

# DATA SCIENCE PRACTICUM: PREDICTING CHURN

JAMES MWAKICHAKO & MANOJ KUMAR

## CONTENTS

1	Introduction	6
1.1	Dataset . . . . .	6
2	List of Features	7
2.1	Gender . . . . .	7
2.2	Sports . . . . .	8
2.3	Search Hits . . . . .	8
2.4	Duration . . . . .	8
3	Evaluation Metrics	9
3.1	Precision . . . . .	9
3.2	Recall . . . . .	9
3.3	F1 Score . . . . .	9
4	One Month Churn Prediction Approaches	10
4.1	Using a Specific Month's Transactions Data . . . . .	10
4.2	Using Aggregate Transactions Data . . . . .	13
4.3	Using Monthly Difference in Transactions Data . . . . .	17
5	Life-time Churn Prediction	20
5.1	Dataset Attributes . . . . .	20
5.2	Churn in both males and Females . . . . .	20
5.3	Churn in Males . . . . .	23
5.4	Churn in Females . . . . .	24
6	Results and Discussion	26

## LIST OF FIGURES

Figure 1	Lostic Regression Confusion Matrix . . . . .	7
Figure 2	Lostic Regression Confusion Matrix . . . . .	8
Figure 3	Lostic Regression Confusion Matrix . . . . .	8
Figure 4	Lostic Regression Confusion Matrix . . . . .	8
Figure 5	Lostic Regression Confusion Matrix . . . . .	11
Figure 6	Lostic Regression Confusion Matrix . . . . .	11
Figure 7	Lostic Regression Confusion Matrix . . . . .	12
Figure 8	Lostic Regression Confusion Matrix . . . . .	13
Figure 9	Lostic Regression Confusion Matrix . . . . .	14
Figure 10	Lostic Regression Confusion Matrix . . . . .	14
Figure 11	Lostic Regression Confusion Matrix . . . . .	15
Figure 12	Lostic Regression Confusion Matrix . . . . .	16
Figure 13	Lostic Regression Confusion Matrix . . . . .	17
Figure 14	Lostic Regression Confusion Matrix . . . . .	18
Figure 15	Lostic Regression Confusion Matrix . . . . .	18
Figure 16	Lostic Regression Confusion Matrix . . . . .	19
Figure 17	Lostic Regression Confusion Matrix . . . . .	21
Figure 18	Lostic Regression Confusion Matrix . . . . .	21
Figure 19	Lostic Regression Confusion Matrix . . . . .	22
Figure 20	Lostic Regression Confusion Matrix . . . . .	23
Figure 21	Lostic Regression Confusion Matrix . . . . .	24
Figure 22	Lostic Regression Confusion Matrix . . . . .	25

## LIST OF TABLES

Table 1	Final List of Features . . . . .	7
Table 2	Dataset Attributes . . . . .	10
Table 3	Precision Recall Values . . . . .	10
Table 4	Precision Recall Values . . . . .	13
Table 5	Precision Recall Values . . . . .	17
Table 6	My caption . . . . .	20
Table 7	Precision Recall Values . . . . .	20
Table 8	Dataset attributes . . . . .	23
Table 9	Precision Recall Values . . . . .	24
Table 10	My caption . . . . .	24
Table 11	Precision Recall Values . . . . .	25

## ABSTRACT

Churn rate, according to the dictionary is the annual percentage rate at which customers stop subscribing to a service or employees leave a job. In the context of CaptainU, and specifically from the perspective of high school athletes, an athlete is considered to have churned if they cancel their subscription before making a college team or canceling before the spring semester of their senior year.

Therefore, the following scenarios are not considered churn:

1. When an athlete makes a team and then cancels his/her subscription
2. When an athlete cancels his/her subscription in the spring of their senior year

---

\* *Department of Data Science, Illinois Institute of Technology, Chicago, United States*

<sup>1</sup> *Department of Data Science, Illinois Institute of Technology, Chicago, United States*

## 1 INTRODUCTION

The primary goal of the practicum was to predict athletes who are most likely to churn early enough so that steps could be taken to mitigate churn. To aid in answering this question, a two pronged approach was taken.

1. Predicting lifetime churn - This method sought to predict the likelihood of a athlete churning at some point in their high school career. This was an easier approach to take and it helped us understand important features and what machine learning models to implement. The main drawback to this approach is that it doesn't have a strong business usecase. Saying ' Athlete A will churn at some point is not as actionable as the same athlete churning in the next month or two. '
2. Predicting one month churn - In this approach we sought to answer, given the monthly following transaction of athlete A, what is his/ her probability of churning in the next month ? Modeling this problem is slightly more challenging than the first but more beneficial

### 1.1 Dataset

To train and test our machine learning models, we used data provided by CaptainU. Specifically we used MSG\_RFM table. We also focused on active subscriptions. Active subscriptions refer to athletes who are paying a monthly fee to be on the system.

## 2 LIST OF FEATURES

In this section we shall briefly discuss the features we used for the machine learning models. For all the models we shall present, we used the same list of features. This list is a subset of the features in MSG\_RFM table. The histograms that accompany gives the mean and standard deviation of each feature given active subscriptions.

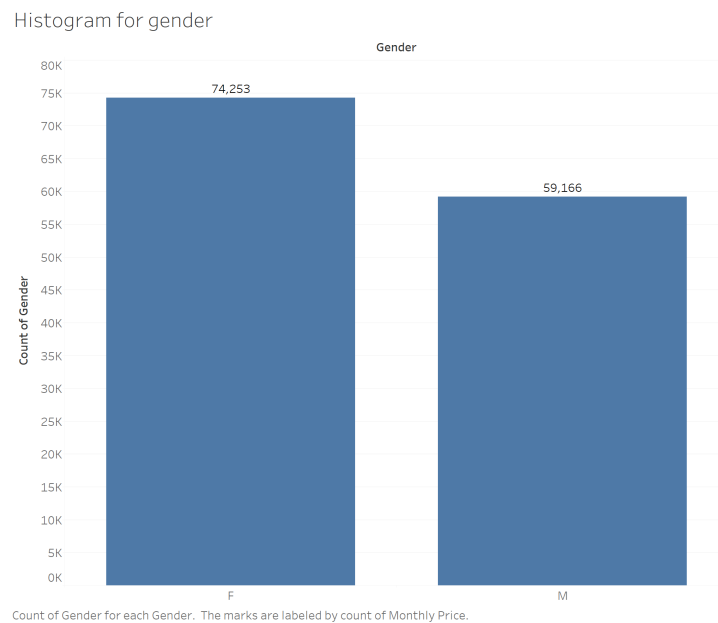
Below is a summary of the features used.

**Table 1: Final List of Features**

gender_F	gender_M
EventsAttended	Hprofileview
Hcoachimport	Hmessage
Hsearchhit	Hcoacheval
Hemailopen	EAthlete newsletter
Eathlete_new	MessagesSent
Ecoach_list_known_updated	ECoachEmailOpen
ECoachEval	MessagesSent
ECoachSearchHit	ECoachVisit
Ecolleges_going_to_the_event	Efailed_subscription
Eparent_new	Eparent_welcome
Epost_event_email	MessagesReceived
CollegeProspects	CaptainU_CHURN

### 2.1 Gender

The histogram below shows that female interacted with CaptainU more than males.



**Figure 1: Gender Histogram**





### 3 EVALUATION METRICS

These are measures we used to determine how good our machine learning models were. We used the testing data to evaluate the model fitted using the training data. We shall briefly explain the measures we used and present the formula. Before then, we shall present two variables we shall use in the proceeding equations.

- True Positive(TP) - These are the athletes who are predicted by the model to be most likely to churn and actually churned.
- True Negative(TN) - These are the athletes who are predicted by the model to be most likely to be retained and actually are retained.
- False Positive(TP) - These are the athletes who are predicted by the model to be most likely to churn and actually are retained.
- False Negative(TN) - These are the athletes who are predicted by the model to be most likely to be retained and actually churned.

#### 3.1 Precision

In the churn context, precision refers to how pure our predicted set is. Given a set of prediction of churners, how many of them are actually churners(True Positive (TP))?

$$\text{Precision} = \frac{TP}{TP + FP}$$

#### 3.2 Recall

Given a set of churners, recall refers to the fraction of churners that the model correctly returns.

$$\text{Recall} = \frac{TP}{TP + FN}$$

#### 3.3 F1 Score

This is the harmonic mean of precision and recall.

$$\text{Recall} = \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

Table 2: Dataset Attributes

Dataset	Rows	columns	Churners	Non-Churners
Full	4360	26	601	3759
Training	3488	26	484	3004
Testing	872	26	117	755

## 4 ONE MONTH CHURN PREDICTION APPROACHES

We used March 2014 to predict churn in the next month. Predicting churn in the next month entailed setting the output variable (CaptainU\_Churn) of the training and testing datasets to the output variable of the next month.

We implemented the following models in predicting churn:

1. Logistic Regression
2. Logistic Regression with (class\_weight = 'balanced') - This mode uses the values of y(CaptainU\_CHURN) to automatically adjusts weights inversely proportional to class frequencies in the test data. Effectively, more attention was paid to churners as they make up the minority class. You can read more about class\_weight [here](#)
3. Gradient Boosting
4. Random Forest Classifier

We then calculated the precision and recall values for churning.

### 4.1 Using a Specific Month's Transactions Data

This is entailed tracking an athlete's monthly transactions and using that information to predict whether the athlete would churn in the next month. In our case for instance, we used March 2014's transaction data to predict if an athlete will churn in April 2014.

#### 4.1.1 Precision Recall Table

Table 3: Precision Recall Values

Model	Precision	Recall	F1-Score
Logistic Regression	0.65	0.26	0.38
Logistic Regression(with class_weight = 'balanced')	0.43	0.56	0.49
Random Forest	0.53	0.35	0.42
Gradient Boosting	0.59	0.33	0.43

#### 4.1.2 Confusion Matrix: Logistic Regression

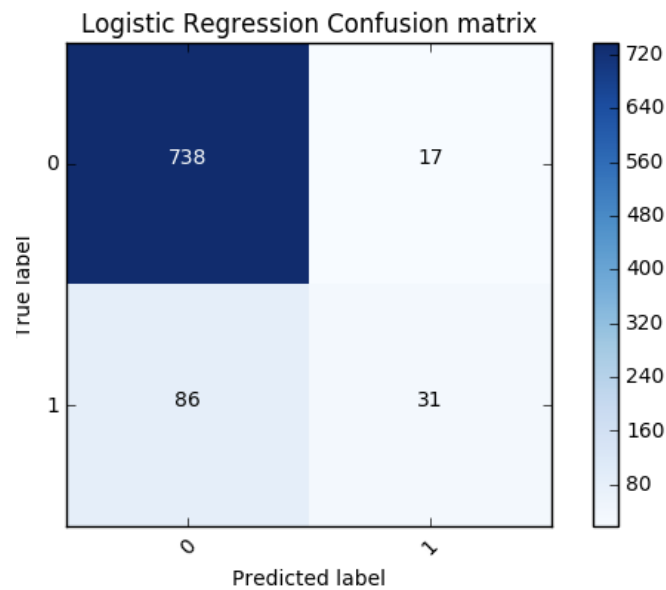


Figure 5: Lostic Regression Confusion Matrix

The above confusion matrix illustrates the following:

Given 117 athletes who are churners the model will:

- correctly predict 31 of the 117 athletes.
- incorrectly predict 17 athlete as having churned
- incorrectly predict 86 athletes as having been retained while in reality they are churners

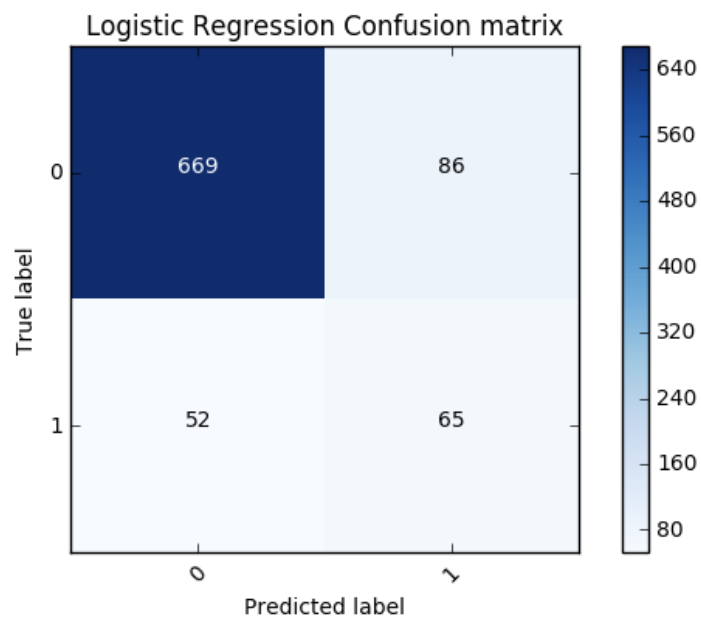


Figure 6: Lostic Regression(class\_weight = 'balanced')

The above confusion matrix illustrates the following: Given 117 athletes who are churners the model will:

- correctly predict 65 of the 117 athletes.
- incorrectly predict 86 athlete as having churned
- incorrectly predict 52 athletes as having been retained while in reality they are churners

#### 4.1.3 *Confusion Matrix: Random Forest*

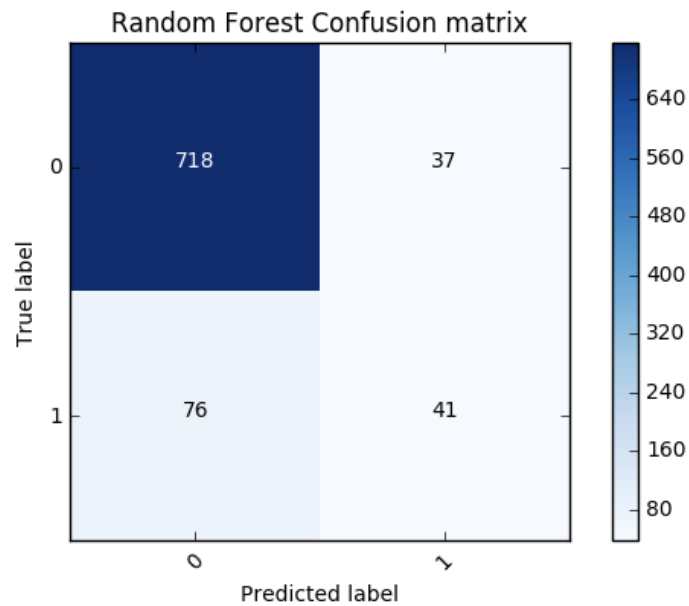


Figure 7: Random Forest Confusion Matrix

The above confusion matrix illustrates the following:  
Given 117 athletes who are churners the model will:

- correctly predict 41 of the 117 athletes.
- incorrectly predict 37 athlete as having churned
- incorrectly predict 76 athletes as having been retained while in reality they are churners

#### 4.1.4 Confusion Matrix: Gradient Boosting

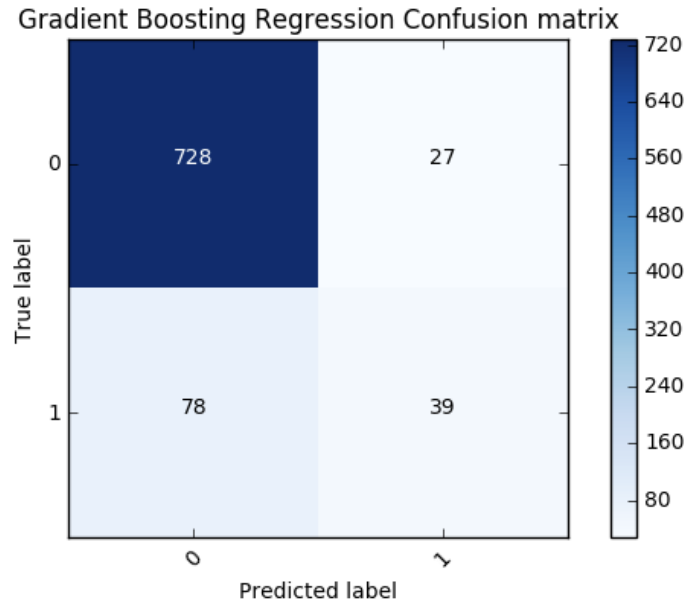


Figure 8: Gradient Boosting Confusion Matrix

The above confusion matrix illustrates the following:

Given 117 athletes who are churners the model will:

- correctly predict 39 of the 117 athletes.
- incorrectly predict 27 athlete as having churned
- incorrectly predict 78 athletes as having been retained while in reality they are churners

#### 4.2 Using Aggregate Transactions Data

This is entailed tracking an athlete's aggregate transactions data up to a certain month and using that information to predict whether the athlete would churn in the next month. In our case, we used the cumulative sum of each feature up to March 2014 to predict churn in April 2014.

##### 4.2.1 Precision Recall Table

Table 4: Precision Recall Values

Model	Precision	Recall	F1-Score
Logistic Regression	0.00	0.00	0.00
Logistic Regression(with class_weight = 'balanced')	0.18	0.61	0.28
Random Forest	0.00	0.00	0.00
Gradient Boosting	0.50	0.02	0.03

#### 4.2.2 Confusion Matrix: Logistic Regression

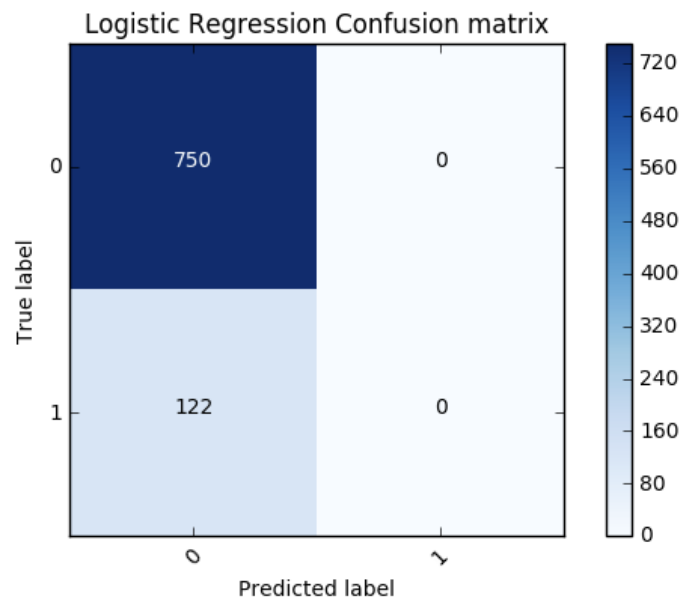


Figure 9: Lostic Regression Confusion Matrix

The above confusion matrix illustrates the following:  
Given 117 athletes who are churners the model will:

- correctly predict none of the 117 athletes.
- incorrectly predict 0 athletes as having churned
- incorrectly predict 122 athletes as having been retained while in reality 117 of them are churners

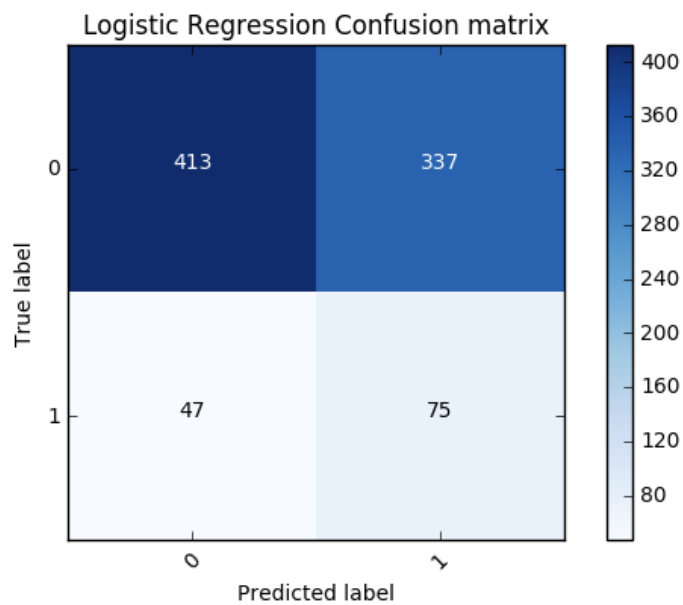


Figure 10: Lostic Regression with class\_weight='balanced'

The above confusion matrix illustrates the following:  
Given 117 athletes who are churners the model will:

- correctly predict 75 of the 117 athletes.
- incorrectly predict 337 athletes as having churned
- incorrectly predict 47 athletes as having been retained while in reality 117 of them are churners

#### 4.2.3 *Confusion Matrix: Random Forest*

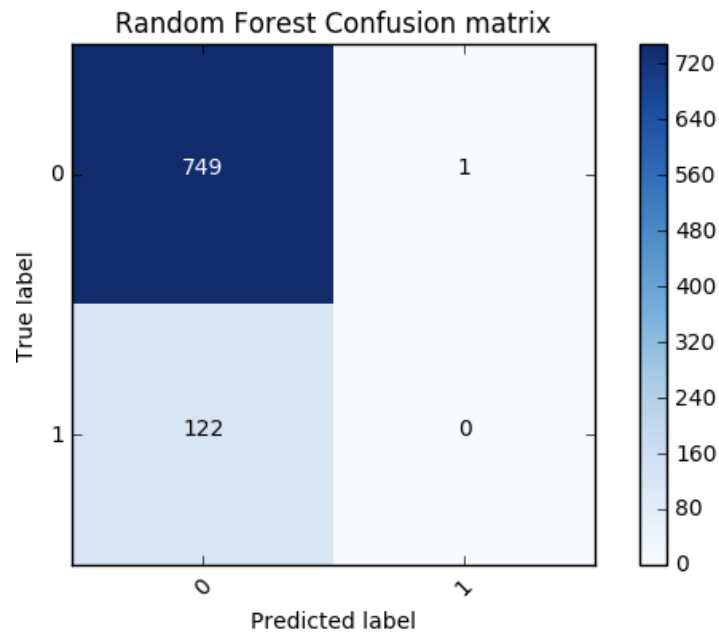


Figure 11: Random Forest Confusion Matrix.

The above confusion matrix illustrates the following:  
Given 117 athletes who are churners the model will:

- correctly predict none of the 117 athletes.
- incorrectly predict 1 athlete as having churned
- incorrectly predict 122 athletes as having been retained while in reality 117 of them are churners

#### 4.2.4 *Confusion Matrix: Gradient Boosting*

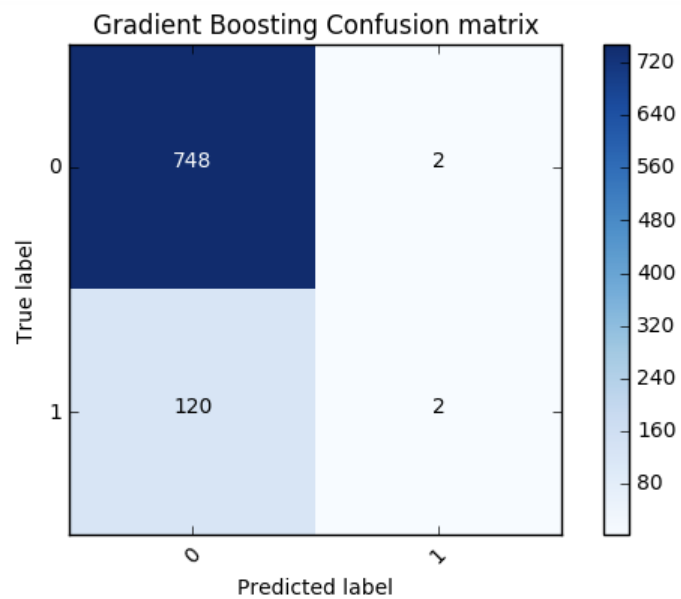


Figure 12: Gradient Boosting Confusion Matrix

The above confusion matrix illustrates the following:

Given 117 athletes who are churners the model will:

- correctly predict 2 of the 117 athletes.
- incorrectly predict 2 athletes as having churned
- incorrectly predict 120 athletes as having been retained while in reality 115 of them are churners



### 4.3 Using Monthly Difference in Transactions Data

This is entailed tracking an athlete's change in transactions data from one month to the next and using that information to predict whether the athlete would churn in the next month. For instance, we would track the change in the number of hits between February 2014 and March 2014 and Creating a Hits\_diff column. For creating our testing and training data, we tracked the difference in transactions(interactions) between February 2014 and March 2014 to predict churn in April 2014.

#### 4.3.1 Precision Recall Table

Table 5: Precision Recall Values

Model	Precision	Recall	F1-Score
Logistic Regression	0.00	0.00	0.00
Logistic Regression(with class_weight = 'balanced')	0.07	0.55	0.13
Random Forest	0.00	0.00	0.00
Gradient Boosting	0.10	0.02	0.03

#### 4.3.2 Confusion Matrix: Logistic Regression

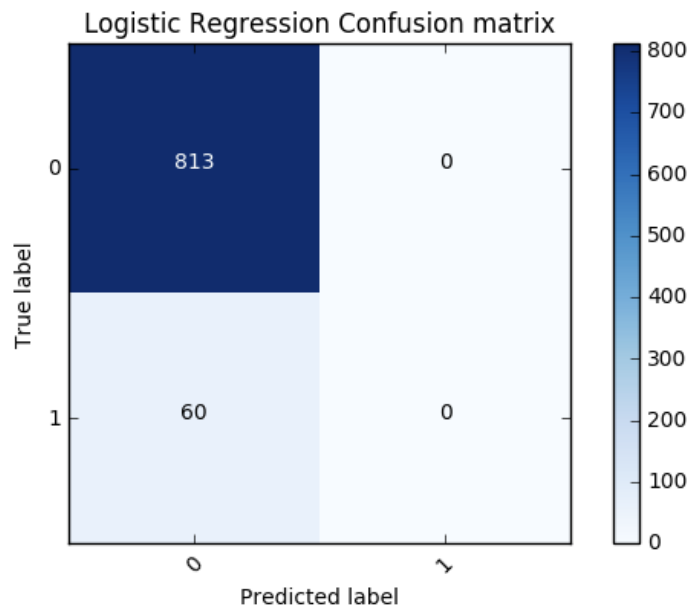


Figure 13: Lostic Regression Confusion Matrix

The above confusion matrix illustrates the following:

Given 117 athletes who are churners the model will:

- correctly predict none of the 117 athletes.
- incorrectly predict 0 athlete as having churned
- incorrectly predict 60 athletes as having been retained

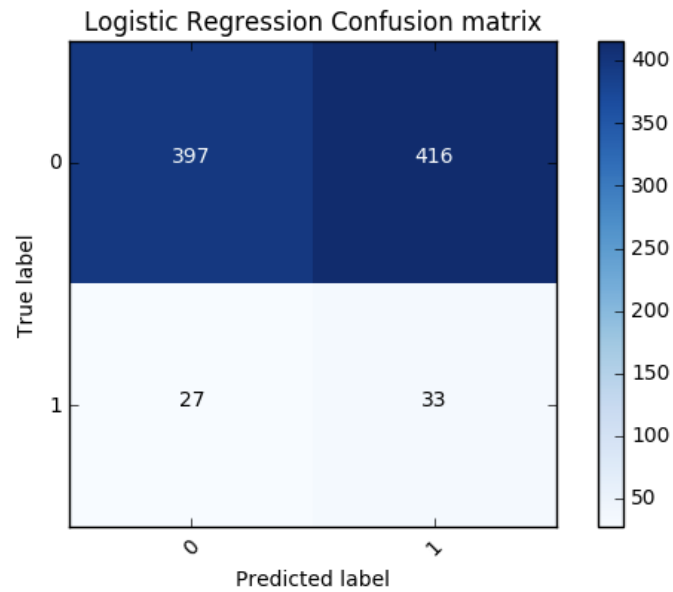


Figure 14: Logistic Regression with `class_weight='balanced'`

The above confusion matrix illustrates the following:  
Given 117 athletes who are churners the model will:

- correctly predict 33 of the 117 athletes.
- incorrectly predict 416 athlete as having churned

#### 4.3.3 Confusion Matrix: Random Forest

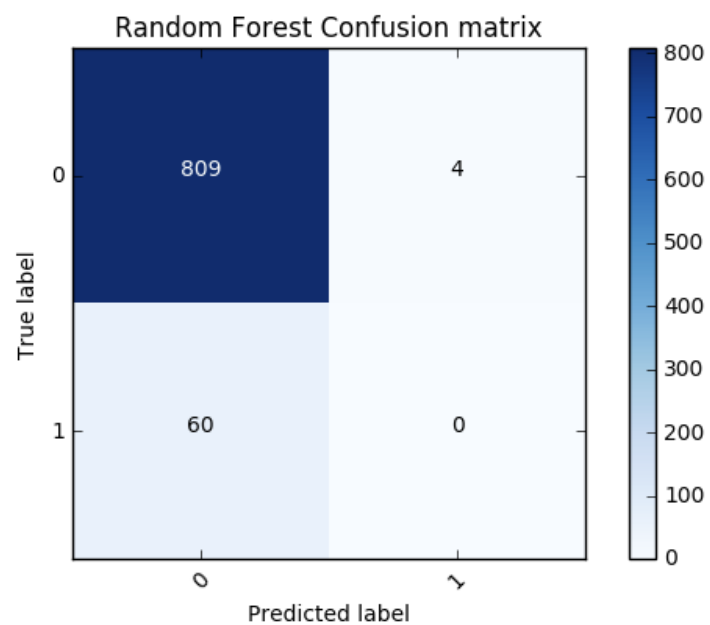


Figure 15: Random Forest Confusion Matrix. We Suspect Overfitting here

#### 4.3.4 *Confusion Matrix: Gradient Boosting*

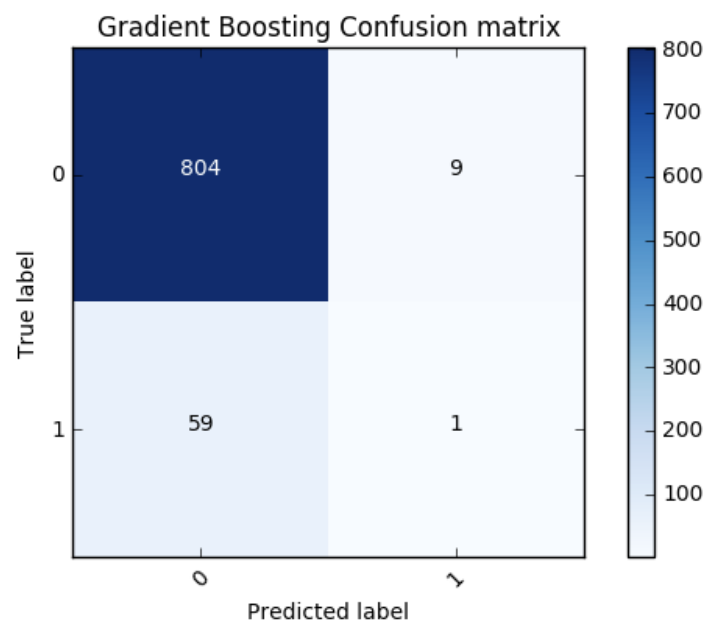


Figure 16: Gradient Boosting Confusion Matrix

The above confusion matrix illustrates the following:

Given 10 athletes who are predicted to have churned, only one actually churns.

## 5 LIFE-TIME CHURN PREDICTION

The models we built here were geared towards answering the question: Will athlete A churn at some point in their high school career before:

- Making a team
- His/her spring semester of his/her senior year

### 5.1 Dataset Attributes

To create the dataset for this model, we aggregated each feature and set the churn value(CaptainU\_CHURN) to the value of the last month the athlete was on the system or the value on December 1st 2014. This is because all athletes in the system had churned in January 2015(Spring semester of their senior year).

Table 6: My caption

Dataset	Rows	columns	Churners	Non-Churners
Full	16117	26	6043	10074
Training	12893	26	4832	8061
Testing	3224	26	1211	2013

### 5.2 Churn in both males and Females

#### 5.2.1 Precision Recall Table

Table 7: Precision Recall Values

Model	Precision	Recall	F1-Score
Logistic Regression	0.75	0.68	0.71
Logistic Regression(with class_weight = 'balanced')	0.65	0.82	0.73
Random Forest	0.79	0.72	0.75
Gradient Boosting	0.77	0.75	0.76
Support Vector Machine	0.75	0.67	0.71

### 5.2.2 Confusion Matrix: Logistic Regression

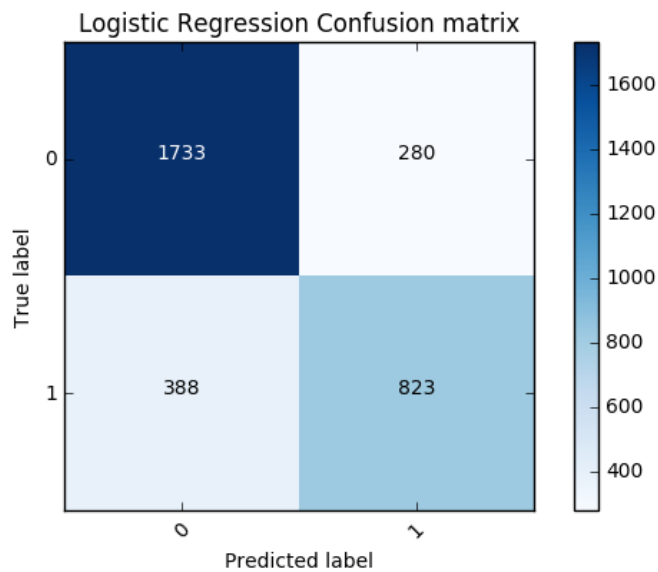


Figure 17: Lostic Regression Confusion Matrix

The above confusion matrix illustrates the following:

Given 1211 athletes who are churners the model will:

- correctly predict 823 of the 1211 athletes.
- incorrectly predict 280 athlete as having churned
- incorrectly predict 388 athletes as having been retained while in reality they churned

### 5.2.3 Confusion Matrix: Logistic Regression with balanced class weight

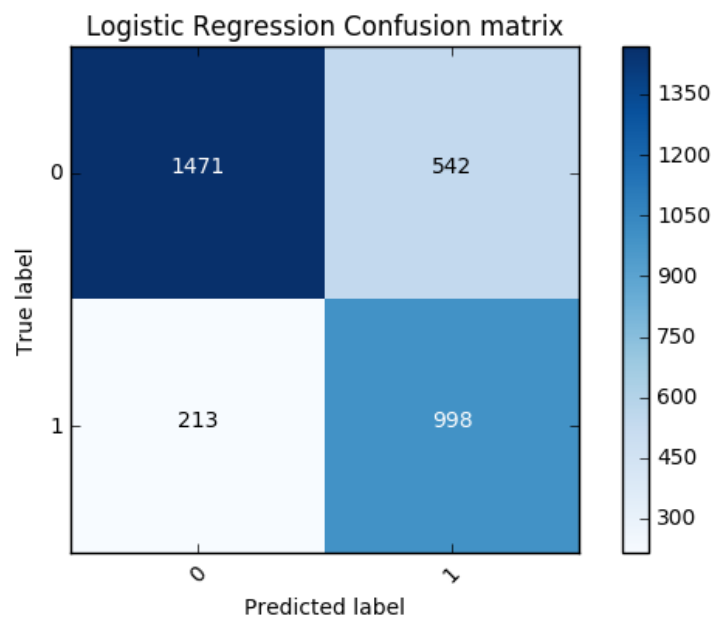


Figure 18: Lostic Regression Confusion Matrix

The above confusion matrix illustrates the following:

Given 1211 athletes who are churners the model will:

- correctly predict 998 of the 1211 athletes.
- incorrectly predict 542 athlete as having churned
- incorrectly predict 213 athletes as having been retained while in reality they churned

Effectively, this means that 2 of the 3 athletes who are predicted to churn at some point in the future are actually churners.

#### 5.2.4 Confusion Matrix: Random Forest

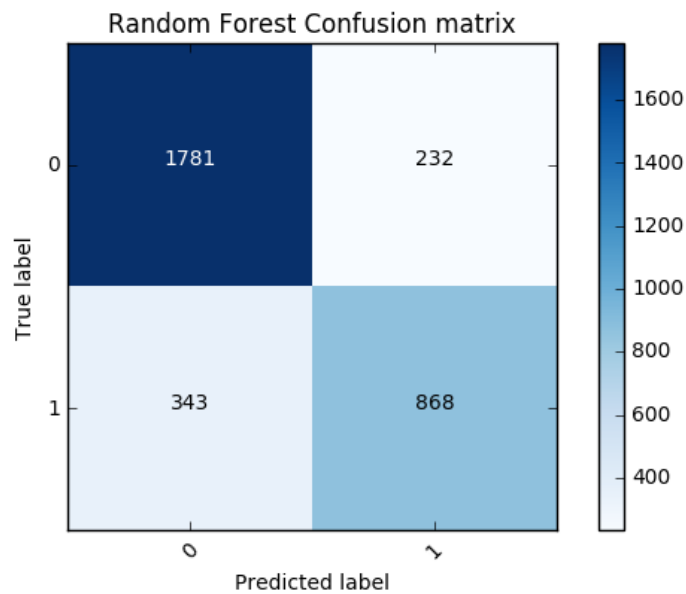


Figure 19: Random Forest Confusion Matrix

The above confusion matrix illustrates the following:

Given 1211 athletes who are churners the model will:

- correctly predict 868 of the 1211 athletes.
- incorrectly predict 232 athlete as having churned
- incorrectly predict 343 athletes as having been retained while in reality they churned

Effectively, this means that about 8 of the 10 athletes who are predicted to churn at some point in the future are actual churners.

### 5.2.5 Confusion Matrix: Gradient Boosting

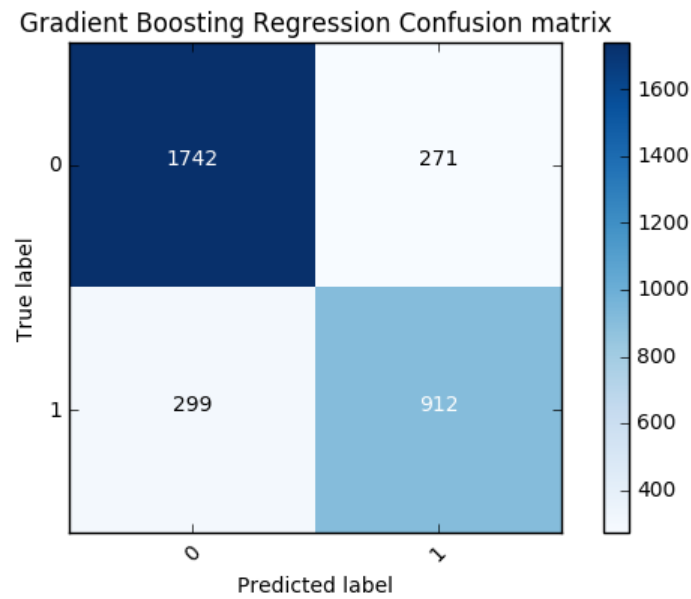


Figure 20: Random Forest Confusion Matrix

The above confusion matrix illustrates the following:

Given 1211 athletes who are churners the model will:

- correctly predict 868 of the 1211 athletes.
- incorrectly predict 232 athlete as having churned
- incorrectly predict 343 athletes as having been retained while in reality they churned

Effectively, this means that about 8 of the 10 athletes who are predicted to churn at some point in the future are actual churners.

## 5.3 Churn in Males

We went ahead and isolated males and tried to understand what factors lead to churn in males.

### 5.3.1 Dataset Attributes

Table 8: Dataset attributes

Dataset	Rows	columns	Churners	Non-Churners
Full	7342	26	2334	5008
Training	5873	26	1869	4004
Testing	1469	26	465	1004

5.3.2 Precision Recall Table

Table 9: Precision Recall Values

Model	Precision	Recall	F1-Score
Logistic Regression	0.72	0.63	0.67
Gradient Boosting	0.77	0.72	0.74

5.3.3 Important features

We visualised the most important features that lead to either churn or retention. The features in blue lead to churn while the the features in read drive retention. This is because positive coefficients increase the log-odds of the response (and thus increase the probability), and negative coefficients decrease the log-odds of the response (and thus decrease the probability).The values on the y-axis represent the weight values of each features based on the logistic regression equation:

$$F(x) = \frac{1}{1 + e^{-\beta_i}}$$

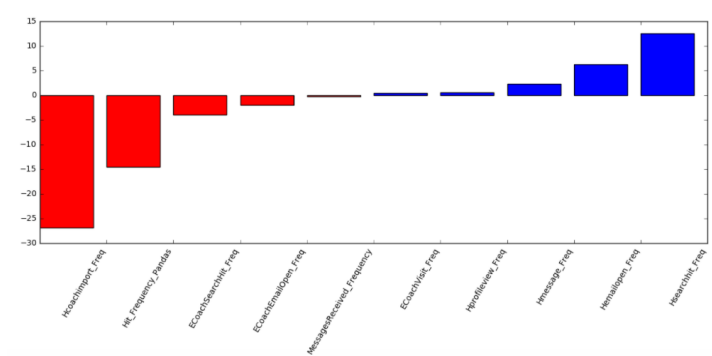


Figure 21: Important Features Affecting Churnin Males

In figure above, we can postulate that the more a coach imports an athlete’s information and the more an athlete gets profile hits, the more he is likely to stay.

5.4 Churn in Females

5.4.1 Dataset Attributes

Table 10: My caption

Dataset	Rows	columns	Churners	Non-Churners
Full	8775	26	3759	5066
Training	7020	26	2994	4076
Testing	1755	26	765	990



5.4.2 Precision Recall Table

Table 11: Precision Recall Values			
Model	Precision	Recall	F1-Score
Logistic Regression	0.72	0.75	0.74
Gradient Boosting	0.76	0.81	0.79

5.4.3 Important features

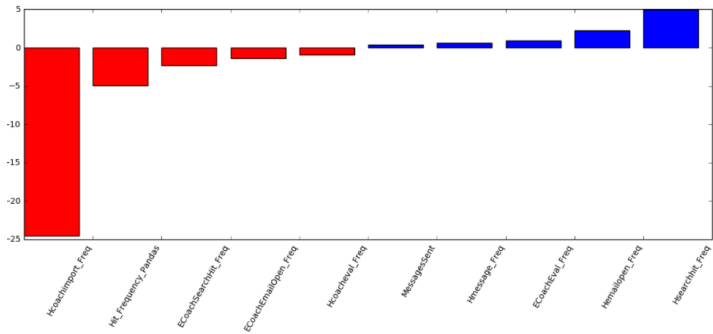


Figure 22: Important Features Affecting Churn in Females

## 6 RESULTS AND DISCUSSION

These preliminary results for predicting monthly and life time churn. With there results, we can infer the following:

- Predicting life time churn has gives better results than predicting one month churn. This is mostly due to the fact that the dataset is larger and the the split of churn and non churners is about 2:1.
- Logistic regression is more consistent in providing credible results compared to random forest and gradient boosting

## 7 FUTURE WORK

With every data science project, good results lie in data engineering. Our results mainly focused on the already assembled MSG\_RFM table, going forward, we would recommend making use of the other tables in the database.