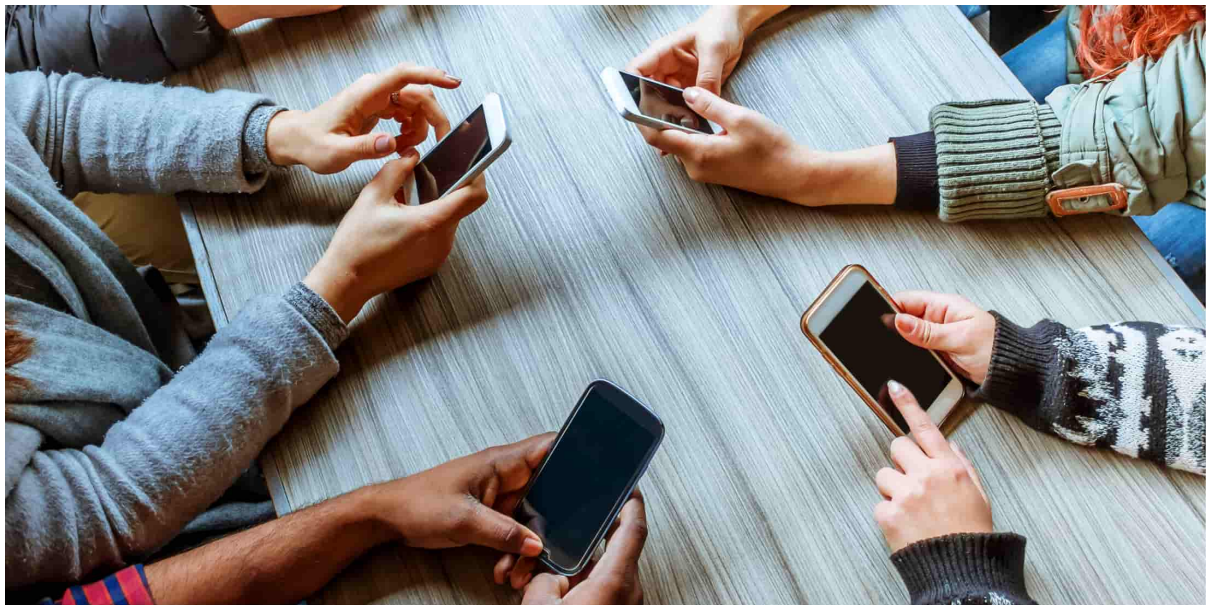


# SyriaTel Customer Churn Analysis

- Student name: Gamze Turan
- Student pace: self paced
- Scheduled project review date/time: Thu, Aug 4, 2022, 2:30 PM - 3:15 PM
- Instructor name: Claude Fried
- Blog post URL: <https://ginaturan.blogspot.com/2022/07/what-are-decision-trees-how-do-they.html> (<https://ginaturan.blogspot.com/2022/07/what-are-decision-trees-how-do-they.html>)



## Overview

I will examine the "SyriaTel Customer Churn" data in this study. The SyriaTel is a telecommunication company. To determine whether a customer will ("soon") discontinue doing business with Syria Tel is the goal of the study.

The best way to determine is to make a predictive model which will classify customers who might stop doing business with Syria Tel, using the data.

I will build a model for classifying whether customer will stop business True or False.

## Business Understanding

This search will detect which customers are likely to leave a service or to cancel a subscription to a service.

Select a model that will be the most accurate in predicting which client will discontinue doing business with SyriaTel.

## Data Understanding

The Data comes from SyriaTel and includes information about their customers. The dataset has customer's state of residence, telephone numbers and length of the account.

There are columns indicating if the customers has an international plan and voicemail plan, how many voice mails they receive.

The dataset includes how many minutes they spend talking, how many calls they make and how much they are charged during day, evening and night periods.

```
In [1]: ▶ 1 # importing necessary libraries
          2 import pandas as pd
          3 import numpy as np
          4
          5 #To help with data visualization
          6 import matplotlib.pyplot as plt
          7 import seaborn as sns
          8 %matplotlib inline
          9
         10 # To suppress warnings
         11 import warnings
         12 warnings.filterwarnings("ignore")
```

## Loading Data

```
In [2]: 1 scc = pd.read_csv("data/bigml_59c28831336c6604c800002a.csv")
        2 scc.head()
```

Out[2]:

	state	account length	area code	phone number	international plan	voice mail plan	number vmail messages	total day minutes	total day calls	total day charge	...
0	KS	128	415	382-4657	no	yes	25	265.1	110	45.07	...
1	OH	107	415	371-7191	no	yes	26	161.6	123	27.47	...
2	NJ	137	415	358-1921	no	no	0	243.4	114	41.38	...
3	OH	84	408	375-9999	yes	no	0	299.4	71	50.90	...
4	OK	75	415	330-6626	yes	no	0	166.7	113	28.34	...

5 rows × 21 columns



```
In [3]: 1 # Checking the number of rows and columns in the data
        2 scc.shape
```

Out[3]: (3333, 21)

- There are a total 21 columns and 3,333 observations in the dataset.

## Data Overview

```
In [4]: 1 # Let's Create a Copy of data
        2 data = scc.copy()
```

In [5]:

1

# Displying first 5 rows of the data

2

data.head()

Out[5]:

	state	account length	area code	phone number	international plan	voice mail plan	number vmail messages	total day minutes	total day calls	total day charge	...
0	KS	128	415	382-4657	no	yes	25	265.1	110	45.07	...
1	OH	107	415	371-7191	no	yes	26	161.6	123	27.47	...
2	NJ	137	415	358-1921	no	no	0	243.4	114	41.38	...
3	OH	84	408	375-9999	yes	no	0	299.4	71	50.90	...
4	OK	75	415	330-6626	yes	no	0	166.7	113	28.34	...

5 rows × 21 columns



In [6]:

1

# Displaying last 5 rows od the data

2

data.tail()

Out[6]:


	state	account length	area code	phone number	international plan	voice mail plan	number vmail messages	total day minutes	total day calls	total day charge
3328	AZ	192	415	414-4276	no	yes	36	156.2	77	26.55
3329	WV	68	415	370-3271	no	no	0	231.1	57	39.29
3330	RI	28	510	328-8230	no	no	0	180.8	109	30.74
3331	CT	184	510	364-6381	yes	no	0	213.8	105	36.35
3332	TN	74	415	400-4344	no	yes	25	234.4	113	39.85

5 rows × 21 columns



In [7]:    
 1 *# Let's check the data types of the columns in the dataset*  
 2 data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3333 entries, 0 to 3332
Data columns (total 21 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   state                                3333 non-null   object
1   account length                       3333 non-null   int64
2   area code                           3333 non-null   int64
3   phone number                        3333 non-null   object
4   international plan                  3333 non-null   object
5   voice mail plan                     3333 non-null   object
6   number vmail messages               3333 non-null   int64
7   total day minutes                   3333 non-null   float64
8   total day calls                     3333 non-null   int64
9   total day charge                    3333 non-null   float64
10  total eve minutes                   3333 non-null   float64
11  total eve calls                     3333 non-null   int64
12  total eve charge                    3333 non-null   float64
13  total night minutes                 3333 non-null   float64
14  total night calls                   3333 non-null   int64
15  total night charge                  3333 non-null   float64
16  total intl minutes                  3333 non-null   float64
17  total intl calls                    3333 non-null   int64
18  total intl charge                   3333 non-null   float64
19  customer service calls              3333 non-null   int64
20  churn                              3333 non-null   bool
dtypes: bool(1), float64(8), int64(8), object(4)
memory usage: 524.2+ KB
```

In [8]:    
 1 *# cheking for dublicates in the data*  
 2 data.duplicated().sum()

Out[8]: 0

```
In [9]: 1 # checking for missing values in the data
        2 data.isnull().sum()
```

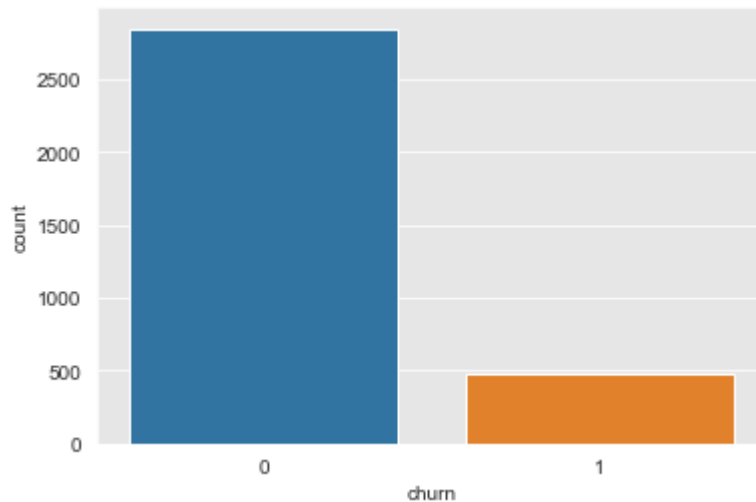
```
Out[9]: state                                0
        account length                      0
        area code                           0
        phone number                        0
        international plan                  0
        voice mail plan                     0
        number vmail messages              0
        total day minutes                   0
        total day calls                     0
        total day charge                    0
        total eve minutes                   0
        total eve calls                     0
        total eve charge                    0
        total night minutes                 0
        total night calls                   0
        total night charge                  0
        total intl minutes                  0
        total intl calls                    0
        total intl charge                   0
        customer service calls             0
        churn                              0
        dtype: int64
```

```
In [10]: 1 # Let's view the statistical summary of the numerical columns in the data
          2 data.describe().T
```

Out[10]:

	count	mean	std	min	25%	50%	75%	max
<b>account length</b>	3333.0	101.064806	39.822106	1.00	74.00	101.00	127.00	243.00
<b>area code</b>	3333.0	437.182418	42.371290	408.00	408.00	415.00	510.00	510.00
<b>number vmail messages</b>	3333.0	8.099010	13.688365	0.00	0.00	0.00	20.00	51.00
<b>total day minutes</b>	3333.0	179.775098	54.467389	0.00	143.70	179.40	216.40	350.80
<b>total day calls</b>	3333.0	100.435644	20.069084	0.00	87.00	101.00	114.00	165.00
<b>total day charge</b>	3333.0	30.562307	9.259435	0.00	24.43	30.50	36.79	59.64
<b>total eve minutes</b>	3333.0	200.980348	50.713844	0.00	166.60	201.40	235.30	363.70
<b>total eve calls</b>	3333.0	100.114311	19.922625	0.00	87.00	100.00	114.00	170.00
<b>total eve charge</b>	3333.0	17.083540	4.310668	0.00	14.16	17.12	20.00	30.91
<b>total night minutes</b>	3333.0	200.872037	50.573847	23.20	167.00	201.20	235.30	395.00
<b>total night calls</b>	3333.0	100.107711	19.568609	33.00	87.00	100.00	113.00	175.00
<b>total night charge</b>	3333.0	9.039325	2.275873	1.04	7.52	9.05	10.59	17.77
<b>total intl minutes</b>	3333.0	10.237294	2.791840	0.00	8.50	10.30	12.10	20.00
<b>total intl calls</b>	3333.0	4.479448	2.461214	0.00	3.00	4.00	6.00	20.00
<b>total intl charge</b>	3333.0	2.764581	0.753773	0.00	2.30	2.78	3.27	5.40
<b>customer service calls</b>	3333.0	1.562856	1.315491	0.00	1.00	1.00	2.00	9.00

```
In [103]: 1 # churn customer countplot
           2 sns.countplot(x = "churn", data = data);
           3 #plt.savefig('images/churn_cusmomer_count.png')
```





## Data Preprocessing

I will remove the column 'phone number' from the dataset because most digit in the phone number is random, and we will not use for modeling.

```
In [11]: 1 data = data.drop("phone number", axis=1)
```

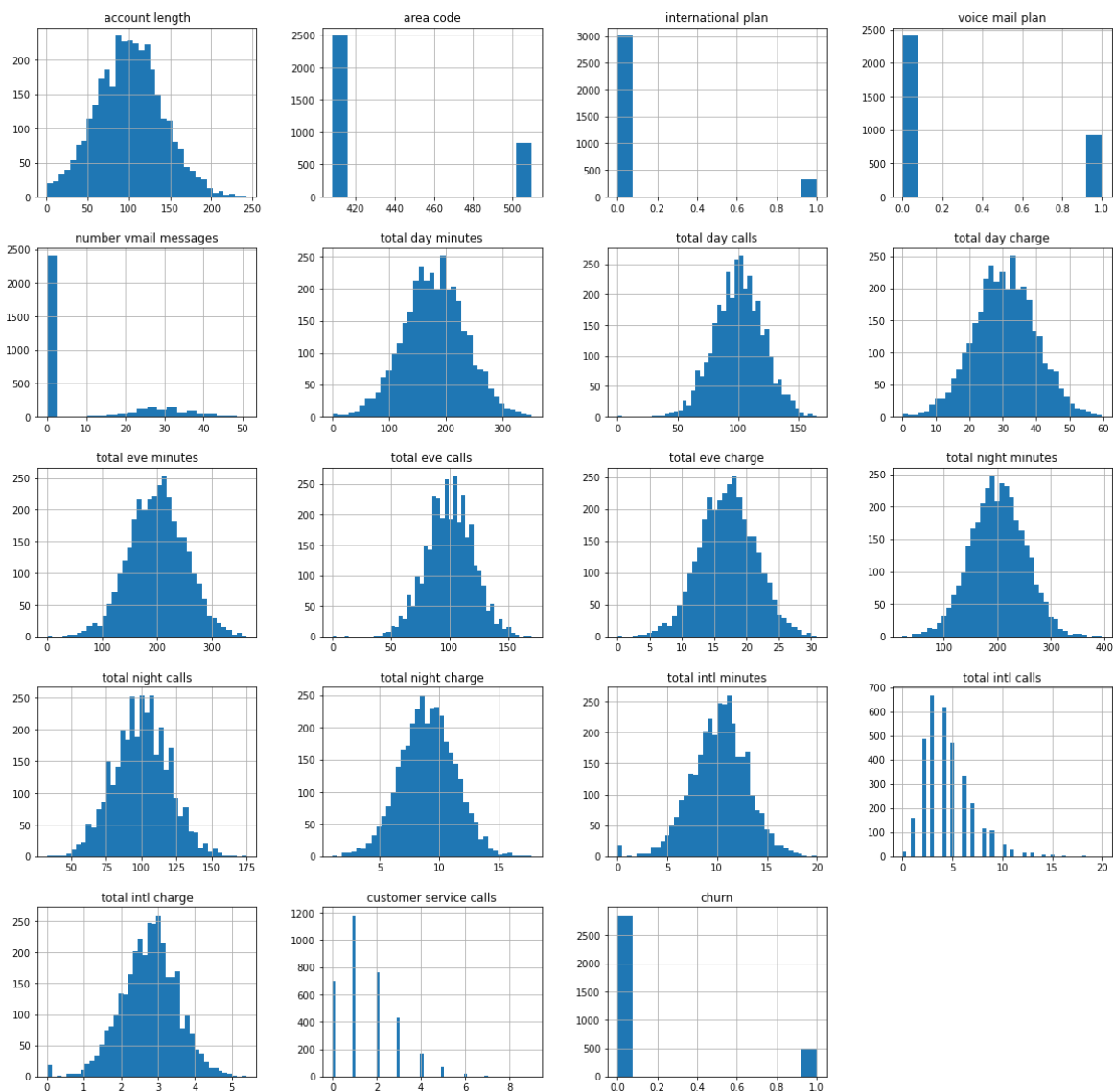
In the dataset international plan and voice mail plan are object to data type. Similarly churn columns bool data type. I will convert these columns to binary.

```
In [12]: 1 # Convert to binary
2 data["international plan"] = data["international plan"].map({"yes": 1, "no": 0})
3 data["voice mail plan"] = data["voice mail plan"].map({"yes": 1, "no": 0})
4 data["churn"] = data["churn"].map({True: 1, False: 0})
5 data.head()
```

Out[12]:

	state	account length	area code	international plan	voice mail plan	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes	total eve charge
0	KS	128	415	0	1	25	265.1	110	45.07	197.4	9.0
1	OH	107	415	0	1	26	161.6	123	27.47	195.5	10.0
2	NJ	137	415	0	0	0	243.4	114	41.38	121.2	11.0
3	OH	84	408	1	0	0	299.4	71	50.90	61.9	8.0
4	OK	75	415	1	0	0	166.7	113	28.34	148.3	12.0

```
In [13]: 1 data.hist(figsize=(20, 20), bins="auto");  
2 #plt.savefig("images/histograms_ALL.png")
```



Now, the binary variables have type int64. I will changed the dtype to object for these variables, to make them available for dummy variable creation.

The variable 'area code' is also dtype int64, however it is a categorical variable. I will also change it to object

```
In [14]: 1 # changing binary variable dtypes int64 to object to create dummy variab
2 # changing categorical variable dtype object
3 data = data.astype({"international plan": "object"})
4 data = data.astype({"voice mail plan": "object"})
5 data = data.astype({"area code": "object"})
```

```
In [15]: 1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3333 entries, 0 to 3332
Data columns (total 20 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   state                                3333 non-null   object
1   account length                       3333 non-null   int64
2   area code                            3333 non-null   object
3   international plan                   3333 non-null   object
4   voice mail plan                     3333 non-null   object
5   number vmail messages               3333 non-null   int64
6   total day minutes                   3333 non-null   float64
7   total day calls                     3333 non-null   int64
8   total day charge                    3333 non-null   float64
9   total eve minutes                   3333 non-null   float64
10  total eve calls                     3333 non-null   int64
11  total eve charge                    3333 non-null   float64
12  total night minutes                 3333 non-null   float64
13  total night calls                   3333 non-null   int64
14  total night charge                  3333 non-null   float64
15  total intl minutes                  3333 non-null   float64
16  total intl calls                    3333 non-null   int64
17  total intl charge                   3333 non-null   float64
18  customer service calls              3333 non-null   int64
19  churn                              3333 non-null   int64
dtypes: float64(8), int64(8), object(4)
memory usage: 520.9+ KB
```

In [16]: 1 data.describe().T

Out[16]:

	count	mean	std	min	25%	50%	75%	max
<b>account length</b>	3333.0	101.064806	39.822106	1.00	74.00	101.00	127.00	243.00
<b>number vmail messages</b>	3333.0	8.099010	13.688365	0.00	0.00	0.00	20.00	51.00
<b>total day minutes</b>	3333.0	179.775098	54.467389	0.00	143.70	179.40	216.40	350.80
<b>total day calls</b>	3333.0	100.435644	20.069084	0.00	87.00	101.00	114.00	165.00
<b>total day charge</b>	3333.0	30.562307	9.259435	0.00	24.43	30.50	36.79	59.64
<b>total eve minutes</b>	3333.0	200.980348	50.713844	0.00	166.60	201.40	235.30	363.70
<b>total eve calls</b>	3333.0	100.114311	19.922625	0.00	87.00	100.00	114.00	170.00
<b>total eve charge</b>	3333.0	17.083540	4.310668	0.00	14.16	17.12	20.00	30.91
<b>total night minutes</b>	3333.0	200.872037	50.573847	23.20	167.00	201.20	235.30	395.00
<b>total night calls</b>	3333.0	100.107711	19.568609	33.00	87.00	100.00	113.00	175.00
<b>total night charge</b>	3333.0	9.039325	2.275873	1.04	7.52	9.05	10.59	17.77
<b>total intl minutes</b>	3333.0	10.237294	2.791840	0.00	8.50	10.30	12.10	20.00
<b>total intl calls</b>	3333.0	4.479448	2.461214	0.00	3.00	4.00	6.00	20.00
<b>total intl charge</b>	3333.0	2.764581	0.753773	0.00	2.30	2.78	3.27	5.40
<b>customer service calls</b>	3333.0	1.562856	1.315491	0.00	1.00	1.00	2.00	9.00
<b>churn</b>	3333.0	0.144914	0.352067	0.00	0.00	0.00	0.00	1.00

In [17]: 1 data['total charge'] = (data['total day charge'] + data['total eve charge'])

## Scatter Plot of Total Customer Spend Over Time

I create a scatter plot to check total customer spending over time. You can see a line of cancelling customers above the staying ones, indicating higher spend for some cancelling customer.

```
In [18]: ▶ 1 plt.figure(figsize=(12, 7), dpi=80)
2 plt.scatter(data['account length'], data['total charge'], data['churn'])
3 plt.scatter(data['account length'], data['total charge'], data['churn'])
4 plt.title("Total Customer Spend Over Time", fontsize = 16)
5 plt.xlabel("Account Length in Months", fontsize = 14)
6 plt.ylabel("Total Dollars Spent By Customer", fontsize = 14)
7 plt.legend(bbox_to_anchor=(1.05, 1.0), loc = 'center')
8 plt.legend(fontsize= 12)
9 plt.tight_layout()
10 #plt.savefig('images/total_customer_spend.png')
11 plt.show()
```



**Heatmap to find the correlation between variables**

```
In [19]: 1 plt.figure(figsize=(15, 7))
2 sns.heatmap(data.corr(), annot=True, vmin=-1, vmax=1, fmt=".2f", cmap="S")
3 plt.savefig('images/heatmap.png')
4 plt.show()
```



Total day minutes and Total day charges, Total evening minutes and Total evening charge, Total night minutes and Total Night charges, Total intl minutes and Total intl charge have positive correlation which make sense that customer take minutes if the amount of charges is high.

Other variables have no significant correlation between them

## Data Preparation for Modeling

## Split Data

```
In [20]: 1 # Seperate data into train and test split
          2 from sklearn.model_selection import train_test_split
          3
          4 df = data.copy()
```

```
In [21]: 1 X = df.drop(["churn"], axis=1)
          2 y = df["churn"]
```

```
In [22]: 1 # creating dummies
          2 X = pd.get_dummies(X)
          3 X
```

Out[22]:

	account length	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes	total eve calls	total eve charge	total night minutes	total night calls	...
0	128	25	265.1	110	45.07	197.4	99	16.78	244.7	91	...
1	107	26	161.6	123	27.47	195.5	103	16.62	254.4	103	...
2	137	0	243.4	114	41.38	121.2	110	10.30	162.6	104	...
3	84	0	299.4	71	50.90	61.9	88	5.26	196.9	89	...
4	75	0	166.7	113	28.34	148.3	122	12.61	186.9	121	...
...	...	...	...	...	...	...	...	...	...	...	...
3328	192	36	156.2	77	26.55	215.5	126	18.32	279.1	83	...
3329	68	0	231.1	57	39.29	153.4	55	13.04	191.3	123	...
3330	28	0	180.8	109	30.74	288.8	58	24.55	191.9	91	...
3331	184	0	213.8	105	36.35	159.6	84	13.57	139.2	137	...
3332	74	25	234.4	113	39.85	265.9	82	22.60	241.4	77	...

3333 rows × 74 columns



```
In [23]: 1 # Splitting data into training, validation and test sets:
2 # First I split data into 2 parts, say temporary and test
3
4 X_temp, X_test, y_temp, y_test = train_test_split(X, y, test_size=0.2, r
5
6 # then we spit the temporary set into train and validation
7
8 X_train, X_val, y_train, y_val = train_test_split(X_temp, y_temp, test_s
9
10 print(X_train.shape, X_val.shape, X_test.shape)
```

(1999, 74) (667, 74) (667, 74)

```
In [24]: 1 print("X_train shape = ", X_train.shape)
2 print("y_train shape = ", y_train.shape)
3 print("X_test shape = ", X_test.shape)
4 print("y_test shape = ", y_test.shape)
```

X\_train shape = (1999, 74)  
y\_train shape = (1999,)  
X\_test shape = (667, 74)  
y\_test shape = (667,)

```
In [25]: 1 print("Number of rows in train data =", X_train.shape[0])
2 print("Number of rows in validation data =", X_val.shape[0])
3 print("Number of rows in test data =", X_test.shape[0])
```

Number of rows in train data = 1999  
Number of rows in validation data = 667  
Number of rows in test data = 667

```
In [26]: 1 print("Train percent :", y_train.value_counts(normalize=True)[1])
2 print("Test percent : ", y_test.value_counts(normalize=True)[1])
```

Train percent : 0.14457228614307155  
Test percent : 0.1454272863568216



```

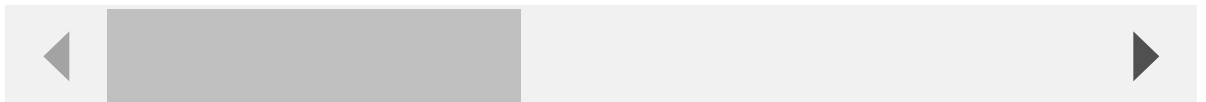
In [27]: 1 # Scale/Normalize the predictor
2 from sklearn.preprocessing import StandardScaler
3
4 scaler = StandardScaler()
5 X_train_scaled = scaler.fit_transform(X_train)
6 X_test_scaled = scaler.transform(X_test)
7
8 # Convert to DataFrame
9 X_train_scaled = pd.DataFrame(X_train_scaled, columns=X_train.columns)
10 X_test_scaled = pd.DataFrame(X_test_scaled, columns=X_test.columns)
11 X_train_scaled.head()

```

Out[27]:

	account length	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes	total eve calls	total eve charge	mi
0	-0.575833	-0.587153	-0.644409	-0.499863	-0.644463	-1.290603	-0.472427	-1.291353	0.20
1	-1.373747	-0.587153	-2.376847	0.990630	-2.376471	-0.448284	1.108706	-0.447538	-0.98
2	-0.451159	-0.587153	1.968657	-0.003032	1.968910	0.624646	-0.175964	0.623901	0.42
3	1.294277	-0.587153	-0.455317	1.089996	-0.454842	0.241597	-1.460635	0.242229	-0.47
4	0.296885	-0.587153	-0.183385	-0.748279	-0.183653	-2.793487	-0.719479	-2.792746	-0.45

5 rows × 74 columns



## Logistic Regression Model

Here, a ROC Curve can be used to show the performance of a classification model at all classification thresholds. This curve plots two parameters: true positive rate and false positive rate.

Area Under the ROC Curve(AUC): An aggregated metric that evaluates how well a logistic regression model classifies positive and negative outcomes at all possible cutoffs. It can range from 0.5 to 1, and the larger the better.

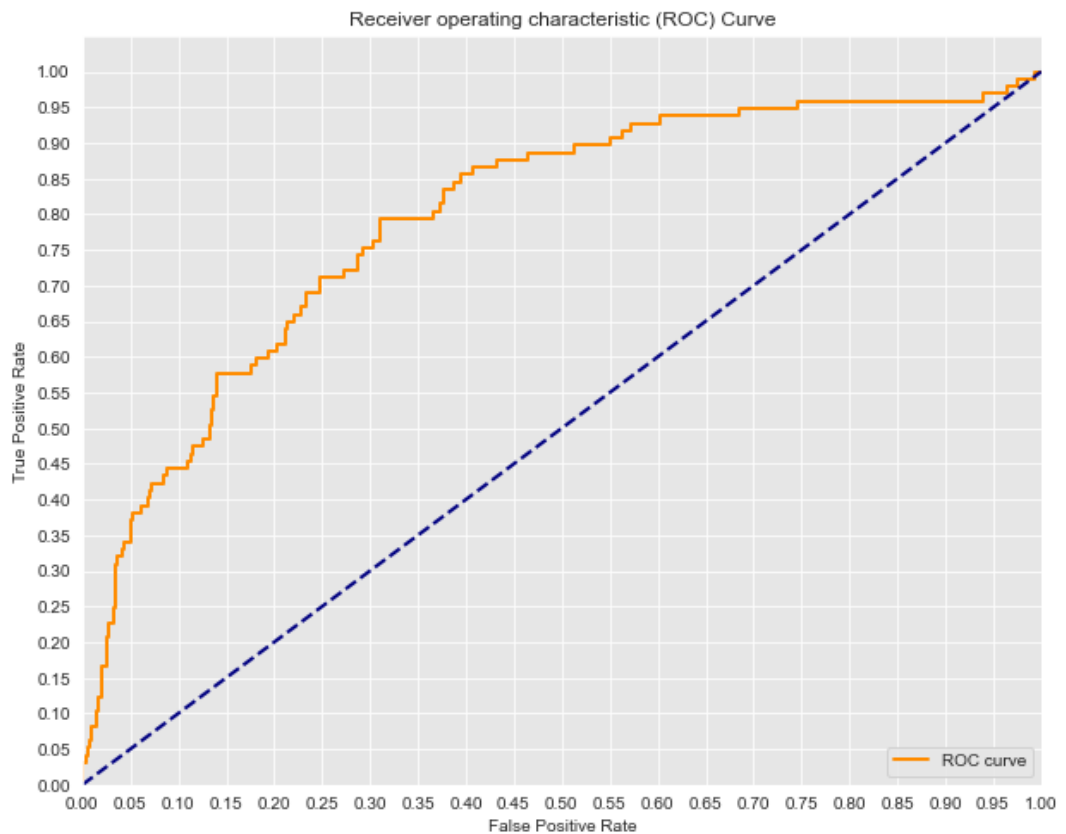
The AUC of .79 indicates that this model is sorting the values at an acceptable way, but not an excellent way.

```

In [28]: ▶ 1 from sklearn.linear_model import LinearRegression, LogisticRegression, R
2 from sklearn.metrics import accuracy_score, roc_auc_score, roc_curve, me
3 from sklearn.metrics import precision_score, recall_score, accuracy_scor
4
5
6 # Initial Model # what does the solver indicate?
7 logreg = LogisticRegression(fit_intercept=False, solver='liblinear')
8
9 # Probability scores for test set
10 y_score = logreg.fit(X_train, y_train).decision_function(X_test)
11
12 # False positive rate and true positive rate
13 fpr, tpr, thresholds = roc_curve(y_test, y_score)
14
15 # Seaborn's beautiful styling
16 sns.set_style('darkgrid', {'axes.facecolor': '0.9'})
17
18 #Print AUC
19 print('AUC: {}'.format(auc(fpr, tpr)))
20
21 # Plot the ROC curve
22 plt.figure(figsize=(10, 8))
23 lw = 2
24 plt.plot(fpr, tpr, color = 'darkorange',
25          lw=lw, label='ROC curve')
26 plt.plot([0,1], [0,1], color = 'navy', lw=lw, linestyle = '--')
27 plt.xlim([0.0,1.0])
28 plt.ylim([0.0, 1.05])
29 plt.yticks([i/20.0 for i in range(21)])
30 plt.xticks([i/20.0 for i in range(21)])
31 plt.xlabel('False Positive Rate')
32 plt.ylabel('True Positive Rate')
33 plt.title('Receiver operating characteristic (ROC) Curve')
34 plt.legend(loc='lower right');
35 #plt.savefig("images/LogisticRegressionModel.png")

```

AUC: 0.7945378911195516



In [29]: 1 `y_pred = logreg.predict(X_test)`

## Resampling

I work on resampling of the data to find the number of perentage if customer stay with SyriaTel or not with the normalized true.

```
In [30]: 1 print('Original whole data class distribution:')
2 print(y.value_counts())
3 print('Original whole data class distribution, normalized:')
4 print(y.value_counts(normalize=True))
```

Original whole data class distribution:

0 2850

1 483

Name: churn, dtype: int64

Original whole data class distribution, normalized:

0 0.855086

1 0.144914

Name: churn, dtype: float64

The analysis shows that 85.5% of customers stay with SyriaTel while 14.5% of customers discontinue doing business with the company. Thus, we will have 85.5% accuracy if we forecast that all consumers will continue. This explains why the model's accuracy score is high despite the other metrics' low values.

By using SMOTE I create a synthetic training sample to take care of imbalance.

```
In [31]: ▶ 1 # Import SMOTE, resample
2 from imblearn.over_sampling import SMOTE
3
4 smote = SMOTE()
5 X_train_scaled_resampled, y_train_resampled = smote.fit_resample(X_train,
6
7 print('Original training data class distribution:')
8 print(y_train.value_counts())
9 print('Synthetic training data class distribution:')
10 print(y_train_resampled.value_counts())
```

Original training data class distribution:

0 1710

1 289

Name: churn, dtype: int64

Synthetic training data class distribution:

1 1710

0 1710

Name: churn, dtype: int64

```
In [32]: 1 # New model after resampling
2 logreg = LogisticRegression(random_state=42)
3 logreg.fit(X_train_scaled_resampled, y_train_resampled)
4
5 print('Training Data:\n', classification_report(y_train_resampled, logreg.predict(X_train_scaled_resampled)))
6 print('Testing Data:\n', classification_report(y_test, logreg.predict(X_test_scaled_resampled)))
```

Training Data:

	precision	recall	f1-score	support
0	0.82	0.80	0.81	1710
1	0.81	0.82	0.82	1710
accuracy			0.81	3420
macro avg	0.81	0.81	0.81	3420
weighted avg	0.81	0.81	0.81	3420

Testing Data:

	precision	recall	f1-score	support
0	0.93	0.78	0.85	570
1	0.34	0.65	0.44	97
accuracy			0.76	667
macro avg	0.63	0.72	0.65	667
weighted avg	0.84	0.76	0.79	667

After resampling, the Logistic Regression Model performance is clearly improved.

The performance in training data is better test data. This is sign of overfitting.

## Parameter Tuning

I initially used the default parameters for the Logistic Regression model. I will now apply parameter tuning with GridSearchCV. It will determine the best parameter combination for the given parameter grid.

```
In [33]: 1 print('Default parameters:')  
2 logreg.get_params()
```

Default parameters:

```
Out[33]: {'C': 1.0,  
          'class_weight': None,  
          'dual': False,  
          'fit_intercept': True,  
          'intercept_scaling': 1,  
          'l1_ratio': None,  
          'max_iter': 100,  
          'multi_class': 'auto',  
          'n_jobs': None,  
          'penalty': 'l2',  
          'random_state': 42,  
          'solver': 'lbfgs',  
          'tol': 0.0001,  
          'verbose': 0,  
          'warm_start': False}
```

```

In [34]: 1 # Tuning Logistic Regression model with GridSearchCV
2 from sklearn.model_selection import GridSearchCV
3
4 logreg_param_grid = {
5     'solver': ['lbfgs', 'liblinear'],
6     'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000, 1e5, 1e20],
7 }
8
9 logreg_gs = GridSearchCV(logreg, logreg_param_grid, cv=5, scoring='f1')
10 #logreg_gs.fit(X_train_scaled, y_train)
11 logreg_gs.fit(X_train_scaled_resampled, y_train_resampled)
12
13 #score_logreg_gs = logreg_gs.score(X_test_scaled, y_test)
14 #print('f1-score for test data:', score_logreg_gs)
15
16 print('Parameter Tuning Results:\n')
17 print("Best Parameter Combination:", logreg_gs.best_params_)
18 print('Training Data:\n', classification_report(y_train_resampled, logreg_gs.predict(X_train_scaled_resampled)))
19 print('Testing Data:\n', classification_report(y_test, logreg_gs.predict(X_test_scaled)))

```

Parameter Tuning Results:

Best Parameter Combination: {'C': 100000.0, 'solver': 'lbfgs'}

Training Data:

	precision	recall	f1-score	support
0	0.82	0.80	0.81	1710
1	0.80	0.83	0.81	1710
accuracy			0.81	3420
macro avg	0.81	0.81	0.81	3420
weighted avg	0.81	0.81	0.81	3420

Testing Data:

	precision	recall	f1-score	support
0	0.93	0.78	0.85	570
1	0.34	0.66	0.45	97
accuracy			0.77	667
macro avg	0.64	0.72	0.65	667
weighted avg	0.85	0.77	0.79	667

It appears that the performance wasn't significantly improved by parameter adjustment using the provided parameter grid. Overfitting has been noticed.

## Logistic Regression Evaluation

```
In [35]: 1 con_matrix = confusion_matrix(y_test, y_pred)
          2 con_matrix
```

```
Out[35]: array([[556, 14],
                [ 79, 18]], dtype=int64)
```

## Confusion Matrix Breakdown in this Order

True Negatives: Predicting that they will not cancel and being correct.

False Positives: Predicting that they will cancel and being wrong.

False Negatives: Predicting that they're not going to cancel and being wrong.

True Positive: Predicting that they will cancel and being correct.

## Recall Calculation

A Type II error would be more detrimental to this project than a Type I error. In the event of a Type II error, SyriaTel would have predicted incorrectly that their customer would not churn, which would suggest a false negative. In contrast to a Type I error or false positive, in which SyriaTel incorrectly predicted that a client would churn but they did not, this is significantly worse. Recall is the best metric to aim for because a Type II Error is a worse case situation when it comes to the practical application of the results. Recall here gauges how well the model can forecast cancellations. Recall at .14 is poor, hence alternative models should be used.

```
In [36]: 1 recall_score(y_test, y_pred)
```

```
Out[36]: 0.18556701030927836
```

## K-Nearest Neighbors Model



```
In [37]: 1 # import K-nearest Neighbor Library
2 from sklearn.neighbors import KNeighborsClassifier
3
4 knn = KNeighborsClassifier()
5 knn.fit(X_train, y_train)
6 test_preds = knn.predict(X_test)
```

```
In [38]: 1 # Complete the function
2 def print_metrics(labels, preds):
3     print("Precision Score: {}".format(precision_score(labels, preds)))
4     print("Recall Score: {}".format(recall_score(labels, preds)))
5     print("Accuracy Score: {}".format(accuracy_score(labels, preds)))
6     print("F1 Score: {}".format(f1_score(labels, preds)))
7
8 print_metrics(y_test, test_preds)
```

Precision Score: 0.47058823529411764  
Recall Score: 0.08247422680412371  
Accuracy Score: 0.8530734632683659  
F1 Score: 0.14035087719298245

## K-Nearest Neighbors Evaluation

```
In [39]: 1 def find_best_k(X_train, y_train, X_test, y_test, min_k=1, max_k=25):
2     best_k = 0
3     best_score = 0.0
4     for k in range(min_k, max_k+1, 2):
5         knn = KNeighborsClassifier(n_neighbors=k)
6         knn.fit(X_train, y_train)
7         preds = knn.predict(X_test)
8         recall = recall_score(y_test, preds)
9         if recall > best_score:
10             best_k = k
11             best_score = recall
12
13     print("Best Value for k: {}".format(best_k))
14     print("Recall: {}".format(best_score))
```

```
In [40]: 1 find_best_k(X_train, y_train, X_test, y_test)
```

Best Value for k: 1  
Recall: 0.17525773195876287

```
In [41]: 1 print('Default parameters:')  
        2 knn.get_params()
```

Default parameters:

```
Out[41]: {'algorithm': 'auto',  
          'leaf_size': 30,  
          'metric': 'minkowski',  
          'metric_params': None,  
          'n_jobs': None,  
          'n_neighbors': 5,  
          'p': 2,  
          'weights': 'uniform'}
```

```

In [42]: ▶ 1 # Tuning KNN model with GridSearchCV
2 # Takes about 10 minutes on my PC
3
4 knn_param_grid = {
5     'n_neighbors': [3, 4, 5, 6, 7, 8],
6     'p': [1, 2, 3, 4]
7 }
8
9 knn_gs = GridSearchCV(knn, knn_param_grid, cv=5, scoring='f1')
10 #knn_gs.fit(X_train_scaled, y_train)
11 knn_gs.fit(X_train_scaled_resampled, y_train_resampled) # Lower performance
12
13 # score_knn_gs = knn_gs.score(X_test_scaled, y_test)
14 #print('f1-score for test data:', score_knn_gs)
15
16 print('Parameter Tuning Results:\n')
17 print("Best Parameter Combination:", knn_gs.best_params_)
18 print('Training Data:\n', classification_report(y_train_resampled, knn_gs.predict(X_train_scaled_resampled)))
19 print('Testing Data:\n', classification_report(y_test, knn_gs.predict(X_test_scaled)))

```

Parameter Tuning Results:

Best Parameter Combination: {'n\_neighbors': 4, 'p': 1}

Training Data:

	precision	recall	f1-score	support
0	0.99	0.93	0.96	1710
1	0.93	0.99	0.96	1710
accuracy			0.96	3420
macro avg	0.96	0.96	0.96	3420
weighted avg	0.96	0.96	0.96	3420

Testing Data:

	precision	recall	f1-score	support
0	0.89	0.87	0.88	570
1	0.31	0.35	0.33	97
accuracy			0.79	667
macro avg	0.60	0.61	0.61	667
weighted avg	0.80	0.79	0.80	667

The performance of the fitting on resampled training data is better. For test data, the f1-score increased from 0.15 to 0.29.(The findings for the resampled data are tested; they are not displayed here.)

Overfitting was noticed.

## Decision Tree Model

```
In [43]: 1 # import library
2 from sklearn.tree import DecisionTreeClassifier
3
4 dt = DecisionTreeClassifier(criterion='entropy', random_state=1)
5 dt.fit(X_train, y_train)
```

Out[43]: DecisionTreeClassifier(criterion='entropy', random\_state=1)

```
In [44]: 1 # Make predictions using test set
2 y_pred = dt.predict(X_test)
3
4 # Check the AUC of predictions
5 false_positive_rate, true_positive_rate, thresholds = roc_curve(y_test,
6 roc_auc = auc(false_positive_rate, true_positive_rate)
7 roc_auc
```

Out[44]: 0.834400434074878

```
In [45]: 1 recall_score(y_test, y_pred)
```

Out[45]: 0.7319587628865979

```
In [46]: 1 dt = DecisionTreeClassifier(criterion='entropy',
2                                     max_features=4,
3                                     max_depth=3,
4                                     min_samples_split=0.7,
5                                     min_samples_leaf=0.25,
6                                     random_state= 1)
7 dt.fit(X_train, y_train)
8 y_pred = dt.predict(X_test)
9 false_positive_rate, true_positive_rate, thresholds = roc_curve(y_test,
10 roc_auc = auc(false_positive_rate, true_positive_rate)
11 roc_auc
12
```

Out[46]: 0.5

```
In [47]: 1 print('Recall: ', recall_score(y_test, y_pred))
```

Recall: 0.0

## Tuning the Decison Tree Model

- Train the AUC, and Test AUC will be indicative of training and test error for learning.

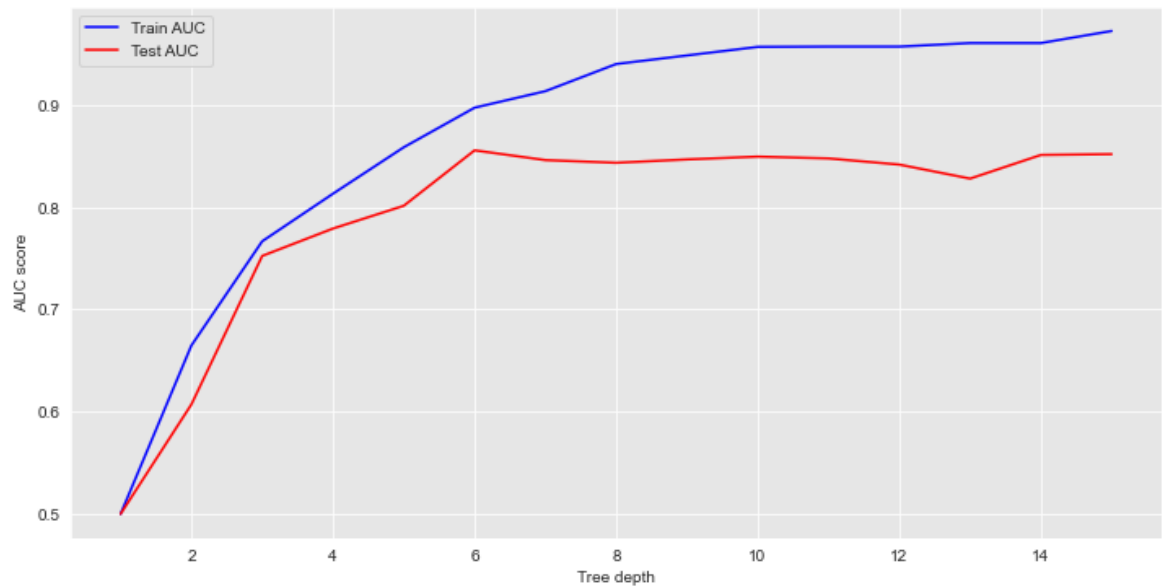
## Identifying Ideal Maximum Tree Depth

- 5 is the ideal maximum tree depth. Greather tree depth is indicative of overfitting as Train AUC soars above Test AUC. At 6, the Test AUC is above the Train AUC.

```

In [48]: ▶ 1 # Identify the optimal tree depth for given data
2 max_depths = np.linspace(1, 15, 15, endpoint=True)
3 train_results = []
4 test_results = []
5 for max_depth in max_depths:
6     dt = DecisionTreeClassifier(criterion='entropy', max_depth=max_depth)
7     dt.fit(X_train, y_train)
8     train_pred = dt.predict(X_train)
9     false_positive_rate, true_positive_rate, thresholds = roc_curve(y_train, train_pred)
10    roc_auc = auc(false_positive_rate, true_positive_rate)
11    # Add auc score to previous train results
12    train_results.append(roc_auc)
13    y_pred = dt.predict(X_test)
14    false_positive_rate, true_positive_rate, thresholds = roc_curve(y_test, y_pred)
15    roc_auc = auc(false_positive_rate, true_positive_rate)
16    # Add auc score to previous test results
17    test_results.append(roc_auc)
18
19
20 plt.figure(figsize=(12,6))
21 plt.plot(max_depths, train_results, 'b', label='Train AUC')
22 plt.plot(max_depths, test_results, 'r', label='Test AUC')
23 plt.ylabel('AUC score')
24 plt.xlabel('Tree depth')
25 plt.legend();
26 #plt.savefig("images/AUC_Score.png")

```



```
In [49]: 1 print(confusion_matrix(y_test, y_pred))
          2 print(classification_report(y_test, y_pred))
```

```
[[536  34]
 [ 23  74]]
```

		precision	recall	f1-score	support
	0	0.96	0.94	0.95	570
	1	0.69	0.76	0.72	97
	accuracy			0.91	667
	macro avg	0.82	0.85	0.84	667
	weighted avg	0.92	0.91	0.92	667

```
In [50]: 1 print('Testing Accuracy for Decision Tree Classifier: {:.4}%'.format(acc
```

Testing Accuracy for Decision Tree Classifier: 91.45%

## Parameter Tuning

```
In [51]: 1 print('Default parameters:')
          2 dt.get_params()
```

Default parameters:

```
Out[51]: {'ccp_alpha': 0.0,
          'class_weight': None,
          'criterion': 'entropy',
          'max_depth': 15.0,
          'max_features': None,
          'max_leaf_nodes': None,
          'min_impurity_decrease': 0.0,
          'min_impurity_split': None,
          'min_samples_leaf': 1,
          'min_samples_split': 2,
          'min_weight_fraction_leaf': 0.0,
          'presort': 'deprecated',
          'random_state': 1,
          'splitter': 'best'}
```

```

In [52]: 1 # Tuning Decision Trees model with GridSearchCV
2 # Takes more then 10 minutes on my PC
3
4 dt_param_grid = {
5     'criterion': ['gini', 'entropy'],
6     'max_depth': [2, 3, 4, 5, 6],
7     'min_samples_split': [2, 3, 4, 5, 6],
8     #'min_samples_leaf': [1, 2, 3, 4, 5, 6]
9 }
10
11 dt_gs = GridSearchCV(dt, dt_param_grid, cv=5, scoring='f1')
12 #dt_gs.fit(X_tarin_scaled, y_train)
13 dt_gs.fit(X_train_scaled_resampled, y_train_resampled)
14
15 #score_dt_gs = dt_gs.score(X_test_scaled, y_test)
16 # print('f1_score for test data:', score_dt_gs)
17
18 print('Parameter Tuning Results:\n')
19 print("Best Parameter Combination:", dt_gs.best_params_)
20 print('Training Data:\n', classification_report(y_train_resampled, dt_gs
21 print('Testing Data:\n', classification_report(y_test, dt_gs.predict(X_t

```

Parameter Tuning Results:

Best Parameter Combination: {'criterion': 'gini', 'max\_depth': 6, 'min\_samples\_split': 3}

Training Data:

	precision	recall	f1-score	support
0	0.87	0.95	0.91	1710
1	0.95	0.85	0.90	1710
accuracy			0.90	3420
macro avg	0.91	0.90	0.90	3420
weighted avg	0.91	0.90	0.90	3420

Testing Data:

	precision	recall	f1-score	support
0	0.97	0.94	0.95	570
1	0.68	0.80	0.74	97
accuracy			0.92	667
macro avg	0.82	0.87	0.85	667
weighted avg	0.92	0.92	0.92	667

- The parameter tuning improved the Decision Tree performance a little.
- Overfitting observed



## Bagged Trees + Decision Tree Model

The bagging classifier is used to reduce variance in the dataset. Decision trees have low bias but high variance which can lead to overfitting and drastic output changes when minute input changes are made

```
In [53]: 1 # import Library
          2 from sklearn.ensemble import BaggingClassifier
          3
          4 bagged_tree = BaggingClassifier(DecisionTreeClassifier(criterion='gini',
          5                                                    n_estimators=20, random_state=1)
```

```
In [54]: 1 bagged_tree.fit(X_train, y_train)
```

```
Out[54]: BaggingClassifier(base_estimator=DecisionTreeClassifier(max_depth=4),
                           n_estimators=20, random_state=1)
```

```
In [55]: 1 bagged_tree.score(X_train, y_train)
```

```
Out[55]: 0.9524762381190596
```

```
In [56]: 1 bagged_tree.score(X_test, y_test)
```

```
Out[56]: 0.9400299850074962
```

```
In [57]: 1 y_pred = bagged_tree.predict(X_test)
```

```
In [58]: 1 print(confusion_matrix(y_test, y_pred))
          2 print(classification_report(y_test, y_pred))
```

```
[[562  8]
 [ 32 65]]

              precision    recall  f1-score   support

         0       0.95      0.99      0.97        570
         1       0.89      0.67      0.76         97

 accuracy          0.94          667
 macro avg       0.92      0.83      0.87          667
 weighted avg    0.94      0.94      0.94          667
```

```
In [59]: 1 recall_score(y_test, y_pred)
```

Out[59]: 0.6701030927835051

```
In [60]: 1 bagged_tree = BaggingClassifier(DecisionTreeClassifier(criterion='gini',
          2                               n_estimators=20, random_state=1))
```

```
In [61]: 1 bagged_tree.fit(X_train, y_train)
```

Out[61]: BaggingClassifier(base\_estimator=DecisionTreeClassifier(max\_depth=5),  
n\_estimators=20, random\_state=1)

```
In [62]: 1 bagged_tree.score(X_test, y_test)
```

Out[62]: 0.9430284857571214

```
In [63]: 1 y_pred = bagged_tree.predict(X_test)
```

```
In [64]: 1 print(confusion_matrix(y_test, y_pred))
          2 print(classification_report(y_test, y_pred))
```

```
[[563  7]
 [ 31 66]]
```

		precision	recall	f1-score	support
	0	0.95	0.99	0.97	570
	1	0.90	0.68	0.78	97
	accuracy			0.94	667
	macro avg	0.93	0.83	0.87	667
	weighted avg	0.94	0.94	0.94	667

```
In [65]: 1 recall_score(y_test, y_pred)
```

Out[65]: 0.6804123711340206

Random Forest Model's performance was not enhanced by parameter adjustment

Overfitting was noticed.

## Gradient Boost Model

To reduce to overall prediction error, Gradient Boost combines the prior models with the next best model that might be used. This classification model's prediction error gauges how accurately it forecast the variable of client churn.

```
In [66]: 1 #import Library
          2 from sklearn.ensemble import GradientBoostingClassifier
          3
          4 gbt_clf = GradientBoostingClassifier(random_state=1)
```

```
In [67]: 1 gbt_clf.fit(X_train, y_train)
```

```
Out[67]: GradientBoostingClassifier(random_state=1)
```

```
In [68]: 1 gbt_clf_train_preds = gbt_clf.predict(X_train)
2 gbt_clf_test_preds = gbt_clf.predict(X_test)
```

## Traning Score

```
In [69]: 1 print_metrics(y_train, gbt_clf_train_preds)
```

```
Precision Score: 1.0
Recall Score: 0.8477508650519031
Accuracy Score: 0.9779889944972486
F1 Score: 0.9176029962546817
```

```
In [70]: 1 gbt_classification_report = classification_report(y_test, gbt_clf_test_p
2 print(gbt_classification_report)
```

	precision	recall	f1-score	support
0	0.94	0.99	0.96	570
1	0.89	0.64	0.74	97
accuracy			0.94	667
macro avg	0.91	0.81	0.85	667
weighted avg	0.93	0.94	0.93	667

```
In [71]: 1 recall_score(y_test, y_pred)
```

```
Out[71]: 0.6804123711340206
```

## Parameter Tuning

```
In [72]: 1 print('Default parameters:')
        2 gbt_clf.get_params()
```

Default parameters:

```
Out[72]: {'ccp_alpha': 0.0,
          'criterion': 'friedman_mse',
          'init': None,
          'learning_rate': 0.1,
          'loss': 'deviance',
          'max_depth': 3,
          'max_features': None,
          'max_leaf_nodes': None,
          'min_impurity_decrease': 0.0,
          'min_impurity_split': None,
          'min_samples_leaf': 1,
          'min_samples_split': 2,
          'min_weight_fraction_leaf': 0.0,
          'n_estimators': 100,
          'n_iter_no_change': None,
          'presort': 'deprecated',
          'random_state': 1,
          'subsample': 1.0,
          'tol': 0.0001,
          'validation_fraction': 0.1,
          'verbose': 0,
          'warm_start': False}
```

## Adaboost Model

```
In [73]: 1 from sklearn.ensemble import AdaBoostClassifier
        2
        3 adaboost_clf = AdaBoostClassifier(random_state= 1)
        4 adaboost_clf.fit(X_train, y_train)
```

```
Out[73]: AdaBoostClassifier(random_state=1)
```

```
In [74]: 1 adaboost_train_preds = adaboost_clf.predict(X_train)
        2 adaboost_test_preds = adaboost_clf.predict(X_test)
```

## Training Score

```
In [75]: 1 print_metrics(y_train, adaboost_train_preds)
```

```
Precision Score: 0.7555555555555555
Recall Score: 0.47058823529411764
Accuracy Score: 0.9014507253626813
F1 Score: 0.5799573560767589
```

## Testing Score

```
In [76]: 1 recall_score(y_test, y_pred)
```

```
Out[76]: 0.6804123711340206
```

## Random Forest Model

```
In [77]: 1 from sklearn.ensemble import RandomForestClassifier
2
3 forest = RandomForestClassifier(n_estimators=20, max_depth= 11, random_s
4 forest.fit(X_train, y_train)
```

```
Out[77]: RandomForestClassifier(max_depth=11, n_estimators=20, random_state=1)
```

```
In [78]: 1 y_pred = forest.predict(X_test)
```

```
In [79]: 1 print(confusion_matrix(y_test, y_pred))
2 print(classification_report(y_test, y_pred))
```

```
[[566  4]
 [ 57 40]]
```

	precision	recall	f1-score	support
0	0.91	0.99	0.95	570
1	0.91	0.41	0.57	97
accuracy			0.91	667
macro avg	0.91	0.70	0.76	667
weighted avg	0.91	0.91	0.89	667

```
In [80]: 1 # training score  
2 forest.score(X_train, y_train)
```

Out[80]: 0.9719859929964982

```
In [81]: 1 # testing score  
2 forest.score(X_test, y_test)
```

Out[81]: 0.9085457271364318

```
In [82]: 1 recall_score(y_test, y_pred)
```

Out[82]: 0.41237113402061853

The parameter tuning didnt improve the performance of Random Forest model.

Overfitting observed.

## XGBoost Model

```
In [83]: 1 from xgboost import XGBClassifier, plot_importance
2
3 xg = XGBClassifier(random_state=1, eval_metric='logloss') #'logloss' is
4 xg.fit(X_train, y_train)
5
6 training_preds = xg.predict(X_train)
7 test_preds = xg.predict(X_test)
8
9 training_accuracy = accuracy_score(y_train, training_preds)
10 test_accuracy = accuracy_score(y_test, test_preds)
11 train_recall = recall_score(y_train, training_preds)
12 test_recall = recall_score(y_test, test_preds)
13
14 print('Training Accuracy: {:.4}%'.format(training_accuracy * 100))
15 print('Validation accuracy: {:.4}%'.format(test_accuracy * 100))
16 print('Training Recall: {:.4}%'.format(train_recall * 100))
17 print('Test Recall: {:.4}%'.format(test_recall * 100))
```

```
Training Accuracy: 100.0%
Validation accuracy: 94.6%
Training Recall: 100.0%
Test Recall: 67.01%
```

## Parameter Tuning



```
In [84]: 1 print('Default parameters:')
         2 xg.get_params()
```

Default parameters:

```
Out[84]: {'objective': 'binary:logistic',
          'use_label_encoder': True,
          'base_score': 0.5,
          'booster': 'gbtree',
          'colsample_bylevel': 1,
          'colsample_bynode': 1,
          'colsample_bytree': 1,
          'enable_categorical': False,
          'gamma': 0,
          'gpu_id': -1,
          'importance_type': None,
          'interaction_constraints': '',
          'learning_rate': 0.300000012,
          'max_delta_step': 0,
          'max_depth': 6,
          'min_child_weight': 1,
          'missing': nan,
          'monotone_constraints': '()',
          'n_estimators': 100,
          'n_jobs': 8,
          'num_parallel_tree': 1,
          'predictor': 'auto',
          'random_state': 1,
          'reg_alpha': 0,
          'reg_lambda': 1,
          'scale_pos_weight': 1,
          'subsample': 1,
          'tree_method': 'exact',
          'validate_parameters': 1,
          'verbosity': None,
          'eval_metric': 'logloss'}
```

```

In [85]: 1 # Tuning XGBClassifier with GridSearchCV
2 # Takes more than 10 minutes on my PC
3
4 from sklearn.model_selection import GridSearchCV
5
6 xgb_param_grid = {
7     'learning_rate': [0.1, 0.2],
8     'max_depth': [2, 3, 4, 5, 6],
9     'min_child_weight': [1, 2],
10    'subsample': [0.5, 0.7],
11    'n_estimators': [30, 100],
12 }
13
14 xgb_gs = GridSearchCV(xg, xgb_param_grid, cv=5, scoring='f1')
15 xgb_gs.fit(X_train_scaled_resampled, y_train_resampled)
16
17 score_xgb_gs = xgb_gs.score(X_test_scaled, y_test)
18 print('f1-score on test data:', score_xgb_gs)
19
20 print('Parameter Tuning Results:\n')
21 print("Best Parameter Combination:", xgb_gs.best_params_)
22 print('Training Data:\n', classification_report(y_train_resampled, xgb_g
23 print('Testing Data:\n', classification_report(y_test, xgb_gs.predict(X_

```

f1-score on test data: 0.7701149425287357

Parameter Tuning Results:

Best Parameter Combination: {'learning\_rate': 0.1, 'max\_depth': 6, 'min\_child\_weight': 1, 'n\_estimators': 100, 'subsample': 0.5}

Training Data:

	precision	recall	f1-score	support
0	0.99	1.00	0.99	1710
1	1.00	0.98	0.99	1710
accuracy			0.99	3420
macro avg	0.99	0.99	0.99	3420
weighted avg	0.99	0.99	0.99	3420

Testing Data:

	precision	recall	f1-score	support
0	0.95	0.98	0.97	570
1	0.97	0.99	0.97	570

## Tuning XGBoost Model with GridSearchCV

```
In [86]: 1 from sklearn.model_selection import cross_val_score, RandomizedSearchCV
2
3 param_grid = {
4     'learning_rate': [0.1, 0.2],
5     'max_depth': [6],
6     'min_child_weight': [1, 2],
7     'subsample': [0.5, 0.7],
8     'n_estimators': [100],
9 }
10
```

```
In [87]: 1 grid_xg = GridSearchCV(xg, param_grid, scoring='accuracy', cv=None, n_jo
2 grid_xg.fit(X_train, y_train)
3
4 best_parameters = grid_xg.best_params_
5
6 print('Grid Search found the following optimal parameters: ')
7 for param_name in sorted(best_parameters.keys()):
8     print('%s: %r' % (param_name, best_parameters[param_name]))
9
10 training_preds = grid_xg.predict(X_train)
11 test_preds = grid_xg.predict(X_test)
12 training_accuracy = accuracy_score(y_train, training_preds)
13 test_accuracy = accuracy_score(y_test, test_preds)
14 train_recall = recall_score(y_train, training_preds)
15 test_recall = recall_score(y_test, test_preds)
16
17 print('')
18 print('Training Accuracy: {:.4}%'.format(training_accuracy * 100))
19 print('Validation accuracy: {:.4}%'.format(test_accuracy * 100))
20 print('Training Recall: {:.4}%'.format(train_recall * 100))
21 print('Test Recall: {:.4}%'.format(test_recall * 100))
22
```

Grid Search found the following optimal parameters:

```
learning_rate: 0.1
max_depth: 6
min_child_weight: 1
n_estimators: 100
subsample: 0.7
```

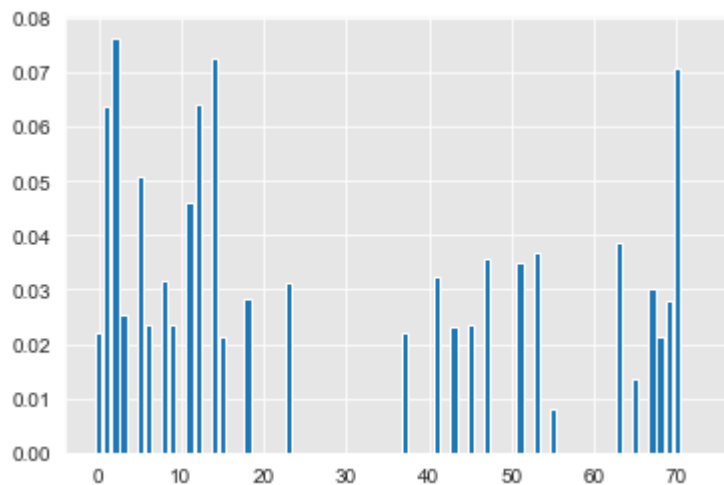
```
Training Accuracy: 99.0%
Validation accuracy: 94.6%
Training Recall: 93.08%
Test Recall: 68.04%
```

## XGBoost with Optimal Parameters

```
In [88]: 1 xg = XGBClassifier(max_depth = 6, learning_rate = .1, n_estimators = 100
2
3 xg.fit(X_train, y_train)
4
5 training_preds = xg.predict(X_train)
6 test_preds = xg.predict(X_test)
7
8 training_accuracy = accuracy_score(y_train, training_preds)
9 test_accuracy = accuracy_score(y_test, test_preds)
10 train_recall = recall_score(y_train, training_preds)
11 test_recall = recall_score(y_test, test_preds)
12
13 print('Training Accuracy: {:.4}%'.format(training_accuracy * 100))
14 print('Validation accuracy: {:.4}%'.format(test_accuracy * 100))
15 print('Training Recall: {:.4}%'.format(train_recall * 100))
16 print('Test Recall: {:.4}%'.format(test_recall * 100))
```

Training Accuracy: 99.0%  
Validation accuracy: 94.6%  
Training Recall: 93.08%  
Test Recall: 68.04%

```
In [89]: 1 plt.bar(range(len(xg.feature_importances_)), xg.feature_importances_);
2 #plt.savefig("images/xgboost.png")
```



## Best Performing Models

Since decision trees can be simpler to read, I have decided to stick with for the rest of my investigation.

```
In [90]: 1 dt = DecisionTreeClassifier(criterion='entropy', random_state=1)
          2 dt.fit(X_train, y_train)
```

```
Out[90]: DecisionTreeClassifier(criterion='entropy', random_state=1)
```

```
In [91]: 1 # Make predictions using test set
          2 y_pred = dt.predict(X_test)
          3
          4 # Check the AUC of predictions
          5 false_positive_rate, true_positive_rate, thresholds = roc_curve(y_test, y.
          6 roc_auc = auc(false_positive_rate, true_positive_rate)
          7 roc_auc
```

```
Out[91]: 0.834400434074878
```

```
In [92]: 1 recall_score(y_test, y_pred)
```

```
Out[92]: 0.7319587628865979
```

## Compare the Models

In this part, I will contrast the available categorization models in order to determine which is the most effective at identifying potential consumers for SyriaTel.

I will now consider evaluation metrics like f1, recall, accuracy and precision.

For each model, I will also calculate AUC and plot ROC curves.

## Optimum parameter sets, with f1-score used for tuning

- Logistic Regression: {'C': 0.01, 'solver': 'liblinear'}
- KNN: Default
- Decision Trees: {'criterion': 'gini', 'max\_depth': 6, 'min\_samples\_split': 2}
- Bagging classifier: {DecisionTreeClassifier {'criterion': 'gini', 'max\_depth': 5}, 'n\_estimators': 20}
- Gradient Boost: Default
- Adaboost: Default
- Random Forest: Default

- XGBoost: {'learning\_rate': 0.1, 'max\_depth': 6, 'min\_child\_weight': 1, 'n\_estimators': 100}

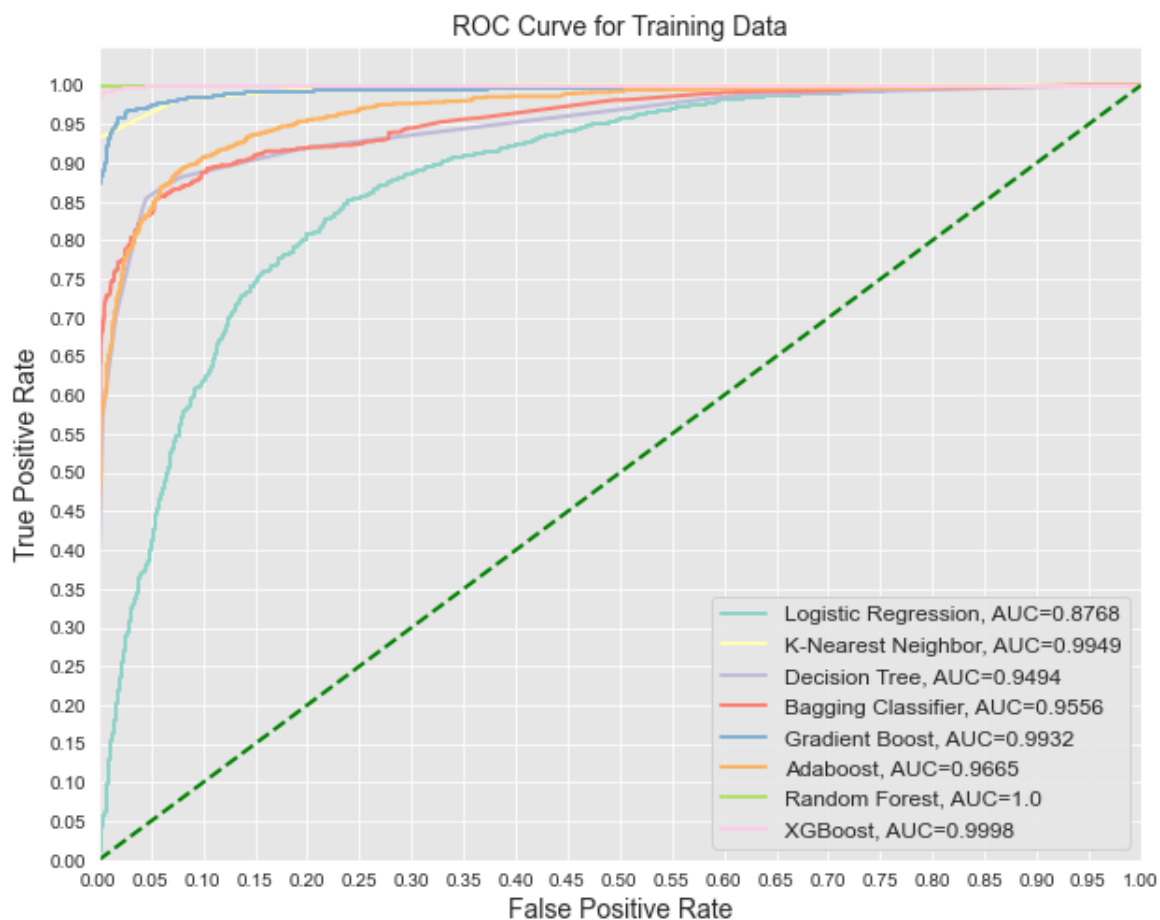
```
In [93]: ▶ 1 # Instantiate models with optimum parameters
2
3 logreg_best = LogisticRegression(C=0.01, solver='liblinear', random_state=42)
4 knn_best = KNeighborsClassifier()
5 dt_best = DecisionTreeClassifier(criterion='gini', max_depth=6, min_samples_split=2,
6 bt_best = BaggingClassifier(DecisionTreeClassifier(criterion='gini', max_depth=6,
7 n_estimators=20, random_state=42)
8 gbt_best = GradientBoostingClassifier(random_state=42)
9 ada_best = AdaBoostClassifier(random_state=42)
10 rf_best = RandomForestClassifier(random_state=42)
11 xgb_best = XGBClassifier(learning_rate=0.1, max_depth=6, min_child_weight=1,
12 random_state=42, eval_metric='logloss')
13
14 total_list = [logreg_best, knn_best, dt_best, bt_best, gbt_best, ada_best, rf_best, xgb_best]
15 model_names = ['Logistic Regression', 'K-Nearest Neighbor', 'Decision Tree', 'Random Forest',
16 'Gradient Boost', 'Adaboost', 'Random Forest', 'XGBoost']
```

```

In [94]: 1 def total_scores(dataset_type, X_scaled, y_true):
2         """
3         dataset_type = 'Testing' or 'Training'
4         X_scaled = X_test_scaled or X_train_scaled
5         y_true = y_train or y_test
6
7         """
8         colors = sns.color_palette('Set3')
9         plt.figure(figsize=(10, 8))
10
11         total_scores_list = []
12
13         for n, clf in enumerate(total_list):
14
15             clf.fit(X_train_scaled_resampled, y_train_resampled)
16
17             y_pred = clf.predict(X_scaled)
18
19             y_prob = clf.predict_proba(X_scaled)
20             fpr, tpr, thresholds = roc_curve(y_true, y_prob[:,1])
21             auc_score = auc(fpr, tpr)
22             plt.plot(fpr, tpr, color=colors[n], lw=2, label=f'{model_names[n]}')
23
24             fit_scores = {'model': model_names[n],
25                           'precision': round(precision_score(y_true, y_pred),3),
26                           'recall': round(recall_score(y_true, y_pred),3),
27                           'accuracy': round(accuracy_score(y_true, y_pred),3),
28                           'f1': round(f1_score(y_true, y_pred),3),
29                           'auc': round(auc_score,3)
30                           }
31
32             total_scores_list.append(fit_scores)
33
34             plt.plot([0, 1], [0, 1], color='green', lw=2, linestyle='--')
35             plt.xlim([0.0, 1.0])
36             plt.ylim([0.0, 1.05])
37             plt.yticks([i/20.0 for i in range(21)])
38             plt.xticks([i/20.0 for i in range(21)])
39             plt.xlabel('False Positive Rate', fontsize=14)
40             plt.ylabel('True Positive Rate', fontsize=14)
41             plt.title(f'ROC Curve for {dataset_type} Data', fontsize=14)
42             plt.legend(loc='lower right', fontsize=12)
43             plt.savefig(f'images/ROC_Curve_{dataset_type}.png')
44             plt.show()
45
46             total_scores_df = pd.DataFrame(total_scores_list)
47             total_scores_df = total_scores_df.set_index('model')
48
49
50         return total_scores_df

```

```
In [95]: 1 total_scores('Training', X_train_scaled_resampled, y_train_resampled)
```

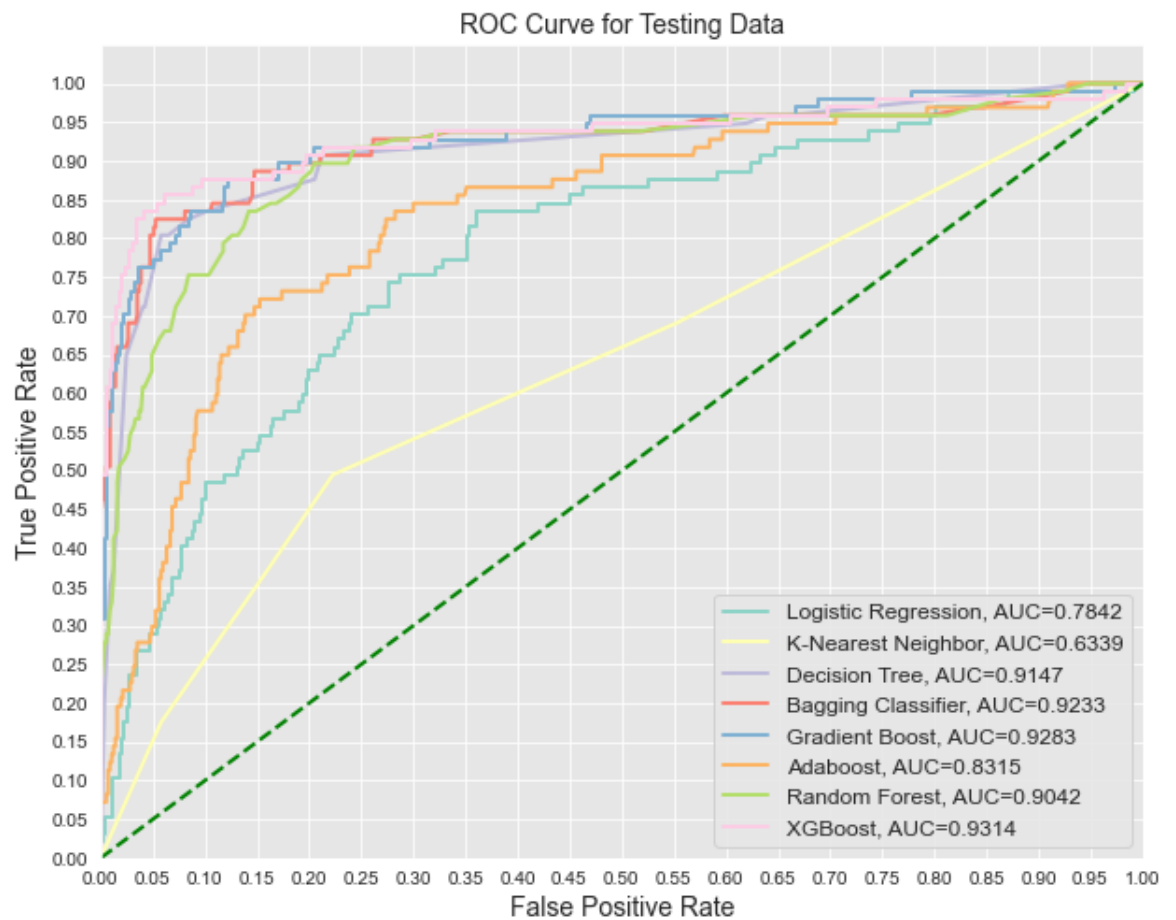


Out[95]:

	precision	recall	accuracy	f1	auc
model					
<b>Logistic Regression</b>	0.783	0.842	0.804	0.811	0.877
<b>K-Nearest Neighbor</b>	0.837	0.995	0.901	0.909	0.995
<b>Decision Tree</b>	0.949	0.854	0.904	0.899	0.949
<b>Bagging Classifier</b>	0.940	0.853	0.899	0.895	0.956
<b>Gradient Boost</b>	0.978	0.958	0.968	0.968	0.993
<b>Adaboost</b>	0.917	0.892	0.906	0.904	0.967
<b>Random Forest</b>	1.000	1.000	1.000	1.000	1.000
<b>XGBoost</b>	0.999	0.984	0.992	0.991	1.000



In [96]: 1 total\_scores('Testing', X\_test\_scaled, y\_test)



Out[96]:

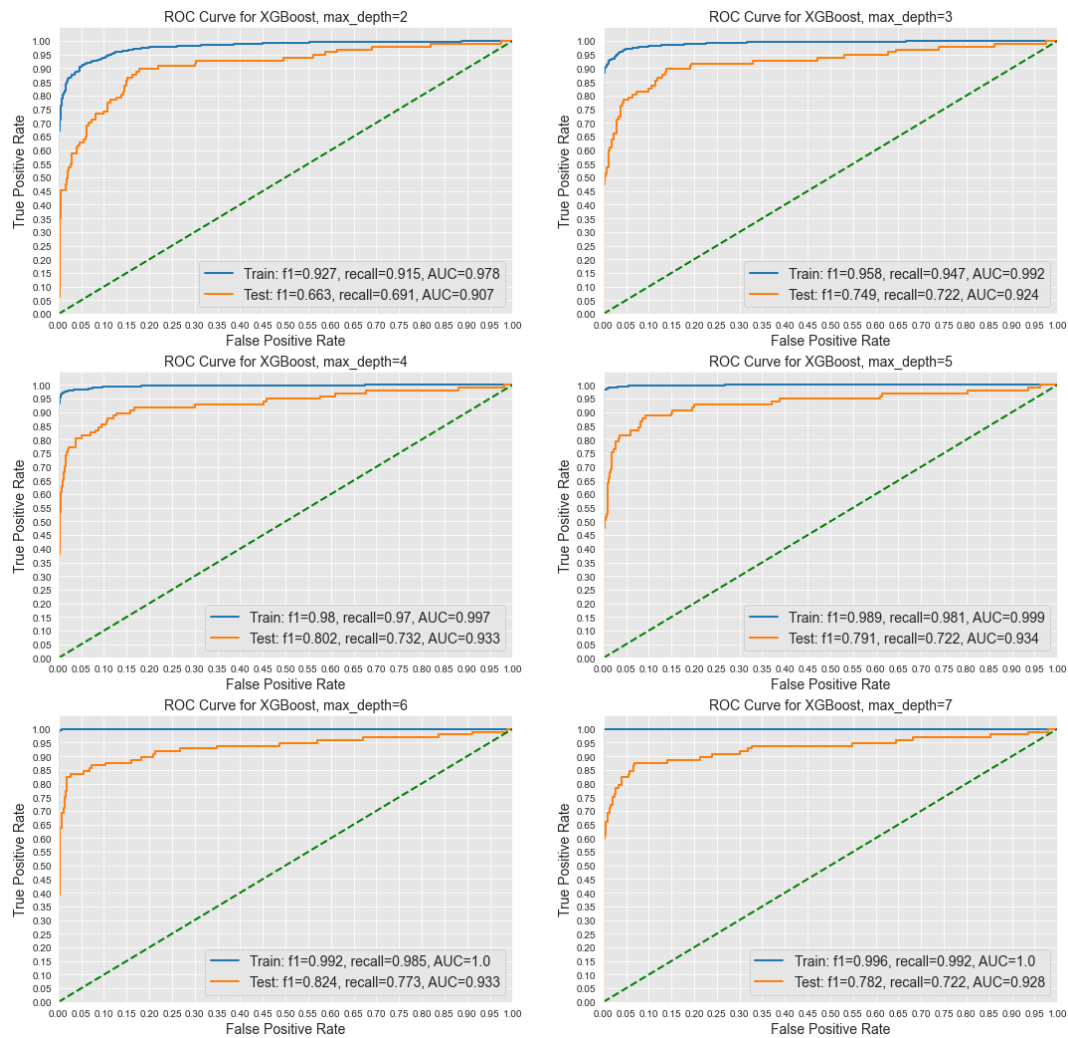
	precision	recall	accuracy	f1	auc
model					
<b>Logistic Regression</b>	0.328	0.680	0.751	0.443	0.784
<b>K-Nearest Neighbor</b>	0.274	0.495	0.736	0.353	0.634
<b>Decision Tree</b>	0.678	0.804	0.916	0.736	0.915
<b>Bagging Classifier</b>	0.690	0.825	0.921	0.751	0.923
<b>Gradient Boost</b>	0.777	0.753	0.933	0.764	0.928
<b>Adaboost</b>	0.514	0.567	0.859	0.539	0.831
<b>Random Forest</b>	0.720	0.608	0.909	0.659	0.904
<b>XGBoost</b>	0.877	0.732	0.946	0.798	0.931

## Final Model

```

In [97]: 1 fig, axes = plt.subplots(3, 2, figsize=(18, 18))
2
3 depths = [2, 3, 4, 5, 6, 7]
4
5 for ax, d in zip(axes.flat, depths):
6
7     xgb_clf = XGBClassifier(learning_rate=0.1, max_depth=d, min_child_weight=1,
8                             subsample=0.7, random_state=42, eval_metric='auc')
9     xgb_clf.fit(X_train_scaled_resampled, y_train_resampled)
10
11     y_train_pred = xgb_clf.predict(X_train_scaled_resampled)
12     y_train_prob = xgb_clf.predict_proba(X_train_scaled_resampled) # each class probability
13     fpr_train, tpr_train, threshold_train = roc_curve(y_train_resampled, y_train_prob[:,1])
14     auc_train = round(auc(fpr_train, tpr_train),3)
15     f1_train = round(f1_score(y_train_resampled, y_train_pred),3)
16     recall_train = round(recall_score(y_train_resampled, y_train_pred),3)
17     ax.plot(fpr_train, tpr_train, lw=2, label=f'Train: f1={f1_train}, recall={recall_train}')
18
19     y_test_pred = xgb_clf.predict(X_test_scaled)
20     y_test_prob = xgb_clf.predict_proba(X_test_scaled) # each class probability
21     fpr_test, tpr_test, thresholds_test = roc_curve(y_test, y_test_prob[:,1])
22     auc_test = round(auc(fpr_test, tpr_test),3)
23     f1_test = round(f1_score(y_test, y_test_pred),3)
24     recall_test = round(recall_score(y_test, y_test_pred),3)
25     ax.plot(fpr_test, tpr_test, lw=2, label=f'Test: f1={f1_test}, recall={recall_test}')
26
27     ax.plot([0, 1], [0, 1], color='green', lw=2, linestyle='--')
28     ax.set_xlim([0.0, 1.0])
29     ax.set_ylim([0.0, 1.05])
30     ax.set_yticks([i/20.0 for i in range(21)])
31     ax.set_xticks([i/20.0 for i in range(21)])
32     ax.set_xlabel('False Positive Rate', fontsize=14)
33     ax.set_ylabel('True Positive Rate', fontsize=14)
34     ax.set_title(f'ROC Curve for XGBoost, max_depth={d}', fontsize=14)
35     ax.legend(loc='lower right', fontsize=14)
36
37 #plt.savefig('images/Roc_curve_XGB_maxd.png')

```



```

In [98]: 1 # Fit Final Model
2
3 xgb_final = XGBClassifier(learning_rate=0.1, max_depth=5, min_child_weight
4                      subsample=0.7, random_state=42, eval_metric='logit')
5
6 xgb_final.fit(X_train_scaled_resampled, y_train_resampled)
7
8 print('Final Model:\n')
9 print('Training Data:\n', classification_report(y_train_resampled, xgb_f
10 print('Testing Data:\n', classification_report(y_test, xgb_final.predict

```

Final Model:

Training Data:

	precision	recall	f1-score	support
0	0.98	1.00	0.99	1710
1	1.00	0.98	0.99	1710
accuracy			0.99	3420
macro avg	0.99	0.99	0.99	3420
weighted avg	0.99	0.99	0.99	3420

Testing Data:

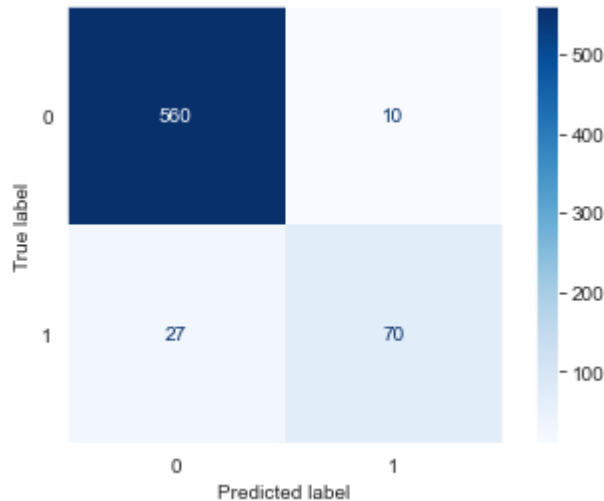
	precision	recall	f1-score	support
0	0.95	0.98	0.97	570
1	0.88	0.72	0.79	97
accuracy			0.94	667
macro avg	0.91	0.85	0.88	667
weighted avg	0.94	0.94	0.94	667

Which model is the best identifying churn customer?

Overall, XGBosst classifier has the best performance, according to test data Evaluation metrics. There is also shown the best Recall and F1-score.

My choice of the best model is XGBoost model.

```
In [117]: 1 from sklearn.metrics import confusion_matrix, plot_confusion_matrix
2
3 plot_confusion_matrix(xgb_final, X_test_scaled, y_test, cmap=plt.cm.Blue
4 plt.grid(False)
5 plt.savefig('images/conf_matrix_XGB.png')
```



## Summary:

### Statistics of the Final Model in concise:

It clearly identifies 72% of the real churning customer. 87% of the customers whose anticipated churn was captured by the algorithm definitely did so (clearly remember).(accuracy)The f1-score's Harmonic Mean of Precision and Recall is 79%.

The experimental database's identification number are:

### Unique identifiers:

70 confirmed positives were found

There are 560 genuine negatives.

10 false alarms were discovered.

27 erroneous alarms were discarded.

70 out of 125 customers who churn are successfully identified.

- Client that churn have higher probability of having a foreign plan than others who stay users.
- Compared with regular users, churn customers are reluctant to get a voicemail subscription.
- Contrary to incumbent users, churn clients have fewer voicemails(as a result of less voicemail plan).
- Users with churn generate more queries to customer service than do customer base.
- In contrast to continuous customers, churn users have greater total dat minutes.

## Business Recommendations:

Improve international plan to attract customers.

For greater satisfaction, revamp its helpdesk(customer service).

Accept a deal at discount with enough cumulative day moments.

## Next Step

Ensure smooth functioning of the XGBT design(completed model).

To understand how parameter influences the performance, browse for it properly

To facilitate a better understanding and familiarity of each parameter exploited in grid search.

Analyze the influence of additional hyperparameter.

To evaluate performance of the model and to alter parameters, use a scaled f1 score that emphasizes recall more accuracy