# Churn Analysis of Telecom Subscriber Service Data

Churn is when a customer with a subscription to a service cancels that subscription, membership, account, etc. Understanding why customers churn is critical to subscriber based businesses. Identifying customers that will potentially churn allow companies to design business strategies to encourage those customers to stay and study the effects of such strategies by the customer subsequent likelihood that they won't churn in the future. For that, however, it would require analysis over time which may involve churn rate and other key measures. However, for this analysis I only have the churn data taken for a single month, a brief description of the variables, and nothing else. I don't have information on the prices of each service or whether having a contract gives customers discounts or other benefits or anything of that sort. All assumptions made are going to be made purely on the finding made during the analysis of the data. The goal is to create a profile of the customer likely to churn based on the variable available. I will use the results from each step of the analysis to make certain assumptions to formulate a Churn Customer Profile. With this analysis I also expect to identify areas that may be affecting customer satisfaction resulting in Churn.

The data to be analyzed was taken from an anonymous telecom service provider. The data has 7043 observations. It has 21 columns. Figure 1.1 shows the top six rows of the data.

```
> head(churn.data)
  customerID gender SeniorCitizen Partner Dependents tenure PhoneService    MultipleLines InternetService OnlineSecurity OnlineBackup DeviceProtection TechSupport StreamingTV
1 7590-VHVEG Female            No     Yes         No      1           No No phone service            DSL             No          Yes               No          No          No
2 5575-GNVDE   Male            No      No         No     34          Yes               No            DSL            Yes           No              Yes          No          No
3 3668-QPYBK   Male            No      No         No      2          Yes               No            DSL            Yes          Yes               No          No          No
4 7795-CFOCW   Male            No      No         No     45           No No phone service            DSL            Yes           No              Yes         Yes          No
5 9237-HQITU Female            No      No         No      2          Yes               No     Fiber optic             No           No               No          No          No
6 9305-CDSKC Female            No      No         No      8          Yes              Yes     Fiber optic             No           No              Yes          No         Yes
  StreamingMovies       Contract PaperlessBilling           PaymentMethod MonthlyCharges TotalCharges Churn
1              No Month-to-month              Yes        Electronic check          29.85        29.85    No
2              No       One year               No            Mailed check          56.95      1889.50    No
3              No Month-to-month              Yes            Mailed check          53.85       108.15   Yes
4              No       One year               No Bank transfer (automatic)        42.30      1840.75    No
5              No Month-to-month              Yes        Electronic check          70.70       151.65   Yes
6             Yes Month-to-month              Yes        Electronic check          99.65       820.50   Yes
```

Figure 1

Here is a brief description of the data:

**Churn**: Yes/No. Customers who left within the last month.

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**: Yes/No.

Customer account information – **tenure:** how long they've been a customer in months, **contract**: one year/two year/month-to-month, **payment method:** Electronic check/Bank Transfer/Mail**, paperless billing:** Yes/No**, monthly charges:** in dollars**, and total charges:** in dollars**.**

Demographic info about customers – **gender:** Yes/No, **Senior Citizen:** Yes/No, and if they have **partners:** Yes/No and **dependents:** Yes/No.

Churn is the dependent variable. There are 19 independent variables of which 3 are continuous and 16 are categorical. Let's first look at the continuous variables tenure, MonthlyCharges and TotalCharges. Figure 1.2 has boxplots showing the distribution of values for the three aforementioned variables for customers that churned and did not churn.
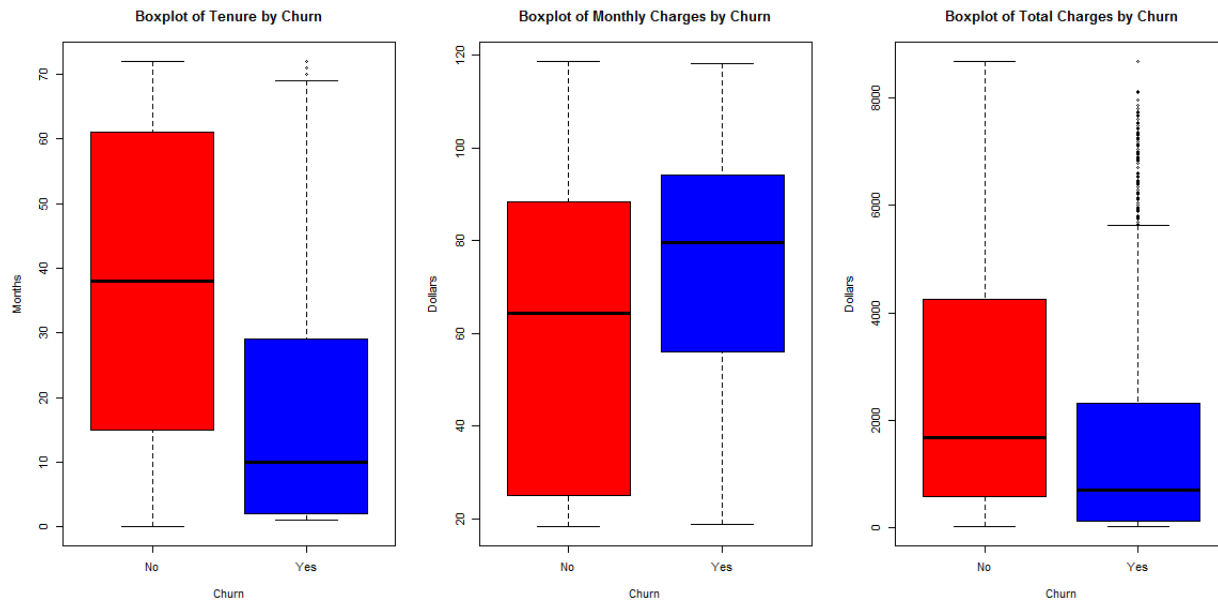
Figure 2

All of the graphs show that the means for each variable differ between those who churned and those who did not. The mean for tenure for customers that churned is 17.98 months while those who did not churn have a much higher mean at 37.57 months. MonthlyCharges also show a considerable difference with those who churned having a mean of $74.44 per month and those who did not having a mean of $61.27 per month. TotalCharges for those who churn is much lower than those who did not churn as it would be expected since those who churn have shorter tenures and therefore less accumulation of monthly payments that make up the total charges. Are these differences statistically significant? I will test that later on but for now let's look at the histogram for these variables in figure 1.3 below.
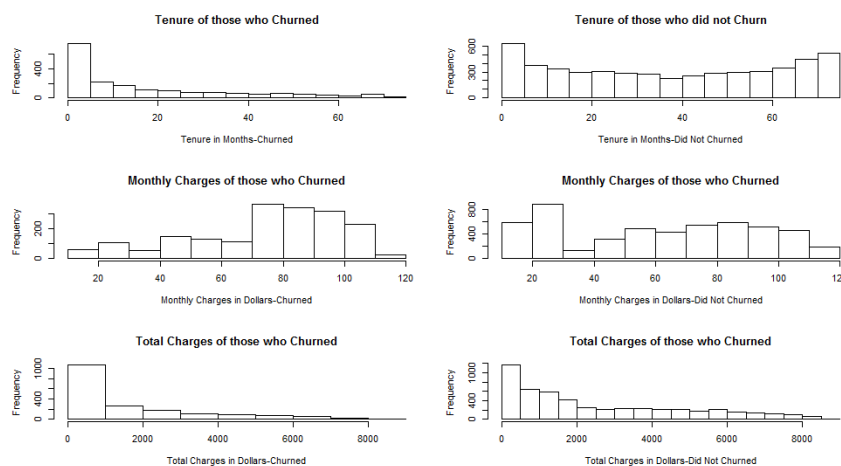


Figure 3

The histograms for the tenure variable show that many of those who churned had short tenures. The contract variable can have 3 values: one year, two year and month-to-month. From the tenure histogram I can assume that those who churn are likely in a month-to-month contract with the freedom to leave at any time while the rest are in contracts for at least a year which will automatically increase

the tenure and the totalcharges. The histogram for Monthly Charges shows that in those who churned there is a higher occurrence of higher monthly bills. That may indicate that since those customers are in a month-to-month contracts they don't get monthly discounts associated with signing one or two year contracts. Again, I am not previewed to that information. I'm making assumptions based on the data. So far my Churn Customer Profile is that of a customer with month-to-month contracts that happen to have higher monthly charges than one or two year contract customers. They don't stay very long and that's why their total charges are so much lower than those who do not churn since those who do not churn tent to have long tenure and therefore accumulate more total charges.

Now let's look at the categorical variables to see if there is something that may interest us. While…



Figure 4

Figure 4 shows that of the 7,043 customers only 1142 are senior citizens or 16.2%. However, of the 1,869 that churned 476 or 25.47% were senior citizens. The difference may be attributed to the higher mortality rate amongst senior citizens.



Figure 5

Even though that in the overall population there are more customers without Partner 3 than with Partner 3,402, of those who did not churn more had partners than did not which may indicate that having a partner may increase the chances of not churning.

**Churn by Having Dependents**                **Churn by Having Dependents**

Figure 6

Something similar can be said about Dependents.  Figure 6 shows that the percentage of customer with dependents in the population of customers that did not churn, 34.5%, is greater than that of the overall population, 30%, indicating that having dependents may increase that the customer will not churn.

**Churn by Having Internet Service**        **Churn by Internet Service**
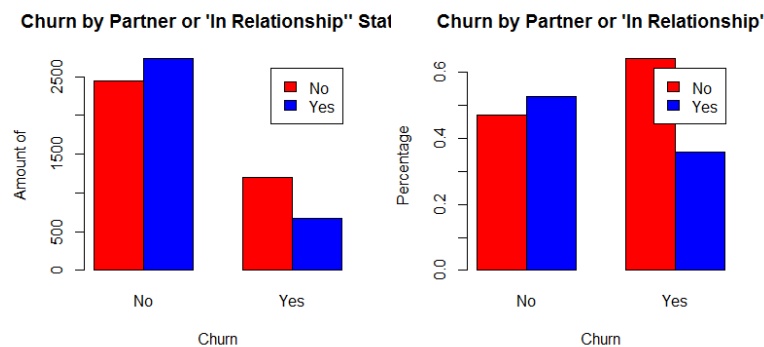
Figure 7

Figure 7 shows how that of those who churned, 94% had internet service either DSL or Fiber optic. 69.4% of those who churned had Fiber optics internet service which may indicate a problem with customer satisfaction with it resulting in customers unhappy enough to churn.

**Churn by Having Online Security**  **Churn by Online Security**  **Churn by Having Online Backup**  **Churn by Having Online Backup**

**Churn by Having Device Protection**  **Churn by Having Device Protection**  **Churn by Having Tech Support**  **Churn by Having Tech Support**

Figure 8

Figure 8 shows how such an abnormal number of customers who churned did not have any security features such as online security, online backup, device protection or tech support.  Figure 7 showed us how so many of those who churned had internet service with fiber optics and figure 2 showed us that

those who churn have shorter tenures and are likely month-to-month contracts.  Could it be that these customers are transient customers looking for a better deal and don't bother with getting extra security or could it be that the customers were persuaded by a great deal in fiber optics internet service and later became unhappy because of lack of security the company neglected to offer them?



Figure 9

Figure 9 confirms our suspicions from figure 3 that many who churned did it while still within the 12 and 24 month periods indicating that they had the freedom to churn because they were in month-to-month contracts.  As seen in the graph, most of those who churned were in month-to-month contracts.



Figure 10

Figure 10 shows us that 75% of those who churned had paperless billing while the percentage of the overall population with paperless billing is 60%.  Paperless billing may be a new feature popular with the customer with short tenures and that were also fiber optic internet customers.



Figure 11

Figure 11 shows that of those who churned most were enrolled in Electronic check.  There seems to be an issue with customer satisfaction with new services such as fiber optics internet service, electronic payments and paperless billing.

After this thorough descriptive analysis of the independent variables I am ready to update my Churn Customer profile.  The Churn customer is one that will opt for the freedom of a month-to-month contract even though it may cost more monthly even without security features.  The customer may or may not have been attracted by an offer in fiber optic internet service at which point the customer

signed up for paperless billing and electronic check payments.  From the description of the customer likely to churn we can deduce the customer not likely to churn.  The customer less likely to churn is one that prefers the long commitment of a one or two year contract and the lower monthly payments.  This customer is more likely to have dependents and a partner.  This customer has been with the service for a long time and was not persuaded by new products and services such as fiber optic internet, paperless billing and electronic payments.  Instead, the customer continued satisfied with their service as usual.

Now that I have gone through all the data and formulated a Churn Customer Profile, I will perform some further analysis to determine if the differences detected in the analysis are statistically significant or mere coincidence.  Figure 12 shows tables with the values of the continuous independent variables: tenure, MonthlyCharges and TotalCharges by churn.

| tenure (in months) | | | | Monthly Charges | | | | TotalCharges | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| No | Yes | | | No | Yes | | | No | Yes |
| 37. 57 | 17.97913 | | | 61. 2651 | 74.44133 | | | 2555. 34 | 1531.796 |

Figure 12

To see if there is a significant difference in the mean tenure between customers that churn and customers that did not churn, I apply a two-sample t-test.  The results of the test are below in figure 13.

```
> t.test(tenure~Churn)

        Welch Two Sample t-test

data:  tenure by Churn
t = 34.824, df = 4048.3, p-value < 2.2e-16
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 18.48789 20.69378
sample estimates:
 mean in group No mean in group Yes
         37.56997          17.97913
```

Figure 13

The p-value indicates that there is a statistically significant difference between the mean tenure for customers that churned versus those who did not.  It also shows that the mean tenure for those who churned, 17.98, is not even in the 95 percent confidence interval, 18.49 to 20.70.

```
> t.test(MonthlyCharges~Churn)

        Welch Two Sample t-test

data:  MonthlyCharges by Churn
t = -18.408, df = 4135.8, p-value < 2.2e-16
alternative hypothesis: true difference in means
95 percent confidence interval:
 -14.57957 -11.77284
sample estimates:
 mean in group No mean in group Yes
         61.26512          74.44133
```
```
> t.test(TotalCharges~Churn)

        Welch Two Sample t-test

data:  TotalCharges by Churn
t = 18.801, df = 4042.9, p-value < 2.2e-16
alternative hypothesis: true difference in means
95 percent confidence interval:
  916.8121 1130.2840
sample estimates:
 mean in group No mean in group Yes
         2555.344          1531.796
```

Figure 14

Figure 14 shows that like tenure, MonthlyCharges and TotalCharges both have statistically significant means for customers that churn versus customer that did not churn.  The fact that these variables are statistically significantly different between the populations of customer that churn and those that did not churn may indicate that these variables help to explain the possibility that a customer will churn.

After analyzing the data to this point I can gather up a number of factor I think help determine whether a customer will churn. I have the three continuous variables: tenure, MonthlyCharges and TotalCharges as well as the categorical variables that caught my attention: SeniorCitizen, Dependents, Partner, Fiber optics Internet Service, lack of security features, Contracts, PaperlessBilling and PaymentMethod.

Next, I will create a logistic regression model to identify the can statistically significantly help to explain the probability that a customer will churn. I will see if the variables I have so far identified as important are indeed statistically significant. The model will help me to more accurately construct a profile of the customer. Figure 15 shows the generalized linear model.

```
> summary(churn.logistic.model)

Call:
glm(formula = Churn ~ gender + SeniorCitizen + Partner + Dependents +
    tenure + MultipleLines + InternetService + OnlineSecurity +
    OnlineBackup + DeviceProtection + TechSupport + StreamingTV +
    StreamingMovies + Contract + PaperlessBilling + PaymentMethod +
    MonthlyCharges + TotalCharges, family = "binomial")

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.9180  -0.6791  -0.2855   0.7282   3.4300

Coefficients: (6 not defined because of singularities)
                                      Estimate Std. Error z value Pr(>|z|)
(Intercept)                          1.337e+00  1.439e+00   0.929  0.35276
genderMale                          -2.183e-02  6.480e-02  -0.337  0.73619
SeniorCitizenYes                     2.168e-01  8.453e-02   2.564  0.01033 *
PartnerYes                          -3.840e-04  7.783e-02  -0.005  0.99606
DependentsYes                       -1.485e-01  8.973e-02  -1.655  0.09796 .
tenure                              -6.059e-02  6.236e-03  -9.716  < 2e-16 ***
MultipleLinesNo phone service       -1.715e-01  6.487e-01  -0.264  0.79153
MultipleLinesYes                     4.484e-01  1.773e-01   2.530  0.01142 *
InternetServiceFiber optic           1.747e+00  7.981e-01   2.190  0.02855 *
InternetServiceNo                   -1.786e+00  8.073e-01  -2.213  0.02691 *
OnlineSecurityNo internet service          NA         NA      NA       NA
OnlineSecurityYes                   -2.054e-01  1.787e-01  -1.150  0.25031
OnlineBackupNo internet service            NA         NA      NA       NA
OnlineBackupYes                      2.604e-02  1.754e-01   0.148  0.88197
DeviceProtectionNo internet service        NA         NA      NA       NA
DeviceProtectionYes                  1.474e-01  1.764e-01   0.836  0.40339
TechSupportNo internet service             NA         NA      NA       NA
TechSupportYes                      -1.805e-01  1.806e-01  -0.999  0.31759
StreamingTVNo internet service             NA         NA      NA       NA
StreamingTVYes                       5.905e-01  3.263e-01   1.810  0.07035 .
StreamingMoviesNo internet service         NA         NA      NA       NA
StreamingMoviesYes                   5.993e-01  3.267e-01   1.834  0.06658 .
ContractOne year                    -6.608e-01  1.076e-01  -6.142 8.15e-10 ***
ContractTwo year                    -1.357e+00  1.764e-01  -7.691 1.46e-14 ***
PaperlessBillingYes                  3.424e-01  7.450e-02   4.596 4.31e-06 ***
PaymentMethodCredit card (automatic) -8.779e-02  1.141e-01  -0.770  0.44156
PaymentMethodElectronic check        3.045e-01  9.450e-02   3.222  0.00127 **
PaymentMethodMailed check           -5.759e-02  1.149e-01  -0.501  0.61627
MonthlyCharges                      -4.034e-02  3.176e-02  -1.270  0.20392
TotalCharges                         3.289e-04  7.063e-05   4.657 3.20e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 8143.4  on 7031  degrees of freedom
Residual deviance: 5826.3  on 7008  degrees of freedom
  (11 observations deleted due to missingness)
AIC: 5874.3
Number of Fisher Scoring iterations: 6
```
Figure 15

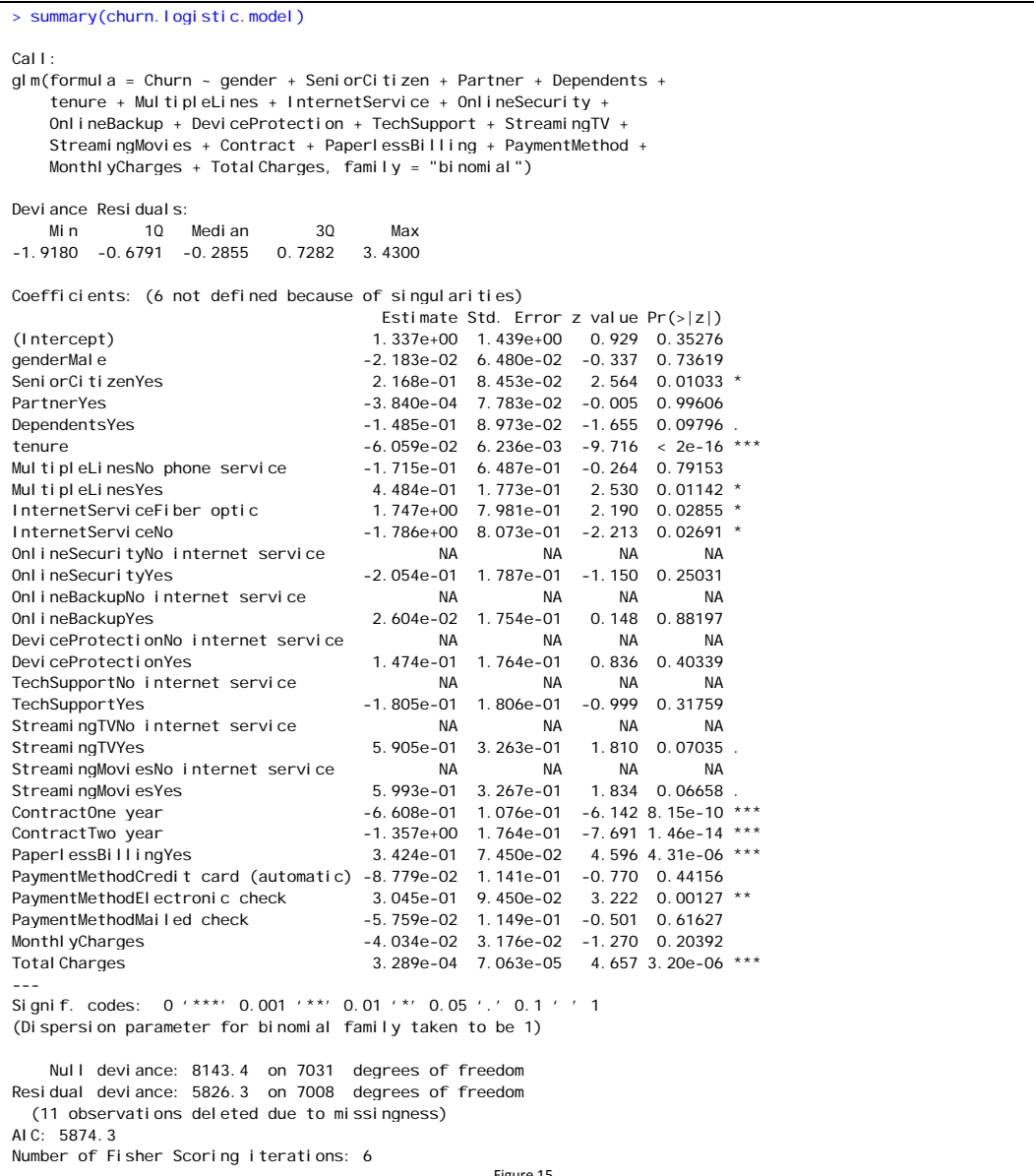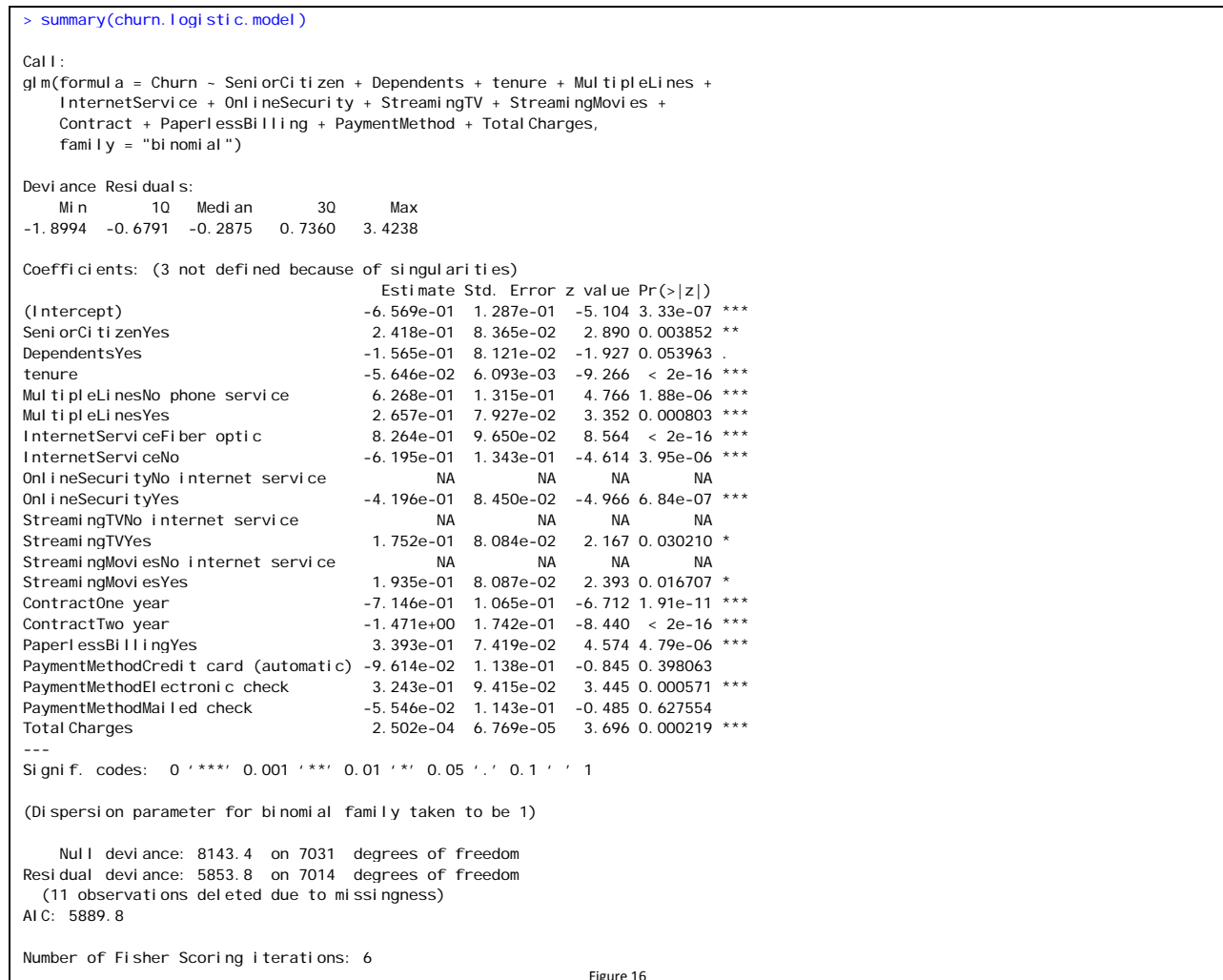The first iteration of the model which included all independent variable did the job of identifying which variable are statistically significant when explaining the possibilities a customer will churn by the one that have a p-value less than 0.1. These variables are: SeniorCitizen, Dependents, tenure, MultipleLines, InternetService, OnlineSecurity, StreamingTV, StreamingMovies, Contract, PaperlessBilling,

PaymentMethod and TotalCharges.  Another iteration of the generalized linear regression model but this time with only these variables yields the results in table 16.

```
> summary(churn.logistic.model)

Call:
glm(formula = Churn ~ SeniorCitizen + Dependents + tenure + MultipleLines +
    InternetService + OnlineSecurity + StreamingTV + StreamingMovies +
    Contract + PaperlessBilling + PaymentMethod + TotalCharges,
    family = "binomial")

Deviance Residuals:
    Min      1Q   Median      3Q      Max
-1.8994  -0.6791  -0.2875   0.7360   3.4238

Coefficients: (3 not defined because of singularities)
                                      Estimate Std. Error z value Pr(>|z|)
(Intercept)                         -6.569e-01  1.287e-01  -5.104 3.33e-07 ***
SeniorCitizenYes                     2.418e-01  8.365e-02   2.890 0.003852 **
DependentsYes                       -1.565e-01  8.121e-02  -1.927 0.053963 .
tenure                              -5.646e-02  6.093e-03  -9.266  < 2e-16 ***
MultipleLinesNo phone service        6.268e-01  1.315e-01   4.766 1.88e-06 ***
MultipleLinesYes                     2.657e-01  7.927e-02   3.352 0.000803 ***
InternetServiceFiber optic           8.264e-01  9.650e-02   8.564  < 2e-16 ***
InternetServiceNo                   -6.195e-01  1.343e-01  -4.614 3.95e-06 ***
OnlineSecurityNo internet service          NA         NA      NA       NA
OnlineSecurityYes                   -4.196e-01  8.450e-02  -4.966 6.84e-07 ***
StreamingTVNo internet service             NA         NA      NA       NA
StreamingTVYes                       1.752e-01  8.084e-02   2.167 0.030210 *
StreamingMoviesNo internet service         NA         NA      NA       NA
StreamingMoviesYes                   1.935e-01  8.087e-02   2.393 0.016707 *
ContractOne year                    -7.146e-01  1.065e-01  -6.712 1.91e-11 ***
ContractTwo year                    -1.471e+00  1.742e-01  -8.440  < 2e-16 ***
PaperlessBillingYes                  3.393e-01  7.419e-02   4.574 4.79e-06 ***
PaymentMethodCredit card (automatic) -9.614e-02  1.138e-01  -0.845 0.398063
PaymentMethodElectronic check        3.243e-01  9.415e-02   3.445 0.000571 ***
PaymentMethodMailed check           -5.546e-02  1.143e-01  -0.485 0.627554
TotalCharges                         2.502e-04  6.769e-05   3.696 0.000219 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 8143.4  on 7031  degrees of freedom
Residual deviance: 5853.8  on 7014  degrees of freedom
  (11 observations deleted due to missingness)
AIC: 5889.8

Number of Fisher Scoring iterations: 6
```
<div align="center">Figure 16</div>

With the information on figure 16 I can create a model to calculate the probability that a customer will churn.  The model is the following, where P is the probability that a customer will churn:

$$P = \frac{\exp(\text{Intercept} + \text{SeniorCitizenYes}(x_1) + \text{DependentsYes}(x_2) + \text{tenure}(x_3) + \text{MultipleLinesNoPhoneService}(x_4) + \text{MultipleLinesYes}(x_5) + \text{InternetServiceFiberoptic}(x_6) + \text{InternetServiceNo}(x_7) + \text{OnlineSecurityYes}(x_8) + \text{StreamTVYes}(x_9) + \text{StreamMoviesYes}(x_{10}) + \text{ContractOneyear}(x_{11}) + \text{Contracttwoyear}(x_{12}) + \text{PaperlessBillingYes}(x_{13}) + \text{PaymentMethodElectronicCheck}(x_{14}) + \text{TotalCharges}(x_{15}))}{1 + \exp(\text{Intercept} + \text{SeniorCitizenYes}(x_1) + + \text{DependentsYes}(x_2) + \text{tenure}(x_3) + \text{MultipleLinesNoPhoneService}(x_4) + \text{MultipleLinesYes}(x_5) + \text{InternetServiceFiberoptic}(x_6) + \text{InternetServiceNo}(x_7) + \text{OnlineSecurityYes}(x_8) + \text{StreamTVYes}(x_9) + \text{StreamMoviesYes}(x_{10}) + \text{ContractOneyear}(x_{11}) + \text{Contracttwoyear}(x_{12}) + \text{PaperlessBillingYes}(x_{13}) + \text{PaymentMethodElectronicCheck}(x_{14}) + \text{TotalCharges}(x_{15}))}$$

I substitute the variable name with the corresponding Estimate value from figure 16.

$$P = \frac{\exp(-0.6569 + 0.2418(x_1) + -0.1565(x_2) + -0.05646(x_3) + 0.6268(x_4) + 0.2657(x_5) + 0.8264(x_6) + -0.6195(x_7) + -0.4196(x_8) + 0.1752(x_9) + 0.1935(x_{10}) + -0.7146(x_{11}) + -1.471(x_{12}) + 0.3393(x_{13}) + 0.3243(x_{14}) + 0.0002502(x_{15}))}{1 + \exp(-0.6569 + 0.2418(x_1) + -0.1565(x_2) + -0.05646(x_3) + 0.6268(x_4) + 0.2657(x_5) + 0.8264(x_6) + -0.6195(x_7) + -0.4196(x_8) + 0.1752(x_9) + 0.1935(x_{10}) + -0.7146(x_{11}) + -1.471(x_{12}) + 0.3393(x_{13}) + 0.3243(x_{14}) + 0.0002502(x_{15}))}$$

Now, I take a random customer and plug in the values for the variables in the model.

```
> churn.data[sample(nrow(churn.data), 1),]
    customerID gender SeniorCitizen Partner Dependents tenure PhoneService MultipleLines
762 1894-IGFSG Female          No      No         No     22          Yes            No
    InternetService OnlineSecurity OnlineBackup DeviceProtection TechSupport StreamingTV
762     Fiber optic             No           No               No          No         Yes
    StreamingMovies     Contract PaperlessBilling    PaymentMethod MonthlyCharges TotalCharges Churn
762             Yes Month-to-month              No Electronic check          89.25      1907.85   Yes
```

The random sample is of row 762 where we can see that the customer Churned.  Let's see what my model tells me.

$$P = \frac{\exp(-0.6569 + 0.2418(0) + -0.1565(0) + -0.05646(22) + 0.6268(0) + 0.2657(0) + 0.8264(1) + -0.6195(0) + -0.4196(0) + 0.1752(1) + 0.1935(1) + -0.7146(0) + -1.471(0) + 0.3393(0) + 0.3243(1) + 0.0002502(1907.85))}{1 + \exp(-0.6569 + 0.2418(0) + -0.1565(0) + -0.05646(22) + 0.6268(0) + 0.2657(0) + 0.8264(1) + -0.6195(0) + -0.4196(0) + 0.1752(1) + 0.1935(1) + -0.7146(0) + -1.471(0) + 0.3393(0) + 0.3243(1) + 0.0002502(1907.85))}$$

Simplified:

$$P = \frac{\exp(-0.6569 -0.05646(22) + 0.8264(1) + 0.1752(1) + 0.1935(1) + 0.3243(1) + 0.0002502(1907.85))}{1 + \exp(-0.6569 -0.05646(22) + 0.8264(1) + 0.1752(1) + 0.1935(1) + 0.3243(1) + 0.0002502(1907.85))}$$

Solved:

$$P = \frac{\exp(-0.6569 - 1.24212 + 0.8264 + 0.1752 + 0.1935 + 0.3243 + 0.47734407)}{1 + \exp(-0.6569 - 1.24212 + 0.8264 + 0.1752 + 0.1935 + 0.3243 + 0.47734407)}$$

$$P = \frac{\exp(0.09772407)}{1 + \exp(0.09772407)}$$

*P = 1.10265848657/(1 + 1.10265848657) = 1.10265848657/(2.10265848657)*
*P = 52.44%*

The model tells me that there is a 52.44% probability that the customer will churn and as we can see in the value for Churn that the customer did indeed churn.  The model has accurately predicted the outcome.

In conclusion, I used different methods of identifying variables that may help me explain the probability that a customer will churn.  With the descriptive statistics I was able to identify some factor I believe contribute to the probability a customer would churn.  These factors were: tenure, MonthlyCharges and TotalCharges as well as the categorical variables that caught my attention: SeniorCitizen, Dependents, Partner, Fiber optics Internet Service, security features, Contracts, PaperlessBilling and PaymentMethod.  Then I used a generalized linear model to identify variable that can statistically significantly explain the probability that a customer would churn.  The model identified many of the same variables I identified while describing the data.  The variables identified by the regression model were: SeniorCitizen, Dependents, tenure, MultipleLines, InternetService, OnlineSecurity, StreamingTV, StreamingMovies, Contract, PaperlessBilling, PaymentMethod and TotalCharges.  The variables coded in green were identified by both the descriptive statistics and the logistic regression model.  I will use these factors to formulate the final version of the Churn Customer Profile.  The customer likely to churn is a customer that prefers month-to-month contracts with no dependents that seek the latest in internet service technology and prefers electronic forms of billing and payment but in not concern with online security services such as device protection, tech support or online backups.  In contrast, the customer less likely to churn have long tenure because they prefer one or two year contracts, have dependents and a partner and have more traditional forms of internet service, billing and payment.  It appears clear to me that the key to success for a telecom subscriber business such as the one this data was taken from is to get customers to sign one or two year contracts.  These long term contracts will expand the length of the tenure increasing not only TotalCharges but also the likelihood that they will not churn.  Even

customer that may have churned due to dissatisfaction with new services, billing and service methods will have had to stay for a considerable amount of time possibly giving the business an opportunity to address those customers by approaching alternative that may work better for them, except for, of course, the ones (Senior Citizens) that will soon perish.  No marketing strategy can retain those folks.