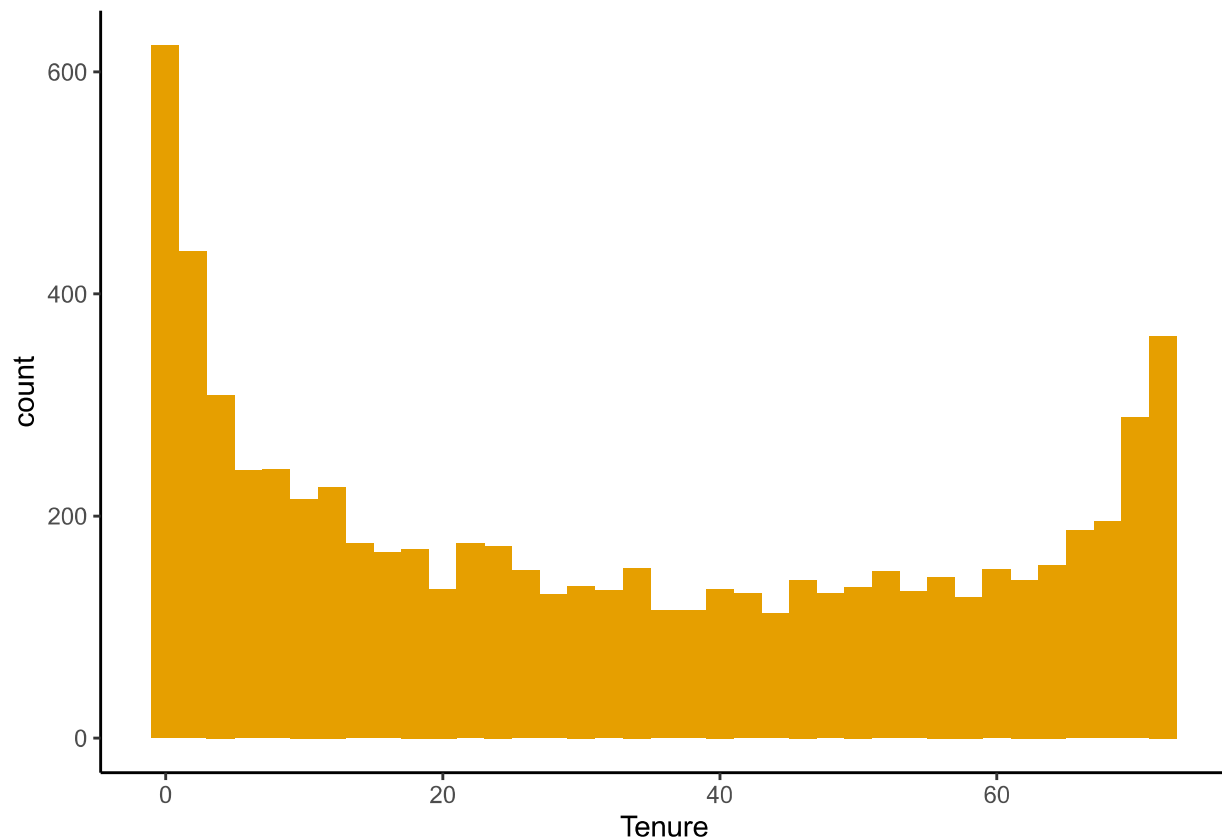
```
///TENURE///
```

Let's take a look at tenure. As per the code output, we have no NAs. The distribution is not Normal. When looking at the violin plot it is evident that there is a higher probability that people with lower tenure will Churn.

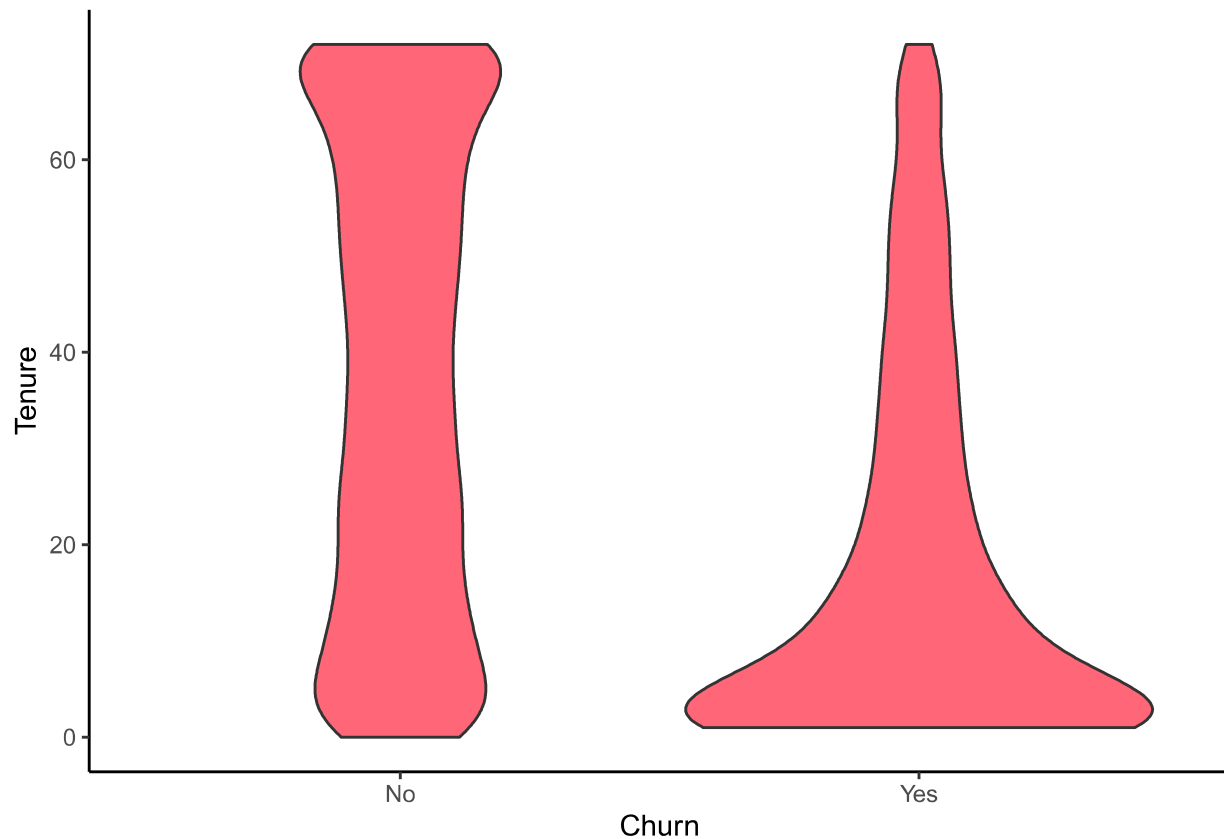
```
sum(is.na(mydata$tenure))
```

```
## [1] 0
```

```
ggplot(mydata, aes(mydata$tenure)) + geom_histogram(binwidth = 2, fill = "#E69F00") + theme_classic() + xlab("Tenure")
```



```
ggplot(mydata, aes(mydata$Churn, mydata$tenure)) + geom_violin(fill = "#FF6677") + xlab("Churn") + ylab("Tenure")
```



Now Lets create a dataframe to split the tenure data into:

1. Values with zero value
2. Values between Between 1 and 20
3. Values between Between 21 and 40
4. Values between Between 41 and 60
5. Values over 61

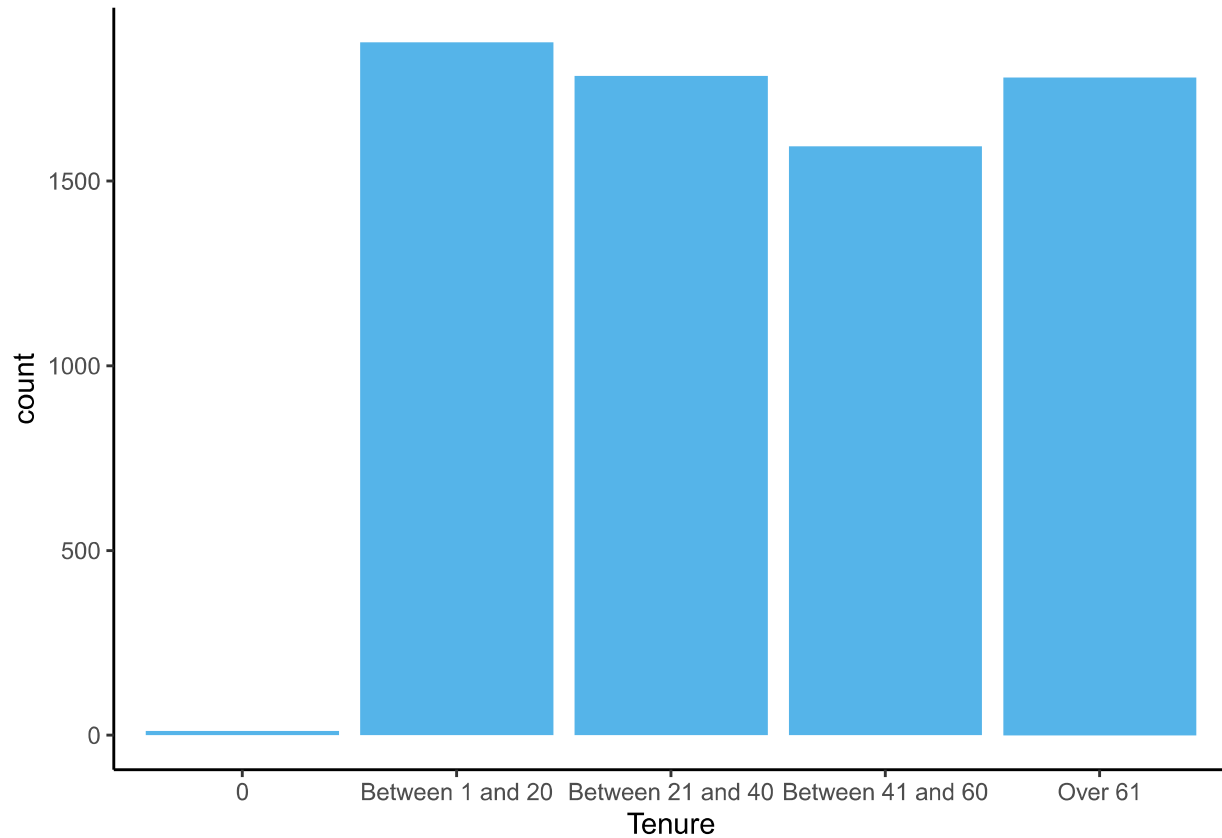
```
TenureFrame <- as.data.frame(mydata$tenure)
x = 1
while (x<=nrow(TenureFrame)){
  if ((TenureFrame$mydata$tenure`[x] > 0) & (TenureFrame$mydata$tenure`[x] <= 20))
    {TenureFrame$mydata$tenure`[x] = "Between 1 and 20"}
  else if ((TenureFrame$mydata$tenure`[x] > 20) & (TenureFrame$mydata$tenure`[x] <= 40))
    {TenureFrame$mydata$tenure`[x] = "Between 21 and 40"}
  else if ((TenureFrame$mydata$tenure`[x] > 40) & (TenureFrame$mydata$tenure`[x] <= 60))
    {TenureFrame$mydata$tenure`[x] = "Between 41 and 60"}
  else if (TenureFrame$mydata$tenure`[x] > 60)
    {TenureFrame$mydata$tenure`[x] = "Over 61"}
  x = x + 1
}
```

```
base::table(TenureFrame$mydata$tenure`)
```

```
##
##          0 Between 1 and 20 Between 21 and 40 Between 41 and 60
##          11          1875          1784          1593
##      Over 61
##      1780
```


Let's plot out Tenure data again.

```
ggplot(mydata, aes(TenureFrame$`mydata$tenure`)) + geom_bar(fill = "#56B4E9") + theme_classic() + xlab('Tenure')
```



Now let's use Dummy Vars to create a dataframe with columns containing binary data representing the tenure periods defined. Loading up this data frame, we can see that we have the 5 columns we defined.

```
# let's create dummy var!
```

```
dmy <- dummyVars(" ~ .", data = TenureFrame)
TenureFrame2 <- data.frame(predict(dmy, newdata = TenureFrame))
head(TenureFrame2)
```

```
##   X.mydata.tenure.0 X.mydata.tenure.Between.1.and.20
## 1                0                                1
## 2                0                                0
## 3                0                                1
## 4                0                                0
## 5                0                                1
## 6                0                                0
##   X.mydata.tenure.Between.21.and.40 X.mydata.tenure.Between.41.and.60
## 1                                0                                0
## 2                                1                                0
## 3                                0                                0
## 4                                0                                1
## 5                                0                                0
## 6                                0                                0
##   X.mydata.tenure.Over.61
## 1                        0
```

```
## 2          0
## 3          0
## 4          0
## 5          0
## 6          1
```

Variable Analysis

///ANALYSIS OF SERVICES AVAILABLE///

The following Variables are all related to the Phone and Internet services:

```
Phone Service
MultipleLines
InternetService
OnlineSecurity
OnlineBackup
DeviceProtection
TechSupport
StreamingTV
SreamingMovies
```

We will first check for NAs in each variable. Then we will take a look at the levels of each Variable. No NAs as per the below output. Also, it can be noted that we have a lot of duplicate information. In each service related to Internet service (such as Online backup and Streaming Movies), we have a level specifying “No Internet Service”. This is the same amount (1526) showing “No” in the “Internet Service” Variable.

```
base::table(mydata$PhoneService)
```

```
##
##   No   Yes
## 682 6361
```

```
sum(is.na(mydata$PhoneService))
```

```
## [1] 0
```

```
base::table(mydata$MultipleLines)
```

```
##
##           No No phone service           Yes
##          3390           682          2971
```

```
sum(is.na(mydata$MultipleLines))
```

```
## [1] 0
```

```
base::table(mydata$InternetService)
```

```
##
##           DSL Fiber optic           No
##          2421          3096          1526
```

```
sum(is.na(mydata$InternetService))
```

```
## [1] 0
```

```
base::table(mydata$OnlineSecurity)
```

```
##
##           No No internet service           Yes
##          3498           1526          2019
```

```
sum(is.na(mydata$OnlineSecurity))
```

```
## [1] 0
```

```
base::table(mydata$OnlineBackup)
```

```
##
```

```
##           No No internet service           Yes
##           3088           1526           2429
```

```
sum(is.na(mydata$OnlineBackup))
```

```
## [1] 0
```

```
base::table(mydata$DeviceProtection)
```

```
##
```

```
##           No No internet service           Yes
##           3095           1526           2422
```

```
sum(is.na(mydata$DeviceProtection))
```

```
## [1] 0
```

```
base::table(mydata$TechSupport)
```

```
##
```

```
##           No No internet service           Yes
##           3473           1526           2044
```

```
sum(is.na(mydata$TechSupport))
```

```
## [1] 0
```

```
base::table(mydata$StreamingTV)
```

```
##
```

```
##           No No internet service           Yes
##           2810           1526           2707
```

```
sum(is.na(mydata$StreamingTV))
```

```
## [1] 0
```

```
base::table(mydata$StreamingMovies)
```

```
##
```

```
##           No No internet service           Yes
##           2785           1526           2732
```

```
sum(is.na(mydata$StreamingMovies))
```

```
## [1] 0
```

We will now use Dummy Variables to create a dataframe with binary data on each of these Variables.

!We must then remove columns containing duplicate data!

```
ServiceFrame <- cbind.data.frame(mydata$PhoneService, mydata$MultipleLines, mydata$InternetService, myda
ncol(ServiceFrame)
```

```
## [1] 9
```

```
dmy <- dummyVars(" ~ .", data = ServiceFrame)
ServiceFrame2 <- data.frame(predict(dmy, newdata = ServiceFrame))
head(ServiceFrame2 )
```

```
##      X.mydata.PhoneService.No X.mydata.PhoneService.Yes
## 1                      1                      0
## 2                      0                      1
## 3                      0                      1
## 4                      1                      0
## 5                      0                      1
## 6                      0                      1
##      X.mydata.MultipleLines.No X.mydata.MultipleLines.No.phone.service
## 1                      0                      1
## 2                      1                      0
## 3                      1                      0
## 4                      0                      1
## 5                      1                      0
## 6                      0                      0
##      X.mydata.MultipleLines.Yes X.mydata.InternetService.DSL
## 1                      0                      1
## 2                      0                      1
## 3                      0                      1
## 4                      0                      1
## 5                      0                      0
## 6                      1                      0
##      X.mydata.InternetService.Fiber.optic X.mydata.InternetService.No
## 1                      0                      0
## 2                      0                      0
## 3                      0                      0
## 4                      0                      0
## 5                      1                      0
## 6                      1                      0
##      X.mydata.OnlineSecurity.No X.mydata.OnlineSecurity.No.internet.service
## 1                      1                      0
## 2                      0                      0
## 3                      0                      0
## 4                      0                      0
## 5                      1                      0
## 6                      1                      0
##      X.mydata.OnlineSecurity.Yes X.mydata.OnlineBackup.No
## 1                      0                      0
## 2                      1                      1
## 3                      1                      0
## 4                      1                      1
## 5                      0                      1
## 6                      0                      1
##      X.mydata.OnlineBackup.No.internet.service X.mydata.OnlineBackup.Yes
## 1                      0                      1
## 2                      0                      0
## 3                      0                      1
## 4                      0                      0
## 5                      0                      0
## 6                      0                      0
##      X.mydata.DeviceProtection.No
```

```

## 1          1
## 2          0
## 3          1
## 4          0
## 5          1
## 6          0
##  X.mydata.DeviceProtection.No.internet.service
## 1          0
## 2          0
## 3          0
## 4          0
## 5          0
## 6          0
##  X.mydata.DeviceProtection.Yes X.mydata.TechSupport.No
## 1          0          1
## 2          1          1
## 3          0          1
## 4          1          0
## 5          0          1
## 6          1          1
##  X.mydata.TechSupport.No.internet.service X.mydata.TechSupport.Yes
## 1          0          0
## 2          0          0
## 3          0          0
## 4          0          1
## 5          0          0
## 6          0          0
##  X.mydata.StreamingTV.No X.mydata.StreamingTV.No.internet.service
## 1          1          0
## 2          1          0
## 3          1          0
## 4          1          0
## 5          1          0
## 6          0          0
##  X.mydata.StreamingTV.Yes X.mydata.StreamingMovies.No
## 1          0          1
## 2          0          1
## 3          0          1
## 4          0          1
## 5          0          1
## 6          1          0
##  X.mydata.StreamingMovies.No.internet.service
## 1          0
## 2          0
## 3          0
## 4          0
## 5          0
## 6          0
##  X.mydata.StreamingMovies.Yes
## 1          0
## 2          0
## 3          0
## 4          0
## 5          0

```

```
## 6
```

```
1
```

It can be observed that there are a lot of columns with duplicate information. The following are examples:
X.mydata.PhoneService.No X.mydata.MultipleLines.No.phone.service We will remove the redunant columns:

```
ServiceFrame2ColVec <-colnames(ServiceFrame2)
ServiceFrame2ColVec

## [1] "X.mydata.PhoneService.No"
## [2] "X.mydata.PhoneService.Yes"
## [3] "X.mydata.MultipleLines.No"
## [4] "X.mydata.MultipleLines.No.phone.service"
## [5] "X.mydata.MultipleLines.Yes"
## [6] "X.mydata.InternetService.DSL"
## [7] "X.mydata.InternetService.Fiber.optic"
## [8] "X.mydata.InternetService.No"
## [9] "X.mydata.OnlineSecurity.No"
## [10] "X.mydata.OnlineSecurity.No.internet.service"
## [11] "X.mydata.OnlineSecurity.Yes"
## [12] "X.mydata.OnlineBackup.No"
## [13] "X.mydata.OnlineBackup.No.internet.service"
## [14] "X.mydata.OnlineBackup.Yes"
## [15] "X.mydata.DeviceProtection.No"
## [16] "X.mydata.DeviceProtection.No.internet.service"
## [17] "X.mydata.DeviceProtection.Yes"
## [18] "X.mydata.TechSupport.No"
## [19] "X.mydata.TechSupport.No.internet.service"
## [20] "X.mydata.TechSupport.Yes"
## [21] "X.mydata.StreamingTV.No"
## [22] "X.mydata.StreamingTV.No.internet.service"
## [23] "X.mydata.StreamingTV.Yes"
## [24] "X.mydata.StreamingMovies.No"
## [25] "X.mydata.StreamingMovies.No.internet.service"
## [26] "X.mydata.StreamingMovies.Yes"

ServiceFrame2 <-subset(ServiceFrame2, select =-X.mydata.MultipleLines.No.phone.service)
ServiceFrame2 <-subset(ServiceFrame2, select =-X.mydata.OnlineSecurity.No.internet.service)
ServiceFrame2 <-subset(ServiceFrame2, select =-X.mydata.TechSupport.No.internet.service)
ServiceFrame2 <-subset(ServiceFrame2, select =-X.mydata.StreamingTV.No.internet.service)
ServiceFrame2 <-subset(ServiceFrame2, select =-X.mydata.StreamingMovies.No.internet.service)
ServiceFrame2 <-subset(ServiceFrame2, select =-X.mydata.OnlineBackup.No.internet.service)
ServiceFrame2 <-subset(ServiceFrame2, select =-X.mydata.DeviceProtection.No.internet.service)

head(ServiceFrame2)

##   X.mydata.PhoneService.No X.mydata.PhoneService.Yes
## 1                        1                        0
## 2                        0                        1
## 3                        0                        1
## 4                        1                        0
## 5                        0                        1
## 6                        0                        1
##   X.mydata.MultipleLines.No X.mydata.MultipleLines.Yes
## 1                        0                        0
## 2                        1                        0
## 3                        1                        0
```

## 4	0	0
## 5	1	0
## 6	0	1
##	X.mydata.InternetService.DSL	X.mydata.InternetService.Fiber.optic
## 1	1	0
## 2	1	0
## 3	1	0
## 4	1	0
## 5	0	1
## 6	0	1
##	X.mydata.InternetService.No	X.mydata.OnlineSecurity.No
## 1	0	1
## 2	0	0
## 3	0	0
## 4	0	0
## 5	0	1
## 6	0	1
##	X.mydata.OnlineSecurity.Yes	X.mydata.OnlineBackup.No
## 1	0	0
## 2	1	1
## 3	1	0
## 4	1	1
## 5	0	1
## 6	0	1
##	X.mydata.OnlineBackup.Yes	X.mydata.DeviceProtection.No
## 1	1	1
## 2	0	0
## 3	1	1
## 4	0	0
## 5	0	1
## 6	0	0
##	X.mydata.DeviceProtection.Yes	X.mydata.TechSupport.No
## 1	0	1
## 2	1	1
## 3	0	1
## 4	1	0
## 5	0	1
## 6	1	1
##	X.mydata.TechSupport.Yes	X.mydata.StreamingTV.No
## 1	0	1
## 2	0	1
## 3	0	1
## 4	1	1
## 5	0	1
## 6	0	0
##	X.mydata.StreamingTV.Yes	X.mydata.StreamingMovies.No
## 1	0	1
## 2	0	1
## 3	0	1
## 4	0	1
## 5	0	1
## 6	1	0
##	X.mydata.StreamingMovies.Yes	
## 1	0	

```
## 2          0
## 3          0
## 4          0
## 5          0
## 6          1
```

```
ncol(ServiceFrame2)
```

```
## [1] 19
```

Variable Analysis

```
///CONTRACT///
```

Let's take a look at contract. We have no NAs, and we have 3 levels. Again, lets use Dummy Variables to creata a dataframe with binary data.

```
base::table(mydata$Contract)
```

```
##
## Month-to-month      One year      Two year
##           3875           1473           1695
```

```
sum(is.na(mydata$Contract))
```

```
## [1] 0
```

```
ContractFrame <- cbind.data.frame(mydata$Contract)
```

```
dmy <- dummyVars(" ~ .", data = ContractFrame)
```

```
ContractFrame2 <- data.frame(predict(dmy, newdata = ContractFrame))
```

```
head(ContractFrame2)
```

```
##      X.mydata.Contract.Month.to.month X.mydata.Contract.One.year
## 1                                1                                0
## 2                                0                                1
## 3                                1                                0
## 4                                0                                1
## 5                                1                                0
## 6                                1                                0
##      X.mydata.Contract.Two.year
## 1                                0
## 2                                0
## 3                                0
## 4                                0
## 5                                0
## 6                                0
```

Variable Analysis

```
///PAPERLESS BILLING///
```

We have no NAs in this variable. Let's create a column with binary data.

```
base::table(mydata$PaperlessBilling)
```

```
##
##      No  Yes
## 2872 4171
```

```
sum(is.na(mydata$PaperlessBilling))
```

```
## [1] 0
```



```
mydata$PaperlessBillingLog[mydata$PaperlessBilling == "Yes"] = 1
```

```
## Warning: Unknown or uninitialised column: 'PaperlessBillingLog'.
```

```
mydata$PaperlessBillingLog[mydata$PaperlessBilling == "No"] = 0
```

```
head(mydata)
```

```
## # A tibble: 6 x 25
##   customerID gender SeniorCitizen Partner Dependents tenure PhoneService
##   <chr>      <chr>          <dbl> <chr>    <chr>          <dbl> <chr>
## 1 7590-VHVEG Female            0 Yes     No              1 No
## 2 5575-GNVDE Male              0 No     No             34 Yes
## 3 3668-QPYBK Male              0 No     No              2 Yes
## 4 7795-CFOCW Male              0 No     No             45 No
## 5 9237-HQITU Female            0 No     No              2 Yes
## 6 9305-CDSKC Female            0 No     No              8 Yes
## # ... with 18 more variables: MultipleLines <chr>, InternetService <chr>,
## #   OnlineSecurity <chr>, OnlineBackup <chr>, DeviceProtection <chr>,
## #   TechSupport <chr>, StreamingTV <chr>, StreamingMovies <chr>,
## #   Contract <chr>, PaperlessBilling <chr>, PaymentMethod <chr>,
## #   MonthlyCharges <dbl>, TotalCharges <dbl>, Churn <chr>, ChurnLog <dbl>,
## #   PartnerLog <dbl>, DependentsLog <dbl>, PaperlessBillingLog <dbl>
```

Variable Analysis

///PAYMENT METHOD///

Let's take a look at payment method. We have no NAs, and 4 levels. Let's create another dataframe with binary values.

```
base::table(mydata$PaymentMethod)
```

```
##
## Bank transfer (automatic)    Credit card (automatic)
##                1544                1522
##      Electronic check          Mailed check
##                2365                1612
```

```
sum(is.na(mydata$PaymentMethod))
```

```
## [1] 0
```

```
PaymentFrame <- as.data.frame(mydata$PaymentMethod)
```

```
dmy <- dummyVars(" ~ .", data = PaymentFrame)
PaymentFrame2 <- data.frame(predict(dmy, newdata = PaymentFrame))
head(PaymentFrame2)
```

```
##   X.mydata.PaymentMethod.Bank.transfer..automatic.
## 1                                                    0
## 2                                                    0
## 3                                                    0
## 4                                                    1
## 5                                                    0
## 6                                                    0
##   X.mydata.PaymentMethod.Credit.card..automatic.
## 1                                                    0
```

```
## 2 0
## 3 0
## 4 0
## 5 0
## 6 0
## X.mydata.PaymentMethod.Electronic.check
## 1 1
## 2 0
## 3 0
## 4 0
## 5 1
## 6 1
## X.mydata.PaymentMethod.Mailed.check
## 1 0
## 2 1
## 3 1
## 4 0
## 5 0
## 6 0
```

Variable Analysis

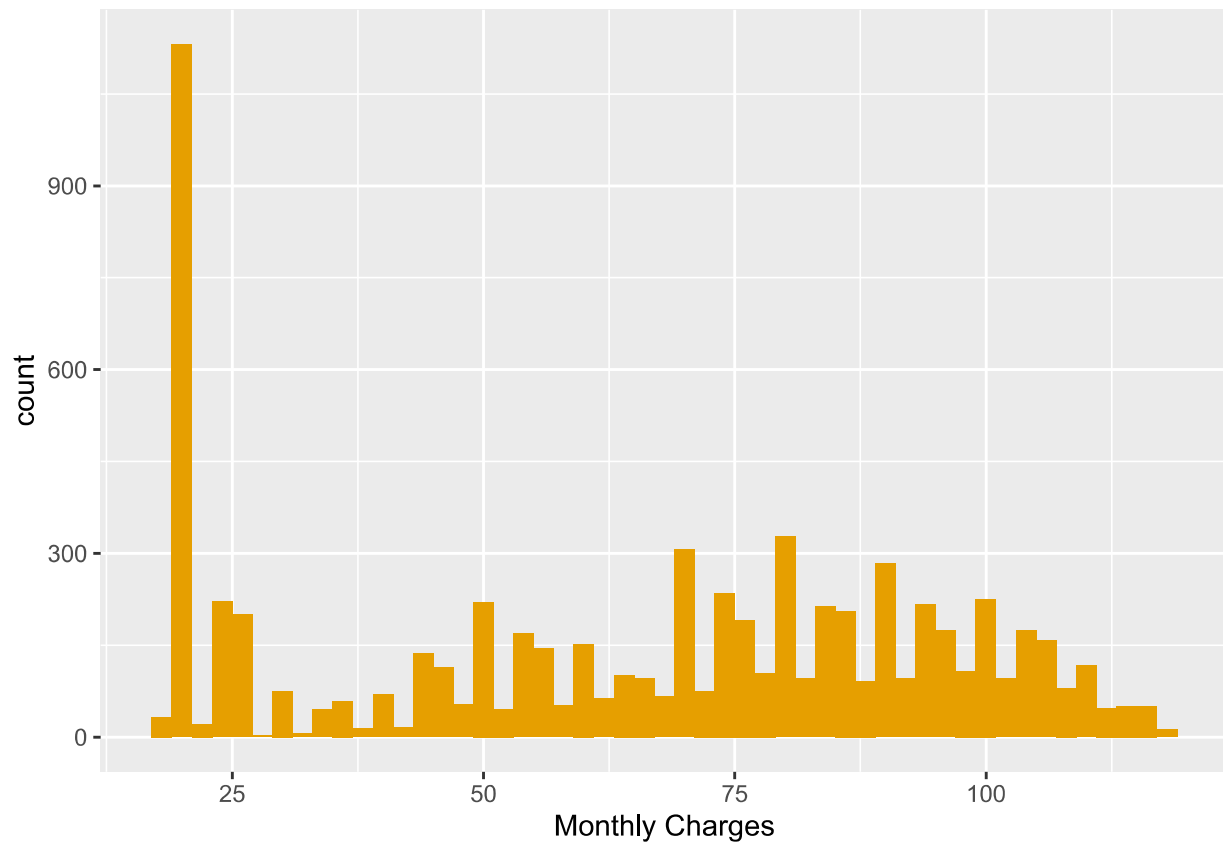
///MONTHLY CHARGES///

Let's look at monthly charges. We have no NAs. Talking a look at the Violin graph, it appears that those with high monthly charges are more likely to churn.

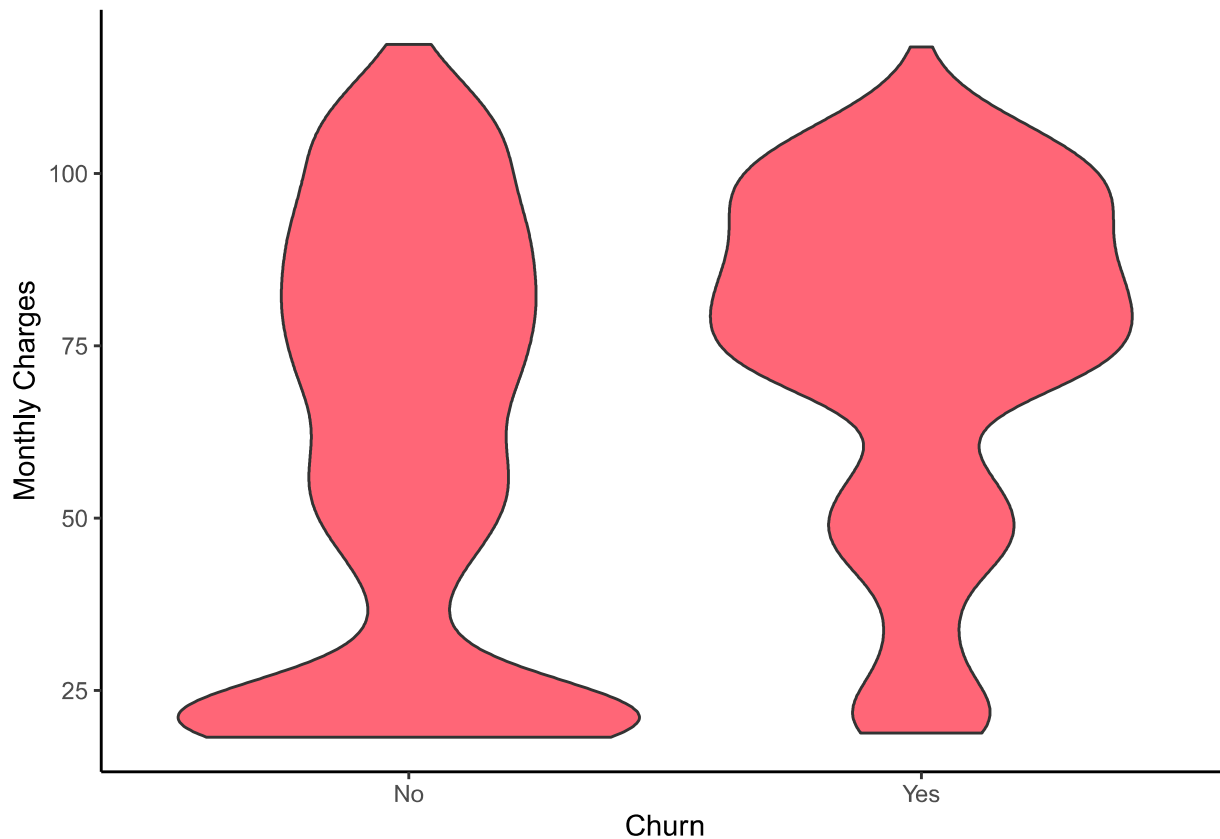
```
sum(is.na(mydata$MonthlyCharges))
```

```
## [1] 0
```

```
ggplot(mydata, aes(mydata$MonthlyCharges)) + geom_histogram(binwidth = 2, fill = "#E69F00") + xlab("Mont
```



```
ggplot(mydata, aes(mydata$Churn, mydata$MonthlyCharges)) + geom_violin(fill = "#FF6677") + xlab("Churn")
```



Variable Analysis

///TOTAL CHARGES///

Let's take a look at total charges. We have 11 NAs.

```
sum(is.na(mydata$TotalCharges))
```

```
## [1] 11
```

-> *VALUE IMPUTATION*

We need to replace the 11 Na values that we have in the "Total Charges" Column.

Let's take a look at the relationship between Monthly Charges and Total Charges. Let's create a column to Analyse this relationship. We will divide the Total Charges by the tenure and verify if it is similar to the Monthly Charges. Indeed they are very similar. In this regard, for the values which are missing from "Total Charges" we can simply multiply the monthly charges with the tenure.

```
mydata$TotalChargesDivTenure <- (mydata$TotalCharges/mydata$tenure)
```

```
test <-cbind.data.frame(mydata$MonthlyCharges,mydata$TotalChargesDivTenure)
head(test)
```

```
##   mydata$MonthlyCharges mydata$TotalChargesDivTenure
## 1                29.85                29.85000
## 2                56.95                55.57353
## 3                53.85                54.07500
## 4                42.30                40.90556
## 5                70.70                75.82500
## 6                99.65                102.56250
```

```
## Missing value replacement
```

```
x=1
```

```
while (x <= nrow(mydata))
```

```
{
```

```
  if ((is.na(mydata$TotalCharges[x])) == TRUE)
```

```
  { mydata$TotalCharges[x] <- (mydata$MonthlyCharges[x]*mydata$tenure[x])}
```

```
  x = x + 1
```

```
}
```

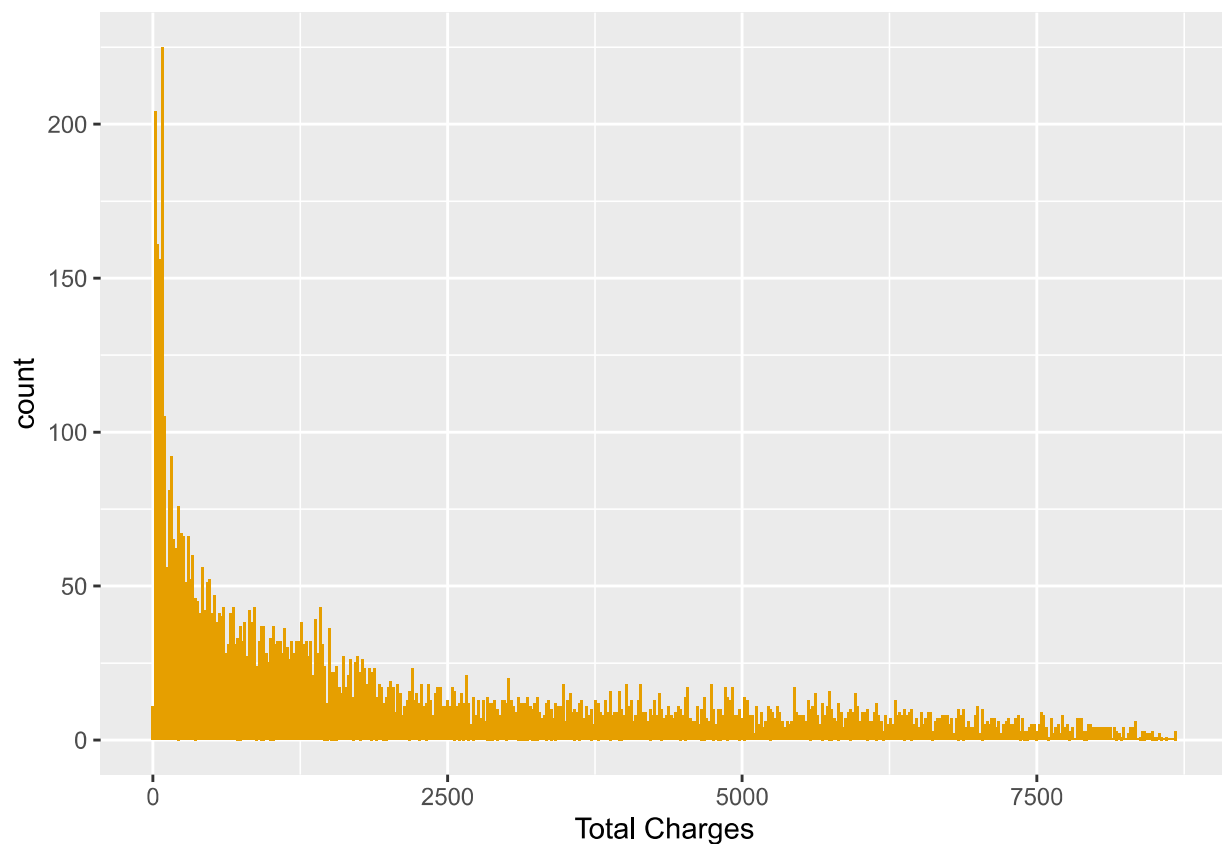
Check NAs. Now that NAs are gone, we can take a better look at our Total Charges Variable.

```
sum(is.na(mydata$TotalCharges))
```

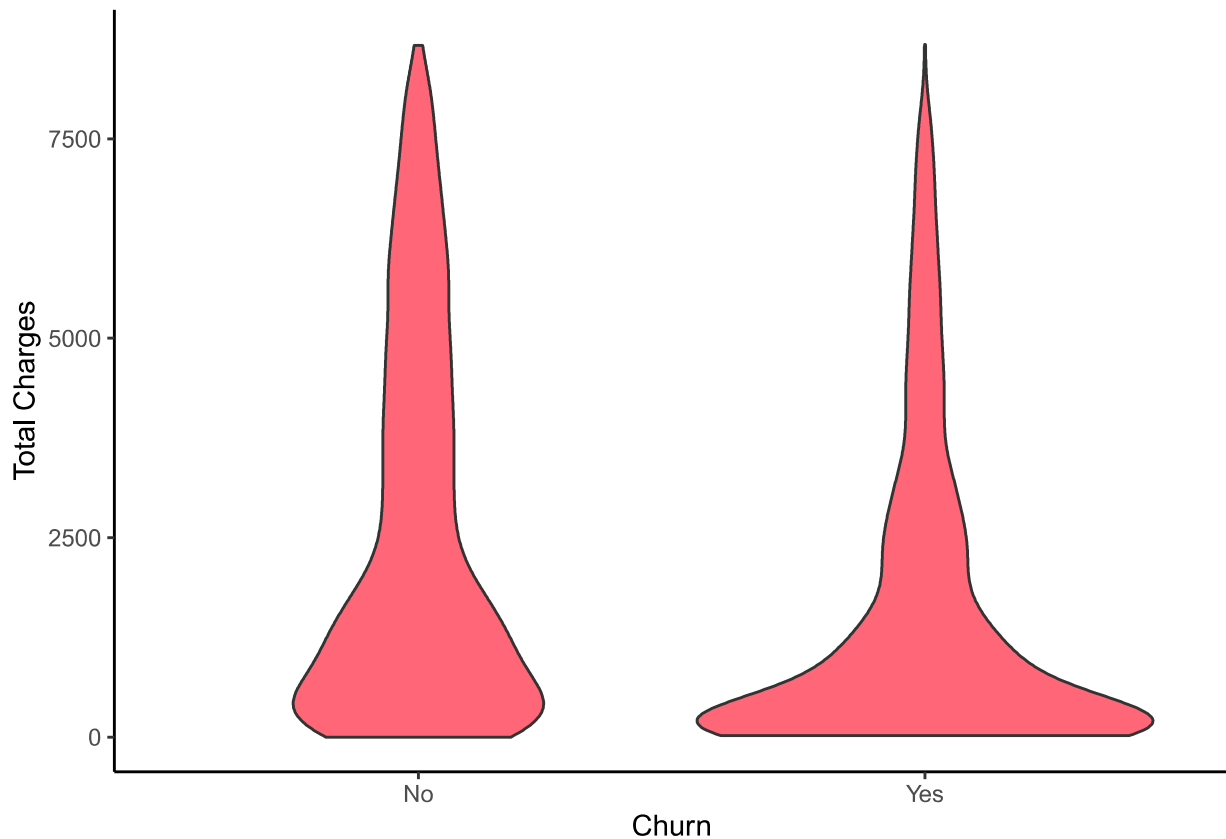
```
## [1] 0
```

Let's take a look at our Total Charges:

```
ggplot(mydata, aes(mydata$TotalCharges)) + geom_histogram(binwidth = 20, fill = "#E69F00") + xlab("Total
```



```
ggplot(mydata, aes(mydata$Churn, mydata$TotalCharges)) + geom_violin(fill = "#FF6677") + xlab("Churn") +
```



Feature Selection

We will now select which features to include into our model. Let's check the column names from the original dataframe "mydata".

- 'customerID' - No inclusion
- 'gender' - No inclusion
- 'SeniorCitizen' - Include
- 'Partner' - No inclusion
- 'Dependents' - No inclusion
- 'tenure' - Include
- 'PhoneService' - No inclusion
- 'MultipleLines' - No inclusion
- 'InternetService' - No inclusion
- 'OnlineSecurity' - No inclusion
- 'OnlineBackup' - No inclusion
- 'DeviceProtection' - No inclusion
- 'TechSupport' - No inclusion
- 'StreamingTV' - No inclusion
- 'StreamingMovies' - No inclusion
- 'Contract' - No inclusion
- 'PaperlessBilling' - No inclusion
- 'PaymentMethod' - No inclusion
- 'MonthlyCharges' - Include
- 'TotalCharges' - Include
- 'Churn' - No inclusion
- 'ChurnLog' - Include
- 'PartnerLog' - Include

- 'DependentsLog' - Include
- 'PaperlessBillingLog' - Include
- 'TotalChargesDivTenure' - No inclusion

We will bind these columns with the dataframes we created: * PaymentFrame2 * TenureFrame2 * Contract-Frame2 * Genderframe2 * ServiceFrame2

We will name our data with the selected features: newdata

```
head(mydata)
```

```
## # A tibble: 6 x 26
##   customerID gender SeniorCitizen Partner Dependents tenure PhoneService
##   <chr>      <chr>          <dbl> <chr>    <chr>          <dbl> <chr>
## 1 7590-VHVEG Female            0 Yes     No            1 No
## 2 5575-GNVDE Male              0 No      No           34 Yes
## 3 3668-QPYBK Male              0 No      No            2 Yes
## 4 7795-CFOCW Male              0 No      No           45 No
## 5 9237-HQITU Female            0 No      No            2 Yes
## 6 9305-CDSKC Female            0 No      No            8 Yes
## # ... with 19 more variables: MultipleLines <chr>, InternetService <chr>,
## #   OnlineSecurity <chr>, OnlineBackup <chr>, DeviceProtection <chr>,
## #   TechSupport <chr>, StreamingTV <chr>, StreamingMovies <chr>,
## #   Contract <chr>, PaperlessBilling <chr>, PaymentMethod <chr>,
## #   MonthlyCharges <dbl>, TotalCharges <dbl>, Churn <chr>, ChurnLog <dbl>,
## #   PartnerLog <dbl>, DependentsLog <dbl>, PaperlessBillingLog <dbl>,
## #   TotalChargesDivTenure <dbl>
```

```
ColVec <- colnames(mydata)
ColVec
```

```
## [1] "customerID"      "gender"
## [3] "SeniorCitizen"   "Partner"
## [5] "Dependents"      "tenure"
## [7] "PhoneService"    "MultipleLines"
## [9] "InternetService" "OnlineSecurity"
## [11] "OnlineBackup"    "DeviceProtection"
## [13] "TechSupport"     "StreamingTV"
## [15] "StreamingMovies" "Contract"
## [17] "PaperlessBilling" "PaymentMethod"
## [19] "MonthlyCharges"  "TotalCharges"
## [21] "Churn"           "ChurnLog"
## [23] "PartnerLog"      "DependentsLog"
## [25] "PaperlessBillingLog" "TotalChargesDivTenure"
```

```
newdata <- cbind.data.frame(mydata$SeniorCitizen, mydata$PartnerLog, mydata$DependentsLog, mydata$tenure,
head(newdata)
```

```
##   mydata$SeniorCitizen mydata$PartnerLog mydata$DependentsLog
## 1                    0                    1                    0
## 2                    0                    0                    0
## 3                    0                    0                    0
## 4                    0                    0                    0
## 5                    0                    0                    0
## 6                    0                    0                    0
##   mydata$tenure X.mydata.gender.Female X.mydata.gender.Male
## 1              1                    1                    0
```

## 2	34	0	1
## 3	2	0	1
## 4	45	0	1
## 5	2	1	0
## 6	8	1	0
## X.mydata.PaymentMethod.Bank.transfer..automatic.			
## 1		0	
## 2		0	
## 3		0	
## 4		1	
## 5		0	
## 6		0	
## X.mydata.PaymentMethod.Credit.card..automatic.			
## 1		0	
## 2		0	
## 3		0	
## 4		0	
## 5		0	
## 6		0	
## X.mydata.PaymentMethod.Electronic.check			
## 1		1	
## 2		0	
## 3		0	
## 4		0	
## 5		1	
## 6		1	
## X.mydata.PaymentMethod.Mailed.check X.mydata.PhoneService.No			
## 1		0	1
## 2		1	0
## 3		1	0
## 4		0	1
## 5		0	0
## 6		0	0
## X.mydata.PhoneService.Yes X.mydata.MultipleLines.No			
## 1	0		0
## 2	1		1
## 3	1		1
## 4	0		0
## 5	1		1
## 6	1		0
## X.mydata.MultipleLines.Yes X.mydata.InternetService.DSL			
## 1	0		1
## 2	0		1
## 3	0		1
## 4	0		1
## 5	0		0
## 6	1		0
## X.mydata.InternetService.Fiber.optic X.mydata.InternetService.No			
## 1		0	0
## 2		0	0
## 3		0	0
## 4		0	0
## 5		1	0
## 6		1	0

	X.mydata.OnlineSecurity.No	X.mydata.OnlineSecurity.Yes
## 1	1	0
## 2	0	1
## 3	0	1
## 4	0	1
## 5	1	0
## 6	1	0

	X.mydata.OnlineBackup.No	X.mydata.OnlineBackup.Yes
## 1	0	1
## 2	1	0
## 3	0	1
## 4	1	0
## 5	1	0
## 6	1	0

	X.mydata.DeviceProtection.No	X.mydata.DeviceProtection.Yes
## 1	1	0
## 2	0	1
## 3	1	0
## 4	0	1
## 5	1	0
## 6	0	1

	X.mydata.TechSupport.No	X.mydata.TechSupport.Yes	X.mydata.StreamingTV.No
## 1	1	0	1
## 2	1	0	1
## 3	1	0	1
## 4	0	1	1
## 5	1	0	1
## 6	1	0	0

	X.mydata.StreamingTV.Yes	X.mydata.StreamingMovies.No
## 1	0	1
## 2	0	1
## 3	0	1
## 4	0	1
## 5	0	1
## 6	1	0

	X.mydata.StreamingMovies.Yes	X.mydata.tenure.0
## 1	0	0
## 2	0	0
## 3	0	0
## 4	0	0
## 5	0	0
## 6	1	0

	X.mydata.tenure.Between.1.and.20	X.mydata.tenure.Between.21.and.40
## 1	1	0
## 2	0	1
## 3	1	0
## 4	0	0
## 5	1	0
## 6	0	0

	X.mydata.tenure.Between.41.and.60	X.mydata.tenure.Over.61
## 1	0	0
## 2	0	0
## 3	0	0
## 4	1	0

```
## 5          0          0
## 6          0          1
##  mydata$PaperlessBillingLog X.mydata.Contract.Month.to.month
## 1          1          1
## 2          0          0
## 3          1          1
## 4          0          0
## 5          1          1
## 6          1          1
##  X.mydata.Contract.One.year X.mydata.Contract.Two.year
## 1          0          0
## 2          1          0
## 3          0          0
## 4          1          0
## 5          0          0
## 6          0          0
##  mydata$MonthlyCharges mydata$TotalCharges mydata$ChurnLog
## 1          29.85          29.85          0
## 2          56.95         1889.50          0
## 3          53.85          108.15          1
## 4          42.30         1840.75          0
## 5          70.70          151.65          1
## 6          99.65          820.50          1
```

Standardizing the Data

This will give the variables zero-mean and unit-variance. This is required for interpretation for a Logistic Regression model.

```
x = 1

while (x<=length(colnames(newdata))) {
  newdata[,x] <- scale(as.numeric(newdata[,x]), center = TRUE, scale = TRUE))
  x = x+1
}

newdata$`mydata$ChurnLog` <-mydata$ChurnLog

head(newdata)

##  mydata$SeniorCitizen mydata$PartnerLog mydata$DependentsLog
## 1          -0.4398853          1.0344568         -0.6539655
## 2          -0.4398853         -0.9665537         -0.6539655
## 3          -0.4398853         -0.9665537         -0.6539655
## 4          -0.4398853         -0.9665537         -0.6539655
## 5          -0.4398853         -0.9665537         -0.6539655
## 6          -0.4398853         -0.9665537         -0.6539655
##  mydata$tenure X.mydata.gender.Female X.mydata.gender.Male
## 1   -1.27735389          1.0094870         -1.0094870
## 2    0.06632271         -0.9904615          0.9904615
## 3   -1.23663642         -0.9904615          0.9904615
## 4    0.51421491         -0.9904615          0.9904615
## 5   -1.23663642          1.0094870         -1.0094870
## 6   -0.99233158          1.0094870         -1.0094870
##  X.mydata.PaymentMethod.Bank.transfer..automatic.
## 1                                     -0.5298476
```

```

## 2 -0.5298476
## 3 -0.5298476
## 4 1.8870673
## 5 -0.5298476
## 6 -0.5298476
## X.mydata.PaymentMethod.Credit.card..automatic.
## 1 -0.5250101
## 2 -0.5250101
## 3 -0.5250101
## 4 -0.5250101
## 5 -0.5250101
## 6 -0.5250101
## X.mydata.PaymentMethod.Electronic.check
## 1 1.4063185
## 2 -0.7109755
## 3 -0.7109755
## 4 -0.7109755
## 5 1.4063185
## 6 1.4063185
## X.mydata.PaymentMethod.Mailed.check X.mydata.PhoneService.No
## 1 -0.5447682 3.0537936
## 2 1.8353823 -0.3274151
## 3 1.8353823 -0.3274151
## 4 -0.5447682 3.0537936
## 5 -0.5447682 -0.3274151
## 6 -0.5447682 -0.3274151
## X.mydata.PhoneService.Yes X.mydata.MultipleLines.No
## 1 -3.0537936 -0.9632614
## 2 0.3274151 1.0379924
## 3 0.3274151 1.0379924
## 4 -3.0537936 -0.9632614
## 5 0.3274151 1.0379924
## 6 0.3274151 -0.9632614
## X.mydata.MultipleLines.Yes X.mydata.InternetService.DSL
## 1 -0.8541155 1.3816141
## 2 -0.8541155 1.3816141
## 3 -0.8541155 1.3816141
## 4 -0.8541155 1.3816141
## 5 -0.8541155 -0.7236884
## 6 1.1706356 -0.7236884
## X.mydata.InternetService.Fiber.optic X.mydata.InternetService.No
## 1 -0.8855969 -0.52589
## 2 -0.8855969 -0.52589
## 3 -0.8855969 -0.52589
## 4 -0.8855969 -0.52589
## 5 1.1290216 -0.52589
## 6 1.1290216 -0.52589
## X.mydata.OnlineSecurity.No X.mydata.OnlineSecurity.Yes
## 1 1.0066242 -0.633888
## 2 -0.9932783 1.577342
## 3 -0.9932783 1.577342
## 4 -0.9932783 1.577342
## 5 1.0066242 -0.633888
## 6 1.0066242 -0.633888

```

```

## X.mydata.OnlineBackup.No X.mydata.OnlineBackup.Yes
## 1 -0.883557 1.3781427
## 2 1.131628 -0.7255112
## 3 -0.883557 1.3781427
## 4 1.131628 -0.7255112
## 5 1.131628 -0.7255112
## 6 1.131628 -0.7255112
## X.mydata.DeviceProtection.No X.mydata.DeviceProtection.Yes
## 1 1.1293470 -0.7239161
## 2 -0.8853417 1.3811794
## 3 1.1293470 -0.7239161
## 4 -0.8853417 1.3811794
## 5 1.1293470 -0.7239161
## 6 -0.8853417 1.3811794
## X.mydata.TechSupport.No X.mydata.TechSupport.Yes X.mydata.StreamingTV.No
## 1 1.013797 -0.6393932 1.2272701
## 2 1.013797 -0.6393932 1.2272701
## 3 1.013797 -0.6393932 1.2272701
## 4 -0.986251 1.5637607 1.2272701
## 5 1.013797 -0.6393932 1.2272701
## 6 1.013797 -0.6393932 -0.8147009
## X.mydata.StreamingTV.Yes X.mydata.StreamingMovies.No
## 1 -0.7900756 1.2364011
## 2 -0.7900756 1.2364011
## 3 -0.7900756 1.2364011
## 4 -0.7900756 1.2364011
## 5 -0.7900756 1.2364011
## 6 1.2655219 -0.8086842
## X.mydata.StreamingMovies.Yes X.mydata.tenure.0
## 1 -0.7960136 -0.03954814
## 2 -0.7960136 -0.03954814
## 3 -0.7960136 -0.03954814
## 4 -0.7960136 -0.03954814
## 5 -0.7960136 -0.03954814
## 6 1.2560815 -0.03954814
## X.mydata.tenure.Between.1.and.20 X.mydata.tenure.Between.21.and.40
## 1 1.660083 -0.5823915
## 2 -0.602294 1.7168143
## 3 1.660083 -0.5823915
## 4 -0.602294 -0.5823915
## 5 1.660083 -0.5823915
## 6 -0.602294 -0.5823915
## X.mydata.tenure.Between.41.and.60 X.mydata.tenure.Over.61
## 1 -0.5406034 -0.5815171
## 2 -0.5406034 -0.5815171
## 3 -0.5406034 -0.5815171
## 4 1.8495221 -0.5815171
## 5 -0.5406034 -0.5815171
## 6 -0.5406034 1.7193958
## mydata$PaperlessBillingLog X.mydata.Contract.Month.to.month
## 1 0.8297386 0.9041196
## 2 -1.2050277 -1.1058913
## 3 0.8297386 0.9041196
## 4 -1.2050277 -1.1058913

```

```
## 5          0.8297386          0.9041196
## 6          0.8297386          0.9041196
##   X.mydata.Contract.One.year X.mydata.Contract.Two.year
## 1          -0.5142129          -0.5629351
## 2           1.9444438          -0.5629351
## 3          -0.5142129          -0.5629351
## 4           1.9444438          -0.5629351
## 5          -0.5142129          -0.5629351
## 6          -0.5142129          -0.5629351
##   mydata$MonthlyCharges mydata$TotalCharges mydata$ChurnLog
## 1          -1.1602405          -0.9925401           0
## 2          -0.2596105          -0.1721525           0
## 3          -0.3626346          -0.9579979           1
## 4          -0.7464825          -0.1936586           0
## 5           0.1973512          -0.9388078           1
## 6           1.1594634          -0.6437435           1
```

```
bound = 0.5
train = newdata[1:(bound*(nrow(newdata))),]
test = newdata[((nrow(mydata))*bound)+1:nrow(newdata),]
print("The number of training rows are")
```

```
## [1] "The number of training rows are"
```

```
nrow(train)
```

```
## [1] 3521
```

```
print("The number of test rows are")
```

```
## [1] "The number of test rows are"
```

```
nrow(test)
```

```
## [1] 3521
```

SECTION 3: MODEL BUILDING

We will use Logistic Regression as the machine learning model.

```
model <- glm( as.numeric(`mydata$ChurnLog`) ~ .
              ,family=binomial(link='logit'),data=train, control = list(maxit = 200))
```

SECTION 4: PREDICTION AND ACCURACY EVALUATION

Make a prediction

We will now use our test data to make a prediction.

```
predictChurn <- predict(model, newdata = test,type='response')
```

```
## Warning in predict.lm(object, newdata, se.fit, scale = 1, type = if (type
## == : prediction from a rank-deficient fit may be misleading
```

Let's round our prediction up, since we want a 1 or 0 as an output.

```
RoundpredictChurn <- round(predictChurn)
head(RoundpredictChurn)
```

```
## 3522 3523 3524 3525 3526 3527
##    0    0    0    1    0    0
```

Measure Accuracy

```
conf_mat <- base::table(test$`mydata$ChurnLog`,RoundpredictChurn)
conf_mat
```

```
##      RoundpredictChurn
##           0         1
##    0 2298    271
##    1   440    512
```

```
cat("Test accuracy: ", sum(diag(conf_mat))/sum(conf_mat))
```

```
## Test accuracy:  0.7980687
```

SECTION 5: CONCLUSIONS AND FUTURE WORK

Several other machine learning models could be applied such as K Nearest Neighbour or Decision Trees. A comparison of the accuracy of these models could be provided, and the most accurate model would be selected. Logistic regression was chosen since it plots the data on a curve which fits a Probabilistic curve with two binary outputs. This is reflective of customer churn behaviour.