Assignment 2: Logistic Regression Using SAS Enterprise Miner

Daanish Ahmed

DATA 640 9041

Fall 2017

dsahmed2334@yahoo.com

Professor Sounak Chakraborty

UMUC

October 16, 2017

**Introduction**

Logistic regression is a valuable predictive model for analyzing a categorical target. It involves using a logarithmic function with values between 0 and 1 to estimate the probability of the dependent variable having an expected class (Knode, 2016a). These models function the best on binary targets, and they often perform better when using variable selection methods such as forward selection, backward elimination, or stepwise selection. In this assignment, I will build three logistic regression models using SAS Enterprise Miner for a dataset on customer churn from a cellular phone company (see Figure 1 in Appendix). My goal is to predict whether a given customer will churn so that the company can identify potential churners and prevent them from leaving. I will create one model for forward selection, one for backward elimination, and one for stepwise selection. I will use appropriate variable transformations for each model to minimize their misclassification rates. And since we are most interested in cases where customers are likely to churn, I will use appropriate cutoff thresholds to balance the number of identified churners with the overall accuracy. I will then compare these models to determine which one will maximize customer retention and minimize financial losses for the company.

The dataset contains 15 columns and 2085 observations. The target variable is churn, which is a binary non-numeric (character) variable. Churn has two possible classes: no (meaning the customer did not churn) and yes (meaning the customer churned). Some of the inputs include account length, customer service calls, day calls, day minutes, evening calls, evening minutes, international plan, international calls, and international minutes (see Figure 2 in Appendix). The customer ID variable (labeled as "VAR1") was set to rejected since it is not useful for the analysis. This leaves us with 14 variables remaining. Of the remaining variables, three are non-numeric: churn, international plan, and voice mail plan. All three of these non-numeric variables are binary.

The remaining inputs are all numeric interval variables. According to the descriptive statistics (see Figure 3 in Appendix), 0.047962% of values are missing for every single variable in this dataset. Since there are 2085 records, this means that there is one missing value in each variable (including the ID). And by examining the data in Figure 2, I found that the last row in the dataset is completely empty. We can therefore remove this row during data preparation. Regarding skewness, we see that the variables with the highest skewness are voice mail messages (1.31), international calls (1.19), and customer service calls (1.15). I will later examine the distributions of each interval input and apply appropriate transformations for the models I will build.

## Data Preparation

I followed the SEMMA process during the implementation of my regression models (see Figure 1 in Appendix). I used nodes for sample (Data Partition, Filter), explore (StatExplore), modify (Replacement, Transform Variables), model (Regression), and assess (Cutoff, Model Comparison). My dataset was originally in .xlsx format, so I saved it as a .csv file in Excel and then imported it into SAS Enterprise Miner using the File Import node. Finally, I converted the file into a .sas7bdat file using the Save Data node. I began data exploration by using the StatExplore node to produce a bar plot of each input variable's worth (see Figure 4 in Appendix). According to the image, the four most significant variables are day minutes, customer service calls, international plan, and evening minutes. Likewise, the least important variables are account length, night calls, evening calls, voice mail plan, night minutes, international calls, and international minutes. During the model creation stage, I will experiment with removing unimportant variables from some of my models. This is to ensure that each of my models has the lowest possible misclassification rate.

By examining the prior probabilities for the output (churn), I found that this variable is skewed such that only 14.2% of customers will churn. To address this issue, I used decision weights to place a cost on incorrect predictions (see Figure 5 in Appendix). In this figure, decision 1 involves the company intervening to prevent a customer from churning, while decision 2 involves the company doing nothing. I set the cost of company intervention to -4.0, the cost of losing a customer to -10.0, and the profit from keeping a customer to 10.0. If a customer is expected to churn and the company appeals to him, then they will make 6.0. But if they appeal to a non-churner, they will waste 4.0. If they do nothing and a customer churns, they lose 10.0. And if they do nothing and the customer doesn't churn, they will gain a 10.0 profit. These weights will add insight to the costs of failing to make accurate predictions.

After this step, I proceeded to handle missing values, outliers, and incorrect entries. Since one of the rows contains every missing value in the dataset, we can remove that row using the Filter node. I set the property "keep missing values" to "no," thus removing all observations with missing values. For both class and interval variables, I set the default filtering method to "none" to prevent the node from filtering records with outliers. Instead, I chose to replace outliers by using the Replacement node. By examining the descriptive statistics, I found that variables such as day minutes and evening minutes have values more than 3 standard deviations from the mean (see Figure 3 in Appendix). Thus, I set the interval variable default limits method to "standard deviations from the mean." By using the default cutoff value of 3.0, this node will compute a replacement for any values more than 3 standard deviations from the mean. Finally, by using the replacement editor, we can verify that there are no incorrect entries in this dataset.

Next, I used the Data Partition node to split the dataset into training, validation, and test sets. I allocated 60% of data to the training set, 30% to the validation set, and 10% to the test set.

Finally, I transformed the variables to ensure they function effectively in each model. I used two types of transformation methods: "best" transform and maximum normal. According to SAS Enterprise Miner Reference Help, the "best" method involves applying several types of transformations to each input and selecting the transformation with the best Chi-squared test results. The Chi-squared test is useful for finding the significance of a relationship between two variables (Ray, 2016), which in this case refers to the relationship between the output (churn) and the inputs. The maximum normal method, as described by SAS Enterprise Miner Reference Help, also performs multiple transformations on each variable, but instead it uses the transformation that maximizes the normality of each variable's distribution. This is useful for reducing the skewness of some inputs such as customer service calls and international calls—both of which have positive skews (see Figure 6 in Appendix). I experimented with matching different models with each of these transformation types and evaluated their misclassification rates. Through this process, I found that the "best" method works better with my forward selection model, while "max normal" is best for the backward and stepwise models. To eliminate non-numeric inputs, as well as handle variables like "voice mail messages" with spiked distributions (see Figure 6 in Appendix), I set the class output for every Transform Variables node to "dummy indicators."

**Model Development**

Once data preparation was finished, I began to implement the logistic regression models. I created one model for forward selection, one for backward elimination, and one for stepwise selection. According to Knode (2016a), forward selection begins with an empty model and adds inputs one at a time, while backward elimination starts with all variables and removes them from the model. Stepwise selection begins similarly as forward selection, but allows for inputs to be

4

added or removed (Knode, 2016a). I built these models using the Regression node by setting the

regression type to "logistic regression" and choosing the selection models "forward," "backward,"

and "stepwise" for the appropriate models. For all three models, I set the selection criterion to

"validation misclassification" to minimize the number of misclassified observations. Before

running these nodes, I experimented with manually filtering out the least significant inputs from

each model. During data exploration, I found that the five least important variables are account

length, night calls, evening calls, voice mail plan, and night minutes (see Figure 4 in Appendix).

These variables were excluded by editing the variables under the Regression node and setting their

usage to "no." Through this process, I found that the forward regression model was the only model

that saw improvement in its validation misclassification and average squared error. Therefore, this

is the only model for which I will manually exclude insignificant inputs.

After building these models, I examined the graphs of the regression coefficients, which

can be used to build the regression equation. The coefficients are also useful for determining how

strongly an input is related to the output and whether the relationship is positive or negative (Frost,

2013). By examining the graph for my forward selection model (see Figure 7 in Appendix), we

see that the most significant inputs are day minutes (three bins with values -4.69, -3.29, and -1.44),

customer service calls (-2.23), and international plan (1.12). For the backward elimination model

(see Figure 12 in Appendix), the important variables are customer service calls (2.45), international

calls (-2.29), and international plan (0.97). Finally, the stepwise selection results (see Figure 16

in Appendix) indicate that the top inputs are customer service calls (2.37), international calls (-

2.16), and international plan (0.95). The backward and stepwise models rely on many of the same

variables with similar coefficients, yet the forward regression model does not. This is likely due

to the backward and stepwise models using the same transformation method of "max normal,"

while the forward model uses the "best" transformation method. Using different transformations can cause the variables to be completely different, and thus the selection methods would have chosen different variables for each of these transformations.

Finally, I will change the cutoff thresholds to maximize each model's predictive performance. Cutoff thresholds are used to decide whether a prediction qualifies as a "1" or a "0," and they can be adjusted to maximize either the overall accuracy or the sensitivity—the number of true positives (Knode, 2016b). Since this dataset's output is heavily skewed, having a higher cutoff threshold will increase the accuracy but lower the sensitivity. Since we are interested in identifying possible churners, we want our models to have high sensitivity while still having acceptable accuracy. By looking at the classification rates for the forward selection model (see Figure 11 in Appendix), we see that 0.1 is a good cutoff threshold because it satisfies the mentioned conditions. I will use 0.1 as the cutoff value for all three models to allow for easier comparison between the models. Finally, I will connect an additional Cutoff node to the stepwise model using a threshold of 0.25 to study the impact of maximizing accuracy over sensitivity.

## Results

I will now describe the results of my analysis to see which model is the best for predicting churn. First, I will examine the iteration plots for my three models. These plots show the misclassification rate for the training and validation sets during each step of the model building process. The plot for my forward regression model reveals that the model required five steps to train, and that the misclassification rate is higher for the validation set (see Figure 9 in Appendix). By moving the cursor over the graph in SAS Enterprise Miner, I found that the training and

validation misclassification rates are 9.76% and 11.86% respectively (see Figure 21 in Appendix). This difference is not large enough to suggest that overfitting took place, especially since the test set's misclassification of 10.95% is lower than that of the validation set. The backward regression model also required five steps to train, but this time the training misclassification rate is higher than the validation rate (see Figure 14 in Appendix). Figure 21 shows that the misclassification for training, validation, and test data are 12.96%, 12.02%, and 14.76%. The stepwise model also has a higher misclassification for training data, but it is also the most complex model since it required six steps to train (see Figure 18 in Appendix). According to Figure 21, the misclassification for training, validation, and test data are 13.36%, 11.86%, and 13.33%. Since the test set is by far the smallest of the three samples with only 10% of the data, I will focus mostly on the validation results. Based on this, the forward and stepwise models appear to be the most accurate models since they have the lowest validation misclassification rates. But we must first examine the classification tables to see which model is better for predicting churn itself.

Classification tables provide us with the number of predicted "yes" and "no" instances compared to the actual number of "yes" and "no" cases for both the training and validation sets. This helps us to find the number of true positives, as well as the number of false positives and negatives. By examining the classification table for the forward selection model (see Figure 8 in Appendix), we see that the model correctly identified 37 out of 88 customers who will churn. This means that only 42% of churning customers were identified. The 51 churners who were not identified can cause notable losses for the company. We also see that 23 non-churners were identified as churners, which may lead to unnecessary costs since these customers do not need to be appealed to. By comparing the validation results to the training results, it is evident that the model is less accurate on validation data. But as mentioned earlier, the difference in accuracy does

not appear to be high enough to suggest overfitting took place. Overall, we see that the model has an accuracy of 87.98%—but this is mainly due to correct identification of non-churners.

By looking at the classification table for the backward regression model (see Figure 13 in Appendix), we find that the model identified 24 out of 88 true positive cases in the validation set. It has a sensitivity of 27%, meaning that it is much weaker at identifying true positives compared to my first model. This model identified 11 non-churners as churners, which is lower than the 23 false positives of my previous model. However, this model failed to recognize a higher number of churners than the previous model (64 compared to 51). Failure to identify churners is much costlier than failure to identify non-churners (see Figure 5 in Appendix), and thus this model will lead to a higher cost for the company. The classification table for the stepwise model (see Figure 17 in Appendix) has similar results to that of the backward model. This model identified 22 true positives, has a sensitivity of 25%, and is altogether weaker than the forward model. However, the backward and stepwise models performed better on the validation data than on the training data, unlike the forward model. These two models also have accuracies similar to that of the first model, but this is once again due to the high number of true negatives identified.

To increase the number of true positives identified, I evaluated each model using a cutoff threshold of 0.1. By examining the cutoff statistics for the forward selection model (see Figure 10 in Appendix), we see that the model's sensitivity has been significantly improved. The model has successfully identified 163 churners in the training set, 74 in the validation set, and 27 in the test set. This leads to a true positive rate of 92.1% for training data, 84.1% for validation data, and 90% for test data. The only downside to this approach is that the overall classification accuracy decreases. However, the forward model's accuracy rates for the training, validation, and test sets are 83.0%, 80.8%, and 83.3%. These rates are only slightly lower than the model's original

accuracy rates. When viewing the cutoff statistics for the backward model, we find the true positive rate to be 81.9% for training data, 78.4% for validation data, and 83.3% for test data (see Figure 15 in Appendix). Though the model now has a higher sensitivity, it is still weaker than the forward regression model. Furthermore, the overall classification rate for validation data is 63.8%, which is significantly worse than the model's original accuracy rate.

The results of the stepwise model are similar to those of the backward model. Its true positive rates for training, validation, and test data are 80.2%, 77.3%, and 86.7% (see Figure 19 in Appendix). Likewise, its overall accuracy is 62.5% on the validation data. By comparing my three models, it is apparent that the forward regression model is by far the strongest model with regards to both accuracy and for identifying churners. Finally, I evaluated the stepwise model using a cutoff threshold of 0.25 to determine if it is better to have a higher overall accuracy. This model now has a classification rate of 82.1% for the validation set (see Figure 20 in Appendix), which is significantly higher than its accuracy when using a cutoff of 0.1. However, its true positive rates decrease to 49.7%, 50%, and 53.3% for the training, validation, and test sets. Therefore, I would recommend using a cutoff threshold of 0.1 to maximize the number of identified churners and minimize financial losses.

**Conclusion**

In my analysis, I developed three logistic regression models to estimate the likelihood of a customer churning from the company. I found that the results differed the most when I used a different transformation type, not when I used a different variable selection method. For instance, the backward and stepwise methods were both impacted most by the variables "customer service

calls," "international calls," and "international plan." In both cases, the variables had similar coefficients and only "international calls" had a negative relationship with churn. However, the forward selection model relied most heavily on "day minutes," "customer service calls," and "international plan"—where only "international plan" had a positive relationship with the output. When analyzing the accuracy and sensitivity, I found that the backward and stepwise models had similar results to each other but differed drastically from those of the forward model. Overall, I found that the forward regression model had the highest number of true positives, and its accuracy was the highest of my three models when setting the cutoff threshold to 0.1. When considering the costs from false predictions, the forward model has the lowest number of false negatives—thus producing the lowest expenses for the company. Therefore, this model is easily the most effective model for maximizing customer retention and minimizing costs.

One of the shortcomings of my analysis is the question of whether my forward regression model is comparable with my other two models. This model not only used a completely different type of transformation from the other models, but it was also the only model in which I manually excluded unimportant input variables prior to model development. Due to these differences, the question remains whether forward selection is truly better than backward or stepwise selection. One suggestion for future improvement would be to implement backward and stepwise models using the "best" transformation type while following similar steps to those I used to build the forward regression model. If one of these new models is more effective at predicting customer churn, then it suggests that using the "best" transformation instead of "max normal" was the primary reason why my forward regression model was so effective. I would recommend further analysis because even though my forward selection model is the most accurate of my models, there is still room for improvement even when using a cutoff threshold of 0.1.

# References

Frost, J. (2013, December 12). Regression Analysis Tutorial and Examples. Retrieved October 4, 2017, from http://blog.minitab.com/blog/adventures-in-statistics/regression-analysis-tutorial-and-examples

Knode, S. (2016a, August 19). *Regression Models*. Lecture presented at UMUC. Retrieved October 3, 2017.

Knode, S. (2016b, October 11). *Adjusting for Skewed Target Distribution*. Lecture presented at UMUC. Retrieved October 14, 2017.

Ray, S. (2016, January 10). A Comprehensive Guide to Data Exploration. Retrieved September 26, 2017, from https://www.analyticsvidhya.com/blog/2016/01/guide-data-exploration/

Welcome to Logistic Regression Analysis. (n.d.). Retrieved October 6, 2017, from http://logisticregressionanalysis.com/

# Appendix

Relevant SAS Enterprise Miner Output Images

*Figure 1.* Process Flow Diagram for Customer Churn Logistic Regression Models.

*Figure 2.* Customer Churn at a Cellular Phone Company Dataset.

| Name | Role | Level | Type | Number of Levels | Percent Missing | Minimum | Maximum | Mean | Standard Deviation | Skewness | Kurtosis |
|------|------|-------|------|------------------|-----------------|---------|---------|------|--------------------|----------|----------|
| Account_Length | Input | Interval | Numeric | . | 0.047962 | 1 | 232 | 100.8714 | 40.90292 | 0.068997 | -0.13497 |
| Churn | Target | Binary | Character | 2 | 0.047962 | . | . | . | . | . | . |
| Customer_Service_Calls | Input | Interval | Numeric | . | 0.047962 | 0 | 9 | 1.571017 | 1.327448 | 1.152134 | 2.047751 |
| Day_Calls | Input | Interval | Numeric | . | 0.047962 | 0 | 163 | 100.2361 | 20.38864 | -0.17372 | 0.418226 |
| Day_Minutes | Input | Interval | Numeric | . | 0.047962 | 0 | 346.8 | 179.9886 | 54.56822 | -0.03005 | 0.095213 |
| Evening_Calls | Input | Interval | Numeric | . | 0.047962 | 0 | 168 | 99.8095 | 20.03431 | -0.11128 | 0.330131 |
| Evening_Minutes | Input | Interval | Numeric | . | 0.047962 | 0 | 351.6 | 201.8364 | 50.12722 | -0.05516 | -0.02615 |
| International_Calls | Input | Interval | Numeric | . | 0.047962 | 0 | 18 | 4.524952 | 2.480428 | 1.197296 | 2.220202 |
| International_Minutes | Input | Interval | Numeric | . | 0.047962 | 0 | 20 | 10.21171 | 2.797119 | -0.25924 | 0.659338 |
| International_Plan | Input | Binary | Character | 2 | 0.047962 | . | . | . | . | . | . |
| Night_Calls | Input | Interval | Numeric | . | 0.047962 | 33 | 175 | 100.2634 | 19.70518 | 0.052944 | 0.004829 |
| Night_Minutes | Input | Interval | Numeric | . | 0.047962 | 43.7 | 395 | 201.6883 | 50.08343 | 0.016151 | 0.052053 |
| VAR1 | Rejected | Interval | Numeric | . | 0.047962 | 11 | 2094 | 1052.5 | 601.7433 | 0 | -1.2 |
| Voice_Mail_Messages | Input | Interval | Numeric | . | 0.047962 | 0 | 51 | 7.797025 | 13.48652 | 1.3136 | 0.074732 |
| Voice_Mail_Plan | Input | Binary | Character | 2 | 0.047962 | . | . | . | . | . | . |

*Figure 3*. Descriptive Statistics of Customer Churn Dataset.



*Figure 4*. StatExplore Plot of Input Variable Worth.



*Figure 5*. Decision Weights for Customer Churn Predictions.

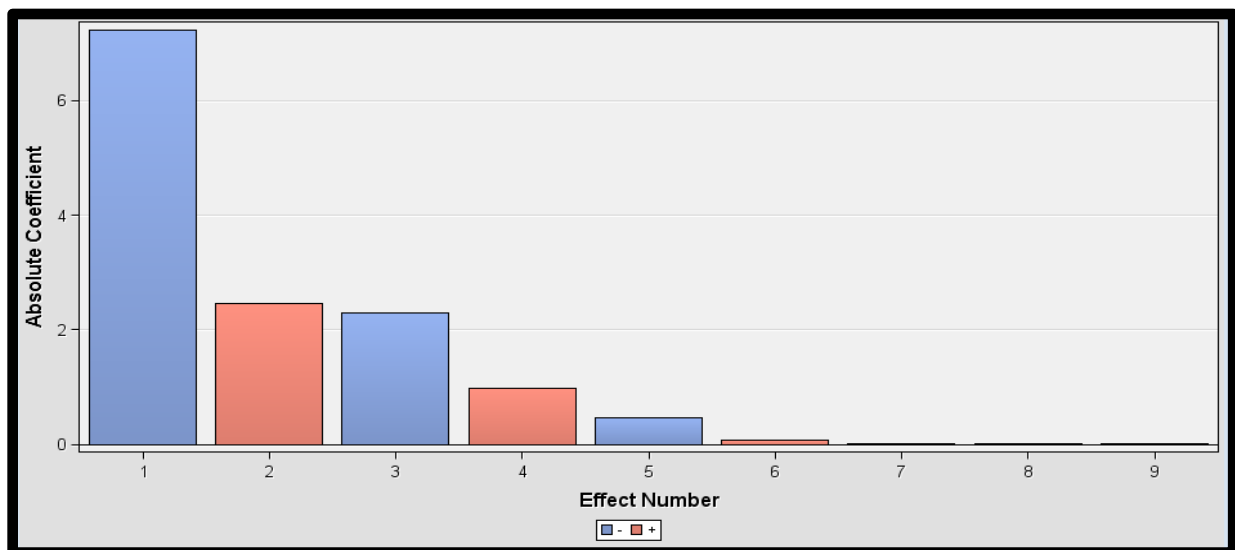*Figure 6.* Initial Distributions of Variables in Customer Churn Dataset.



*Figure 7.* Graph of Regression Coefficients for Forward Selection Model.

```
Classification Table

Data Role=TRAIN Target Variable=REP_Churn Target Label=Replacement: Churn

                        Target       Outcome      Frequency      Total
Target      Outcome     Percentage   Percentage   Count          Percentage

  NO          NO         92.4175      96.5517       1036          82.88
  YES         NO          7.5825      48.0226         85           6.80
  NO          YES        28.6822       3.4483         37           2.96
  YES         YES        71.3178      51.9774         92           7.36


Data Role=VALIDATE Target Variable=REP_Churn Target Label=Replacement: Churn

                        Target       Outcome      Frequency      Total
Target      Outcome     Percentage   Percentage   Count          Percentage

  NO          NO         90.9574      95.7090        513          82.2115
  YES         NO          9.0426      57.9545         51           8.1731
  NO          YES        38.3333       4.2910         23           3.6859
  YES         YES        61.6667      42.0455         37           5.9295
```

*Figure 8.*  Classification Table for Forward Selection Model.



*Figure 9.*  Iteration Plot for Forward Selection Model.



*Figure 10.*  Cutoff Statistics for Forward Selection Model with 0.1 Cutoff Threshold.

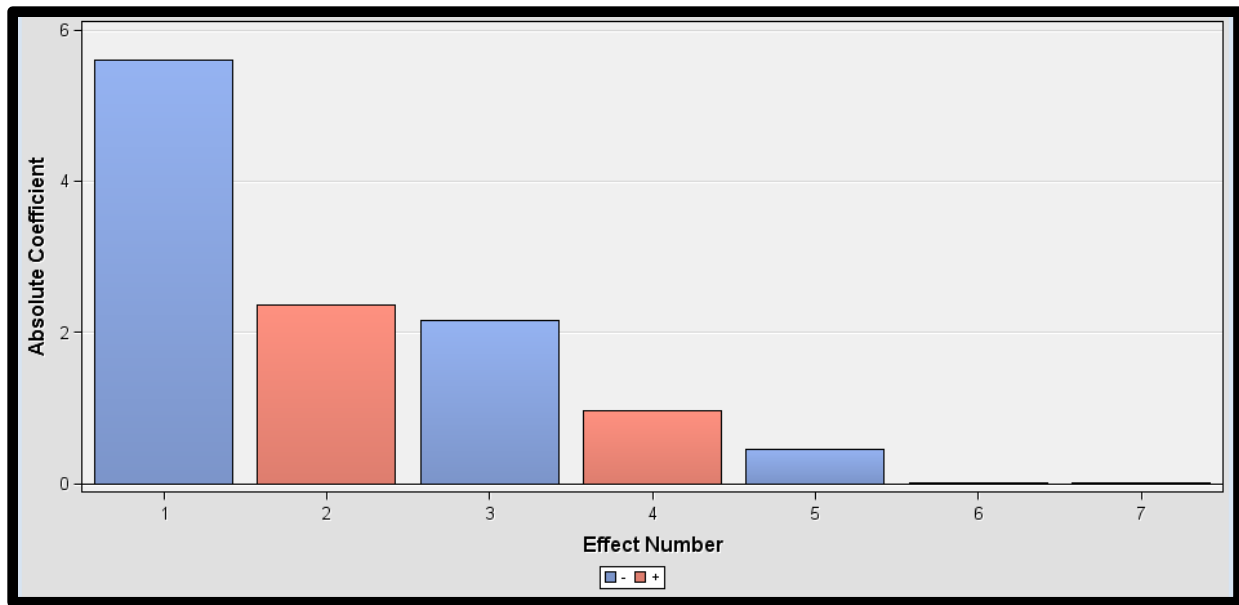*Figure 11.* Classification Rates for Forward Selection Model with 0.1 Cutoff Threshold.



*Figure 12.* Graph of Regression Coefficients for Backward Elimination Model.

```
Classification Table

Data Role=TRAIN Target Variable=REP_Churn Target Label=Replacement: Churn

                       Target        Outcome      Frequency       Total
Target     Outcome    Percentage    Percentage      Count       Percentage

 NO         NO         88.1811       98.0429         1052         84.16
 YES        NO         11.8189       79.6610          141         11.28
 NO         YES        36.8421        1.9571           21          1.68
 YES        YES        63.1579       20.3390           36          2.88


Data Role=VALIDATE Target Variable=REP_Churn Target Label=Replacement: Churn

                       Target        Outcome      Frequency       Total
Target     Outcome    Percentage    Percentage      Count       Percentage

 NO         NO         89.1341       97.9478          525         84.1346
 YES        NO         10.8659       72.7273           64         10.2564
 NO         YES        31.4286        2.0522           11          1.7628
 YES        YES        68.5714       27.2727           24          3.8462
```

*Figure 13.* Classification Table for Backward Elimination Model.
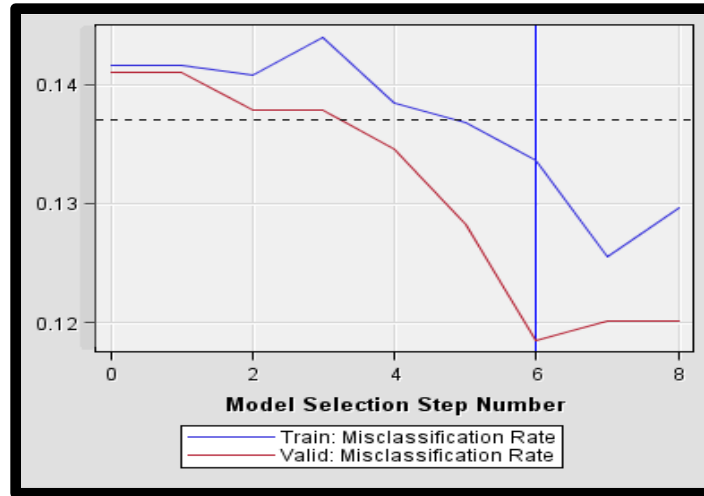


*Figure 14.* Iteration Plot for Backward Elimination Model.



| Cutoff | Counts of True Positives | Counts of False Positives | Counts of True Negatives | Counts of False Negatives | Counts of Predicted Positives | Counts of Predicted Negatives | Counts of False Positives and Negatives | Counts of True Positive and Negatives | Overall Classification Rate | Change Count True Positives | Change Count False Positives | True Positive Rate | True Negative Rate | False Positive Rate | Misscl. cost prior 0.1416 equal cost structure | Misscl. cost prior 0.1 equal cost structure | Misscl. cost prior 0.2 equal cost structure | Misscl. cost prior 0.3 equal cost structure | Misscl. cost prior 0.4 equal cost structure | Misscl. cost prior 0.5 equal cost structure | Event Precision Rate | Non Event Precision Rate | Overall Precision Rate | Data Role |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.11 | 24 | 65 | 115 | 6 | 89 | 121 | 71 | 139 | 66.19048 | 1 | 6 | 80 | 63.88889 | 36.11111 | 0.338298 | 0.345 | 0.328889 | 0.312778 | 0.296667 | 0.280556 | 26.96629 | 95.04132 | 61.00381 | TEST |
| 0.1 | 145 | 412 | 661 | 32 | 557 | 693 | 444 | 806 | 64.48 | 7 | 37 | 81.9209 | 61.60298 | 38.39702 | 0.3552 | 0.363652 | 0.343334 | 0.323016 | 0.302698 | 0.282381 | 26.03232 | 95.3824 | 60.70736 | TRAIN |
| 0.1 | 69 | 207 | 329 | 19 | 276 | 348 | 226 | 398 | 63.78205 | 0 | 21 | 78.40909 | 61.3806 | 38.6194 | 0.362082 | 0.369166 | 0.352137 | 0.335109 | 0.31808 | 0.301052 | 25 | 94.54023 | 59.77011 | VALIDATE |
| 0.1 | 25 | 70 | 110 | 5 | 95 | 115 | 75 | 135 | 64.28571 | 1 | 5 | 83.33333 | 61.11111 | 38.88889 | 0.357422 | 0.366667 | 0.344444 | 0.322222 | 0.3 | 0.277778 | 26.31579 | 95.65217 | 60.98398 | TEST |
| 0.09 | 148 | 442 | 631 | 29 | 590 | 660 | 471 | 779 | 62.32 | 3 | 30 | 83.61582 | 58.80708 | 41.19292 | 0.3768 | 0.38712 | 0.362312 | 0.337503 | 0.312694 | 0.287885 | 25.08475 | 95.60606 | 60.3454 | TRAIN |

*Figure 15.* Cutoff Statistics for Backward Elimination Model with 0.1 Cutoff Threshold.

17

*Figure 16.* Graph of Regression Coefficients for Stepwise Selection Model.



*Figure 17.* Classification Table for Stepwise Selection Model.

*Figure 18.* Iteration Plot for Stepwise Selection Model.



*Figure 19.* Cutoff Statistics for Stepwise Selection Model with 0.1 Cutoff Threshold.



*Figure 20.* Cutoff Statistics for Stepwise Selection Model with 0.25 Cutoff Threshold.

| Selected Model | Predecessor Node | Model Node | Model Description ▲ | Target Variable | Target Label | Train: Misclassifica tion Rate | Selection Criterion: Valid: Misclassifica tion Rate | Test: Misclassifica tion Rate |
|---|---|---|---|---|---|---|---|---|
| Y | CUT | Reg | Regression 1 - Forward | REP_Churn | Replacement: Churn | 0.0976 | 0.11859 | 0.109524 |
| | CUT2 | Reg2 | Regression 2 - Backward | REP_Churn | Replacement: Churn | 0.1296 | 0.120192 | 0.147619 |
| | CUT3 | Reg3 | Regression 3 - Stepwise | REP_Churn | Replacement: Churn | 0.1336 | 0.11859 | 0.133333 |

*Figure 21.* Model Comparison Based on Misclassification Rate.