

## 1 Final Project Submission

- Student name: Sam Oliver
- Student pace: self paced
- Scheduled project review date/time: 7/15/2022 (11:00am EST)
- Instructor name: Abhineet Kulkarni
- Blog post URL: <https://medium.com/@samueloliver303/exploring-how-random-forests-work-in-machine-learning-7515015a6aa8> (<https://medium.com/@samueloliver303/exploring-how-random-forests-work-in-machine-learning-7515015a6aa8>)

## 2 Overview

This notebook represents efforts to predict trends in the telecommunication data that was investigated in the EDA notebook. For initial context, it is recommended to first read the README file. It is helpful to scan through the EDA first before this notebook as well to understand some of the strategies implemented in this notebook.

The problem: SyriaTel (a telecommunications company in Syria) is using a US-based telecom dataset surrounding the idea of customer churn in order to decrease losses involved in losing customers. Main considerations for this problem:

- Business consideration: Identify customers that are likely to churn and implement a strategy, such as discounted rates for x number of months, to retain these customers.
- Evaluation consideration: More money will be lost on a False Negative (customer that is predicted to be retained but actually churned) than on a False Positive (customer that is predicted to churn but actually is retained) because it is assumed that it is more costly to replace the False Negative with a new customer and lose the stream of income from this churned customer than it is to implement a retention strategy (i.e. discounted rate) on a customer that does not need it and simply lose the difference between what they would have paid versus the discounted payment stream. For this reason, an F-2 Score will be used to evaluate model success because this metric penalizes False Negatives more than False Positives.
- It should also be noted that F-2 Score is used instead of another metric like recall because F-2 will penalize both FPs and FNs, which is important because both of these cases negatively affect the company. F-2 Score also allows for some buffer in assumptions made for how to assign a specific cost to FPs and FN's, which will be evaluated later in this notebook.

The F-Score Formula is below. To find F-2, Beta is simply replaced with 2:

$$F_{\beta} = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{(\beta^2 \cdot \text{precision}) + \text{recall}}$$

In terms of type I and type II errors, the formula is as follows:

$$F_{\beta} = \frac{(1 + \beta^2) \cdot \text{true positive}}{(1 + \beta^2) \cdot \text{true positive} + \beta^2 \cdot \text{false negative} + \text{false positive}}$$

Beta is squared, which indicates that if it is lower than 1, then it penalizes true positives more than false negatives. If Beta is greater than 1, then false negatives are penalized more.

F-2 is applicable as an evaluation metric for this project because it penalizes false negatives significantly more than true positives, and it also allows for some room for error to be made with assumptions for how much exactly a False Positive or False Negative costs the company. If something else like ROC-AUC curves were used to optimize the rates of FP's and FN's, then the calculations for cost of these categories would have to be much more specific.

## 3 Modeling Preparation

### 3.1 Imports and Some Dtype Conversions

```
In [1]: 1 # imports
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 import sklearn
7
8 %matplotlib inline
9
10 from sklearn.model_selection import train_test_split
11 from sklearn.impute import SimpleImputer
12 from sklearn.preprocessing import OneHotEncoder
13 from sklearn.preprocessing import MinMaxScaler
14 from sklearn.linear_model import LogisticRegression
15
16 from imblearn.over_sampling import SMOTE
17
18 from sklearn.preprocessing import StandardScaler
19 from sklearn.neighbors import KNeighborsClassifier
20 from sklearn.model_selection import train_test_split, GridSearchCV
21 from sklearn.ensemble import RandomForestClassifier
22 from sklearn.pipeline import Pipeline
23 from sklearn import tree
24 from sklearn.svm import SVC
25 from sklearn import svm
26
27 from sklearn.metrics import classification_report, confusion_matrix
28 from sklearn.metrics import fbeta_score, accuracy_score, make_scorer
29 from sklearn.metrics import plot_confusion_matrix
```

```

In [2]: 1 # import dataset
2 df = pd.read_csv('../data\\telecom_data.csv')
3
4 # Handle object types for international plan and voice mail plan
5 df.loc[df['international plan'] == 'no', 'international plan'] = 0
6 df.loc[df['international plan'] == 'yes', 'international plan'] = 1
7
8 df.loc[df['voice mail plan'] == 'no', 'voice mail plan'] = 0
9 df.loc[df['voice mail plan'] == 'yes', 'voice mail plan'] = 1
10
11 # Change churn to values: 1 (churned/True) 0 (no churn/False)
12 df.loc[df['churn'] == True, 'churn'] = 1
13 df.loc[df['churn'] == False, 'churn'] = 0
14
15 # going to create backup df and drop phone number from original df
16 # phone number could be used as unique id, but it doesn't seem necessary
17 df_backup = df.copy()
18 df = df.drop(['phone number'], axis=1)
19
20 # casting int values to churn, voice mail plan, and international plan cols
21 objs = ['international plan', 'voice mail plan', 'churn']
22
23 for o in objs:
24     df = df.astype({o: int})
25
26 # dropping area code
27 df = df.drop(['area code'], axis=1)
28
29 # dropping various columns due to perfect (or exceptionally high) correlation
30 df = df.drop(['total intl minutes', 'number vmail messages'], axis=1)
31
32 # check df
33 df.head()

```

Out[2]:

	state	account length	international plan	voice mail plan	total day minutes	total day calls	total day charge	total eve minutes	total eve calls	total eve charge	total night minutes	t
0	KS	128	0	1	265.1	110	45.07	197.4	99	16.78	244.7	
1	OH	107	0	1	161.6	123	27.47	195.5	103	16.62	254.4	
2	NJ	137	0	0	243.4	114	41.38	121.2	110	10.30	162.6	
3	OH	84	1	0	299.4	71	50.90	61.9	88	5.26	196.9	
4	OK	75	1	0	166.7	113	28.34	148.3	122	12.61	186.9	

Description of the variables in the dataset:

Feature	Description
State	Which state in the US the customer is located
Account length	Number of months customer has been using the service
Area code	Customer's area code

Feature	Description
Phone number	Customer's phone number without area code
International plan	(Y/N) does customer have an international plan?
Voice mail plan	(Y/N) does the customer have a voice mail plan?
Number vmail messages	Number of voice mail messages
Total day minutes	Total day minutes used
Total day calls	Calls made during the day
Total day charge	Total charge for day calls
Total eve minutes	Total evening minutes used
Total eve calls	Calls made during the evening
Total eve charge	Total charge for evening calls
Total night minutes	Total night minutes used
Total night calls	Calls made during the night
Total night charge	Calls made during the night
Total intl minutes	Total minutes for international calls
Total intl calls	Total number of international calls
Total intl charge	Total charges for international calls
Customer service calls	Total number of calls to customer service
Churn	(T/F) Did customer churn (cancel service) or not?

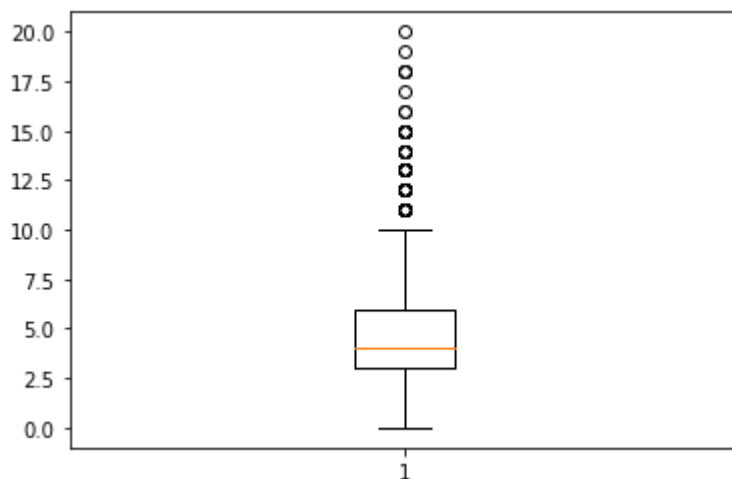
Handle international calls to bin them into categories easier.

```
In [3]: 1 # handle international calls to bin them into categories easier.
        2 df['total intl calls'].value_counts()
```

```
Out[3]: 3    668
        4    619
        2    489
        5    472
        6    336
        7    218
        1    160
        8    116
        9    109
       10     50
       11     28
        0     18
       12     15
       13     14
       15      7
       14      6
       18      3
       16      2
       19      1
       17      1
       20      1
Name: total intl calls, dtype: int64
```

```
In [4]: 1 plt.boxplot(df['total intl calls'])
```

```
Out[4]: {'whiskers': [<matplotlib.lines.Line2D at 0x1dd385ce4c0>,
<matplotlib.lines.Line2D at 0x1dd385ce820>],
'caps': [<matplotlib.lines.Line2D at 0x1dd385ceb80>,
<matplotlib.lines.Line2D at 0x1dd385ceee0>],
'boxes': [<matplotlib.lines.Line2D at 0x1dd385ce160>],
'medians': [<matplotlib.lines.Line2D at 0x1dd385e8280>],
'fliers': [<matplotlib.lines.Line2D at 0x1dd385e85e0>],
'means': []}
```



A few outliers in total international calls, but most customers have made around 3-6 calls.

In [5]: 1 df['total intl calls'].describe()

```
Out[5]: count      3333.000000
mean         4.479448
std          2.461214
min          0.000000
25%          3.000000
50%          4.000000
75%          6.000000
max          20.000000
Name: total intl calls, dtype: float64
```

Treating total international calls as categorical because it can be treated as such in the dataset.  
Binning into low, moderate, and high international call frequency for simplicity.

```
In [6]: 1 # range is 0-20 for international calls with most concentrated from 0-10
2 # I will bin into cats: low, moderate, and high with values <3, 3-6, and >6
3
4 list_tmp = []
5
6 for index, row in df.iterrows():
7     if row['total intl calls'] < 3:
8         list_tmp.append('low')
9     elif row['total intl calls'] > 6:
10        list_tmp.append('high')
11    else:
12        list_tmp.append('moderate')
13
14 df['total_intl_calls'] = list_tmp
15
16 df['total_intl_calls'].describe()
```

```
Out[6]: count      3333
unique         3
top      moderate
freq      2095
Name: total_intl_calls, dtype: object
```

## 3.2 What would be the % chance of guessing correctly if the customer was assumed to not churn?

In [7]: 1 1-df.churn.mean()

```
Out[7]: 0.8550855085508551
```

Guessing "not churned" for every customer would yield about an 85.5% chance of guessing correctly. Notice the class imbalance as well.

## 3.3 Handle datatypes

Handle state and total\_intl\_calls object types. Turn these into integers for later modeling.

```
In [8]: 1 states = ["AL", "AK", "AZ", "AR", "CA", "CO", "CT", "DC", "DE", "FL", "GA",
2           "HI", "ID", "IL", "IN", "IA", "KS", "KY", "LA", "ME", "MD", "MA",
3           "MI", "MN", "MS", "MO", "MT", "NE", "NV", "NH", "NJ", "NM", "NY",
4           "NC", "ND", "OH", "OK", "OR", "PA", "RI", "SC", "SD", "TN", "TX",
5           "UT", "VT", "VA", "WA", "WV", "WI", "WY"]
```

```
In [9]: 1 state_int = []
2 for i, row in df.iterrows():
3     state_int.append(states.index(row['state']))
```

```
In [10]: 1 df['state_int'] = state_int
```

States in the dataset are alphabetically ordered such that Alabama is 0 and Wyoming is 50.

```
In [11]: 1 intl_calls_int = []
2 for i, row in df.iterrows():
3     if row['total_intl_calls'] == 'low':
4         intl_calls_int.append(0)
5     elif row['total_intl_calls'] == 'moderate':
6         intl_calls_int.append(1)
7     else:
8         intl_calls_int.append(2)
9
10 df['intl_calls_bins'] = intl_calls_int
```

Total international calls is now split into 0, 1, and 2 instead of low, moderate, and high. The new columns correspond to the following: 0 represents customers with less than 3 international calls, 1 represents customers with between 3 and 6 international calls, and 2 represents customers with over 6 international calls.

### 3.4 Feature Engineering: Totals for Charge, Minutes, and Calls

To reduce the number of variables and redundancies in the dataset, total calls, total charge, and total minutes will represent totals for each category. These new categories will not include international call data.

```
In [12]: 1 df['total_calls'] = df.apply(lambda x: x['total night calls'] + x['total day calls'],
2                                axis=1)
3 df['total_charge'] = df.apply(lambda x: x['total day charge'] + x['total evening charge'],
4                                axis=1)
5 df['total_minutes'] = df.apply(lambda x: x['total day minutes'] + x['total evening minutes'],
6                                axis=1)
```

It was observed from the EDA that total minutes and total charge are perfectly correlated. Total minutes will be dropped from the dataframe.

```
In [13]: 1 df = df.drop(['total_minutes'], axis=1)
```

### 3.5 Define X and y and Normalize the Data

```
In [14]: 1 y = df['churn']
2
3 # drop redundancies now that there are total charge, minutes, and calls col
4 # dropping any other columns that have been replaced with new ones.
5
6 X = df.drop(['churn', 'total intl calls',
7             'total_intl_calls', 'state', 'total day minutes',
8             'total day calls', 'total day charge', 'total eve minutes',
9             'total eve calls', 'total eve charge', 'total night minutes',
10            'total night calls', 'total night charge'], axis=1)
```

### 3.6 Train-test-split and SMOTE

The target feature has imbalance because customers that churned only represent about 14.5% of the data. Synthetic Minority Oversampling Technique (SMOTE) will be implemented to create balance in the target feature such that the models do not overlook the minority class (churned customers).

```
In [15]: 1 # Split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=27)
3
4 # Fixing class imbalance with SMOTE
5 smote = SMOTE(random_state=42)
6 X_train_resampled, y_train_resampled = smote.fit_sample(X_train, y_train)
```

### 3.7 OHE and Normalize

Avoiding data leakage by transforming after the split. Also avoiding data leakage by fitting OHE and Scaler on training data and then transforming both OHE and Scaler using the training data fit.

```
In [16]: 1 # Define categorical columns for OHE
2 cats = ['state_int', 'international plan', 'intl_calls_bins',
3         'customer service calls', 'voice mail plan']
```



```

In [17]: 1 # Seperate data into categorical for train and test sets
2 X_train_cats = X_train_resampled[cats]
3 X_test_cats = X_test[cats]
4
5 # handle categorical values
6 ohe = OneHotEncoder(handle_unknown="ignore", sparse=False)
7
8 # OHE for training categoricals
9 ohe.fit(X_train_cats)
10 X_train_ohe = pd.DataFrame(
11     ohe.transform(X_train_cats),
12     index=X_train_cats.index,
13     columns=ohe.get_feature_names()
14 )
15
16 # OHE for testing categoricals - transform with train fit
17 X_test_ohe = pd.DataFrame(
18     ohe.transform(X_test_cats),
19     index=X_test_cats.index,
20     columns=ohe.get_feature_names()
21 )
22
23 # Scaling variables to work well with OHE data -scale train and test data
24 X_train_numerics = X_train_resampled.drop(cats, axis=1)
25 X_test_numerics = X_test.drop(cats, axis=1)
26
27 scaler = MinMaxScaler()
28
29 scaler.fit(X_train_numerics)
30 X_train_scaled = pd.DataFrame(
31     scaler.transform(X_train_numerics),
32     index=X_train_numerics.index,
33     columns=X_train_numerics.columns
34 )
35
36 # Scale variables of test data with the train fit (avoids data Leakage)
37 X_test_scaled = pd.DataFrame(
38     scaler.transform(X_test_numerics),
39     index=X_test_numerics.index,
40     columns=X_test_numerics.columns
41 )
42
43
44 # Concatenate and replace X_train and X_test with OHE+scaled data
45 X_train_resampled = pd.DataFrame()
46 X_train_resampled = pd.concat([X_train_scaled, X_train_ohe], axis=1)
47
48 X_test = pd.DataFrame()
49 X_test = pd.concat([X_test_scaled, X_test_ohe], axis=1)

```

Printing a few rows of the training data below. Explanation of variable names is below this output.

In [18]: 1 X\_train\_resampled.head()

Out[18]:

	account length	total intl charge	total_calls	total_charge	x0_0	x0_1	x0_2	x0_3	x0_4	x0_5	...	x3_2	x3_
0	0.223214	0.479630	0.779279	0.794778	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0
1	0.982143	0.120370	0.414414	0.550741	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0
2	0.486607	0.585185	0.612613	0.382357	0.0	0.0	0.0	0.0	0.0	0.0	...	1.0	0
3	0.611607	0.550000	0.211712	0.722089	0.0	0.0	0.0	0.0	0.0	0.0	...	1.0	0
4	0.352679	0.655556	0.436937	0.144954	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0

5 rows × 72 columns

To clarify the new confusing column names ('x0\_0' and so on), the following list is provided to retain what each column refers to:

- First 51 columns (x0\_0-x0\_50) refer to the states
- The next two columns (x1\_0 & x1\_1) refer to whether or not the customer has an international plan.
- Columns x2\_0-x2\_2 refer to the frequency (low, moderate, high) of international calls.
- Columns x3\_0-x3\_9 refer to how many calls to customer service the customer has completed.
- The last two columns indicate whether or not the customer has a voice mail plan.



## 3.8 Confusion Matrix Helper Function

This function return a simple, string of confusion matrix values.

```

In [19]: 1 # This function prints out a string of totals for TP, TN, FP, & FN
          2 def conf_matrix(y_true, y_pred):
          3     cm = {'TP': 0, 'TN': 0, 'FP': 0, 'FN': 0}
          4
          5     for ind, label in enumerate(y_true):
          6         pred = y_pred[ind]
          7         if label == 1:
          8             # CASE: TP
          9             if label == pred:
         10                 cm['TP'] += 1
         11             # CASE: FN
         12             else:
         13                 cm['FN'] += 1
         14         else:
         15             # CASE: TN
         16             if label == pred:
         17                 cm['TN'] += 1
         18             # CASE: FP
         19             else:
         20                 cm['FP'] += 1
         21     return cm

```

## ▼ 4 Baseline Model - Logistic Regression

### ▼ 4.1 Fit data to model

Using a simpler model (Logistic Regression) for the baseline model.

```

In [20]: 1 # model
          2 logreg = LogisticRegression(fit_intercept=False, C=1e12, solver='liblinear',
          3                               random_state=1)
          4 model_log = logreg.fit(X_train_resampled, y_train_resampled)
          5 model_log

```

```

Out[20]: LogisticRegression(C=1000000000000.0, fit_intercept=False, random_state=1,
                             solver='liblinear')

```

### ▼ 4.2 Evaluate

```

In [21]: 1 # Create prediction and residuals for training data
          2 y_hat_train = logreg.predict(X_train_resampled)

```

```
In [22]: 1 # Confusion matrix and classification report for training data.
2 print(confusion_matrix(y_train_resampled, y_hat_train))
3 print(classification_report(y_train_resampled, y_hat_train))
4 print("The accuracy score is" + " " + str(accuracy_score(y_train_resampled,
```

```
[[1637 484]
 [ 460 1661]]
```

	precision	recall	f1-score	support
0	0.78	0.77	0.78	2121
1	0.77	0.78	0.78	2121
accuracy			0.78	4242
macro avg	0.78	0.78	0.78	4242
weighted avg	0.78	0.78	0.78	4242

The accuracy score is 0.7774634606317775

- A baseline accuracy score should be over 85.5% because that is what guessing 'churn' every time will score.
- Baseline training model produces about the same amount of FP's and FN's.
- Precision, recall, and f1 is about the same for churn vs. no churn categories.

Investigate test results

```
In [23]: 1 y_hat_test = logreg.predict(X_test)
2
3 print(confusion_matrix(y_test, y_hat_test))
4 print(classification_report(y_test, y_hat_test))
5 print("The accuracy score is" + " " + str(accuracy_score(y_test, y_hat_test))
```

```
[[549 180]
 [ 17  88]]
```

	precision	recall	f1-score	support
0	0.97	0.75	0.85	729
1	0.33	0.84	0.47	105
accuracy			0.76	834
macro avg	0.65	0.80	0.66	834
weighted avg	0.89	0.76	0.80	834

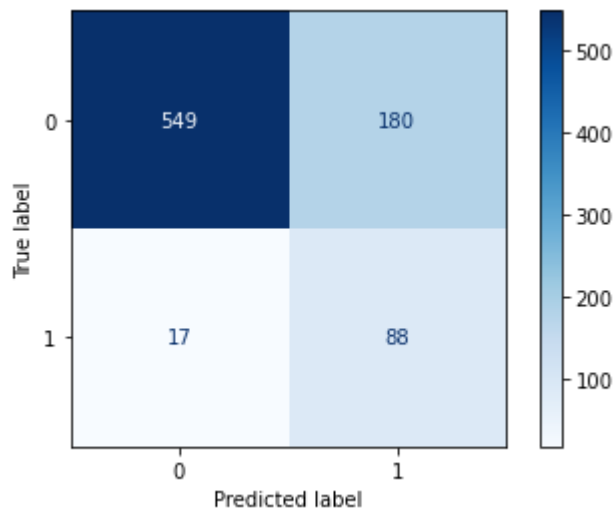
The accuracy score is 0.7637889688249401

- Accuracy score is slightly lower than training data.
- Not many false negatives, which is good. However, a lot of false positives present.
- Precision very low for churn customers. F1 low for churn customers.

```
In [24]: 1 conf_matrix(y_test, y_hat_test)
```

```
Out[24]: {'TP': 88, 'TN': 549, 'FP': 180, 'FN': 17}
```

```
In [25]: 1 plot_confusion_matrix(logreg, X_test, y_test,
2                               cmap=plt.cm.Blues)
3 plt.show()
```



Better visualization than the string printout from above...

It can be seen that there are quite a few false positives in the baseline model.

The large amount of FP's here ultimately decreases the score of this model.

```
In [26]: 1 # create list for storing F-Beta2 scores for each model
2 fbeta2_scores = []
3
4 f = fbeta_score(y_test, y_hat_test, beta=2.0)
5 fbeta2_scores.append(round(f, 3))
6 print(fbeta2_scores)
```

[0.64]

Logistic Regression results:

- Both accuracy scores are below a the strategy of simply guessing "not churned" for every customer.
- It will be key to reduce the number of customers that are identified as retained but actually churned. Ideally everyone that can be identified as churning will be retained with some kind of strategy.
- This model produces a lot of false positives, which ultimately penalized the F2 score.

## 5 Model 2 - RF

Using a random forest model to make use of a more complex model that will likely generate better scores than the baseline. Random forest will allow to make use of a wide range of random parameter selection, which is beneficial for this dataset because of the large amount of variables.

### 5.1 Fit model

```
In [27]: 1 # random forest model
          2 forest = RandomForestClassifier(n_estimators=100, max_depth=13, random_state=1)
          3 forest.fit(X_train_resampled, y_train_resampled)
```

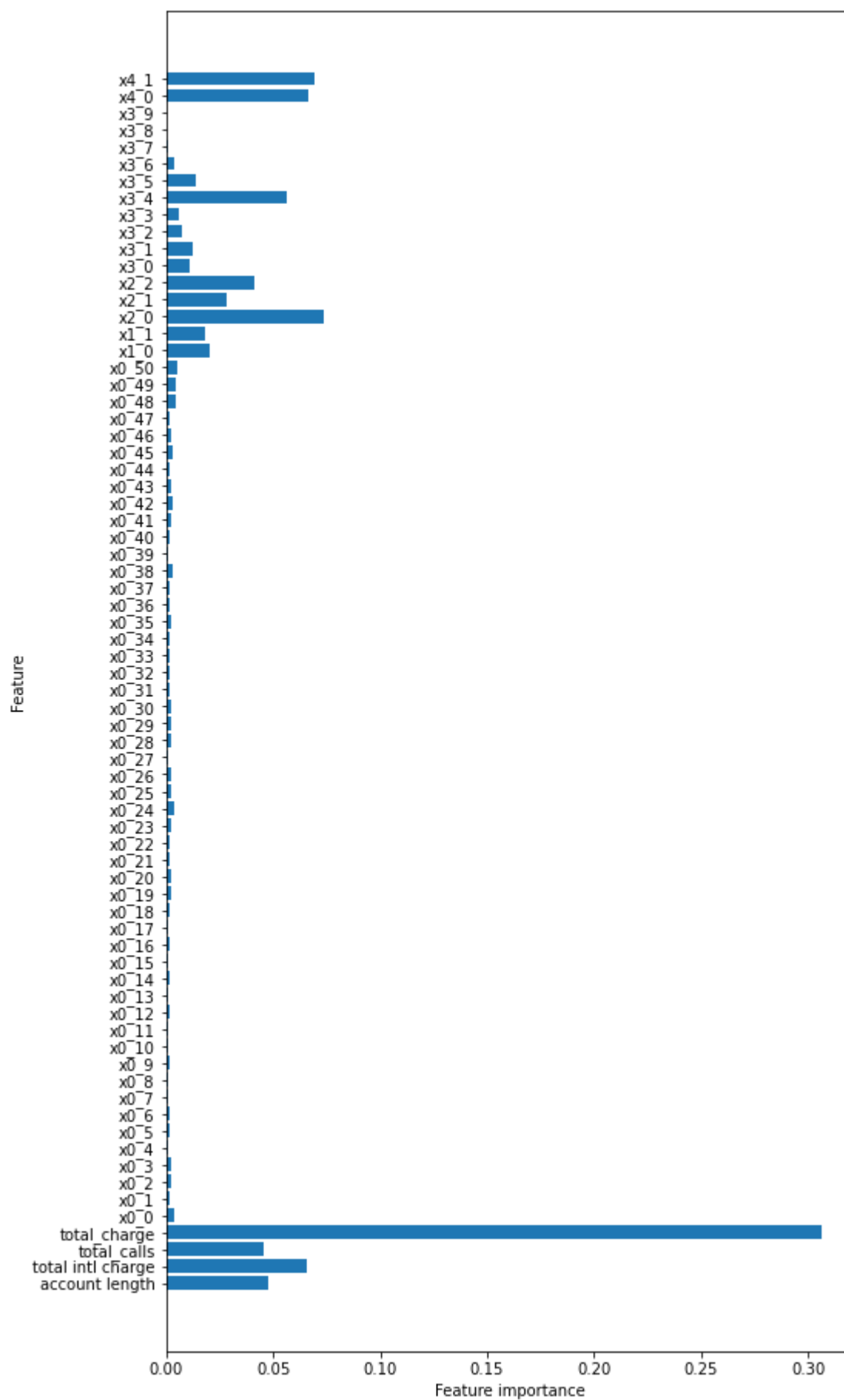
Out[27]: RandomForestClassifier(max\_depth=13, random\_state=1)



## 5.2 Evaluate

```
In [28]: 1 def plot_feature_importances(model):
          2     n_features = X_train_resampled.shape[1]
          3     plt.figure(figsize=(8,16))
          4     plt.barh(range(n_features), model.feature_importances_, align='center')
          5     plt.yticks(np.arange(n_features), X_train_resampled.columns.values)
          6     plt.xlabel('Feature importance')
          7     plt.ylabel('Feature')
```

```
In [29]: 1 plot_feature_importances(forest)
```



- Individual states were not very indicative of anything.
- Total charge is most important feature by far.
- Having a voice mail plan was the most important feature out of any of the categorical features.

```
In [30]: 1 # Training Accuracy  
        2 forest.score(X_train_resampled, y_train_resampled)
```

```
Out[30]: 0.8814238566713815
```

```
In [31]: 1 # Testing Accuracy  
        2 forest.score(X_test, y_test)
```

```
Out[31]: 0.8776978417266187
```

Similar score between train and test sets may indicate lack of overfitting.



In [32]:

```

1 y_hat_train = forest.predict(X_train_resampled)
2
3 print(confusion_matrix(y_train_resampled, y_hat_train))
4 print(classification_report(y_train_resampled, y_hat_train))

```

```

[[1965  156]
 [ 347 1774]]

```

	precision	recall	f1-score	support
0	0.85	0.93	0.89	2121
1	0.92	0.84	0.88	2121
accuracy			0.88	4242
macro avg	0.88	0.88	0.88	4242
weighted avg	0.88	0.88	0.88	4242

- Precision, recall, and f1-score all look pretty good.
- There are many more FN's than FP's, which is not great for F2-Score, which is the metric being used to compare models.

In [33]:

```

1 y_hat_test = forest.predict(X_test)
2
3 print(confusion_matrix(y_test, y_hat_test))
4 print(classification_report(y_test, y_hat_test))

```

```

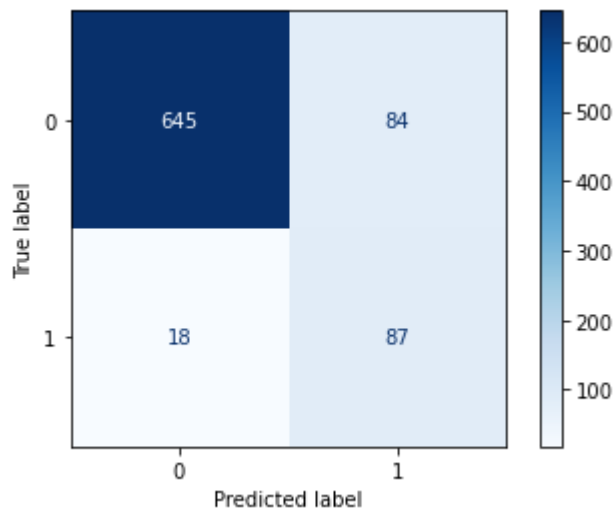
[[645  84]
 [ 18  87]]

```

	precision	recall	f1-score	support
0	0.97	0.88	0.93	729
1	0.51	0.83	0.63	105
accuracy			0.88	834
macro avg	0.74	0.86	0.78	834
weighted avg	0.91	0.88	0.89	834

- F1-score for churned customers is low.
- Precision for churned customers is low.
- Not many FN's, but still quite a few FP's.

```
In [34]: 1 plot_confusion_matrix(forest, X_test, y_test,  
2                               cmap=plt.cm.Blues)  
3 plt.show()
```



Making this plot again to easily visualize the confusion matrix for the test data.

```
In [35]: 1 f = fbeta_score(y_test, y_hat_test, beta=2.0)  
2 fbeta2_scores.append(round(f, 3))  
3 fbeta2_scores
```

```
Out[35]: [0.64, 0.736]
```

## ▼ 6 Model 3 - Optimize RF Model

### ▼ 6.1 Fit and use GridSearchCV to find best parameters

Creating F-2 Scorer for GridSearchCV so that it can optimize that as a score instead of accuracy or something else.

```
In [36]: 1 # FBeta2 Scorer  
2 ftwo_scorer = make_scorer(fbeta_score, beta=2)
```

```

In [37]: 1 # Using a pipeline to select for optimal parameters in the RF Classifier
2 pipe = Pipeline([('clf', RandomForestClassifier(random_state=1))])
3
4 grid_params = [{'clf__n_estimators': [100, 1000],
5                      'clf__criterion': ['gini', 'entropy'],
6                      'clf__max_depth': [None, 1, 5, 7, 11, 13, 17, 19],
7                      'clf__min_samples_split': [2, 5, 10],
8                      'clf__min_samples_leaf': [2, 3, 4, 5]
9                      }]
10
11 # grid search
12 gs = GridSearchCV(estimator=pipe,
13                   param_grid=grid_params,
14                   scoring=f'two_scorer',
15                   cv=5)
16
17 # Fit using grid search
18 gs.fit(X_train_resampled, y_train_resampled)
19
20 # Best F-2 Score
21 print('Best F-2 Score: %.3f' % gs.best_score_)
22
23 # Best params
24 print('\nBest params:\n', gs.best_params_)

```

Best F-2 Score: 0.833

Best params:

```
{'clf__criterion': 'gini', 'clf__max_depth': None, 'clf__min_samples_leaf': 2,
'clf__min_samples_split': 5, 'clf__n_estimators': 1000}
```

This GridSearchCV code took my machine about 40 minutes to run.

```

In [38]: 1 forest = RandomForestClassifier(n_estimators=1000, criterion='gini',
2                                     max_depth=None, min_samples_leaf=2,
3                                     min_samples_split=5, random_state=1)
4 forest.fit(X_train_resampled, y_train_resampled)

```

```
Out[38]: RandomForestClassifier(min_samples_leaf=2, min_samples_split=5,
                                n_estimators=1000, random_state=1)
```

## ▼ 6.2 Evaluate

```

In [39]: 1 # Training Accuracy
2 forest.score(X_train_resampled, y_train_resampled)

```

```
Out[39]: 0.9556812824139557
```

```
In [40]: 1 # Testing Accuracy
        2 forest.score(X_test, y_test)
```

Out[40]: 0.9016786570743405

```
In [41]: 1 print(classification_report(y_train_resampled, y_hat_train))
        2 print(classification_report(y_test, y_hat_test))
```

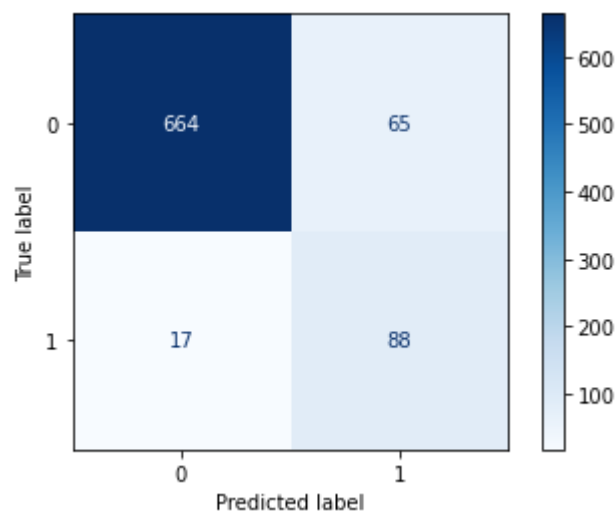
	precision	recall	f1-score	support
0	0.85	0.93	0.89	2121
1	0.92	0.84	0.88	2121
accuracy			0.88	4242
macro avg	0.88	0.88	0.88	4242
weighted avg	0.88	0.88	0.88	4242

	precision	recall	f1-score	support
0	0.97	0.88	0.93	729
1	0.51	0.83	0.63	105
accuracy			0.88	834
macro avg	0.74	0.86	0.78	834
weighted avg	0.91	0.88	0.89	834

- Train report looks good.
- Test report shows low precision in churn column.

```
In [42]: 1 plot_confusion_matrix(forest, X_test, y_test,
        2                          cmap=plt.cm.Blues)
        3 plt.show()
```



High FP rate.

In [43]: 1 plot\_feature\_importances(forest)



Similar feature importances as the first RF model.

```
In [44]: 1 y_hat_test = forest.predict(X_test)
2 f = fbeta_score(y_test, y_hat_test, beta=2.0)
3 fbeta2_scores.append(round(f, 3))
4 fbeta2_scores
```

Out[44]: [0.64, 0.736, 0.768]

## 7 Model 4 - Optimize RF Model (2nd Attempt)

Attempting to reduce runtime by decreasing parameters a bit, but this code still does take a while to run.

To avoid overfitting in the RF Model, the following strategies will be utilized:

- Use less features (use something like 50% of total features for the max features parameter).
- Keep n\_estimators large because the more trees there are, the less likely it is to overfit.

### 7.1 Optimize RF Parameters and Fit

```

In [45]: 1 # Using a pipeline to select for optimal parameters in the RF Classifier
2 pipe = Pipeline([('clf', RandomForestClassifier(random_state=1))])
3
4 grid_params = [{'clf__n_estimators': [100],
5                  'clf__criterion': ['entropy'],
6                  'clf__max_depth': [None, 1, 7, 15],
7                  'clf__min_samples_split': [2, 5, 10],
8                  'clf__min_samples_leaf': [2, 3, 4, 5],
9                  'clf__max_features': [37]      #this is 50% of total features
10             }]
11
12 # grid search
13 gs = GridSearchCV(estimator=pipe,
14                   param_grid=grid_params,
15                   scoring=f2_score,
16                   cv=5)
17
18 # Fit using grid search
19 gs.fit(X_train_resampled, y_train_resampled)
20
21 # Best F-2 Score
22 print('Best F-2 Score: %.3f' % gs.best_score_)
23
24 # Best params
25 print('\nBest params:\n', gs.best_params_)

```

Best F-2 Score: 0.859

Best params:

```
{'clf__criterion': 'entropy', 'clf__max_depth': None, 'clf__max_features': 37,
'clf__min_samples_leaf': 2, 'clf__min_samples_split': 5, 'clf__n_estimators': 100}
```

```

In [46]: 1 forest = RandomForestClassifier(n_estimators=100, criterion='entropy',
2                                     max_depth=None, min_samples_leaf=2,
3                                     min_samples_split=5, random_state=1,
4                                     max_features=37)
5 forest.fit(X_train_resampled, y_train_resampled)

```

```
Out[46]: RandomForestClassifier(criterion='entropy', max_features=37, min_samples_leaf=
2,
                                min_samples_split=5, random_state=1)
```

## 7.2 Evaluate

```

In [47]: 1 # Training Accuracy
2 forest.score(X_train_resampled, y_train_resampled)

```

```
Out[47]: 0.9893917963224894
```

```
In [48]: 1 # Testing Accuracy
         2 forest.score(X_test, y_test)
```

Out[48]: 0.9136690647482014

```
In [49]: 1 print(classification_report(y_train_resampled, y_hat_train))
         2 print(classification_report(y_test, y_hat_test))
```

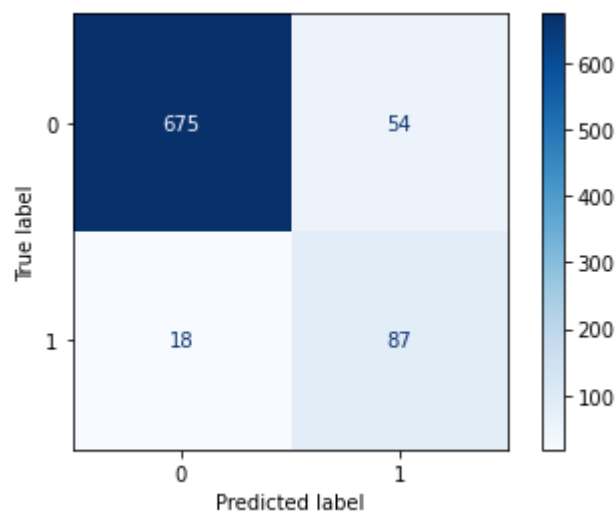
	precision	recall	f1-score	support
0	0.85	0.93	0.89	2121
1	0.92	0.84	0.88	2121
accuracy			0.88	4242
macro avg	0.88	0.88	0.88	4242
weighted avg	0.88	0.88	0.88	4242

	precision	recall	f1-score	support
0	0.98	0.91	0.94	729
1	0.58	0.84	0.68	105
accuracy			0.90	834
macro avg	0.78	0.87	0.81	834
weighted avg	0.92	0.90	0.91	834

- Improved precision slightly from other optimized RF.
- Pretty similar performance as other optimized RF.

```
In [50]: 1 plot_confusion_matrix(forest, X_test, y_test,
         2                          cmap=plt.cm.Blues)
         3 plt.show()
```



- Decrease in amount of FP's from previous model.

```
In [51]: 1 y_hat_test = forest.predict(X_test)
2 f = fbeta_score(y_test, y_hat_test, beta=2.0)
3 fbeta2_scores.append(round(f, 3))
4 fbeta2_scores
```

```
Out[51]: [0.64, 0.736, 0.768, 0.775]
```

## ▼ 8 Model 5 - AdaBoost

The motivation for this model is to reduce runtime expense found in the RF model optimization processes.

### ▼ 8.1 Instantiate Model and Fit

```
In [52]: 1 from sklearn.ensemble import AdaBoostClassifier
2 adaboost_clf = AdaBoostClassifier(random_state=42)
3 adaboost_clf.fit(X_train_resampled, y_train_resampled)
```

```
Out[52]: AdaBoostClassifier(random_state=42)
```

### ▼ 8.2 Evaluate

```
In [53]: 1 # AdaBoost model predictions for test and train sets
2 adaboost_train_preds = adaboost_clf.predict(X_train_resampled)
3 adaboost_test_preds = adaboost_clf.predict(X_test)
```

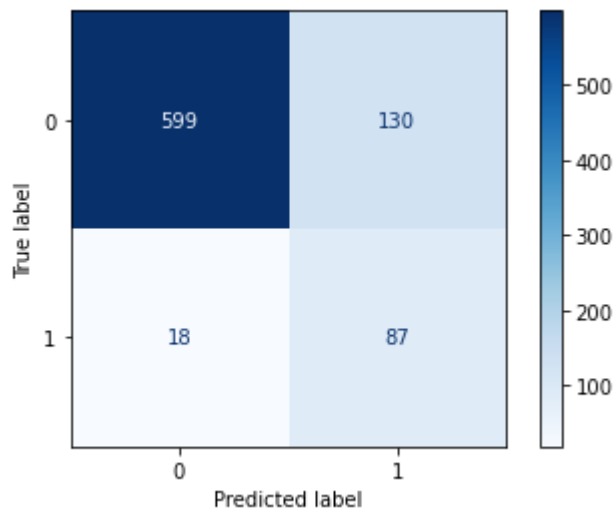
```
In [54]: 1 print(classification_report(y_test, adaboost_test_preds))
```

	precision	recall	f1-score	support
0	0.97	0.82	0.89	729
1	0.40	0.83	0.54	105
accuracy			0.82	834
macro avg	0.69	0.83	0.72	834
weighted avg	0.90	0.82	0.85	834

- Very low precision for churned customers.



```
In [55]: 1 plot_confusion_matrix(adaboost_clf, X_test, y_test,  
2                             cmap=plt.cm.Blues)  
3 plt.show()
```



Lots of FP's!

```
In [56]: 1 f = fbeta_score(y_test, adaboost_test_preds, beta=2.0)  
2 fbeta2_scores.append(round(f, 3))  
3 fbeta2_scores
```

```
Out[56]: [0.64, 0.736, 0.768, 0.775, 0.683]
```

Comments on AdaBoost Model:

- It is very fast
- It has a much higher rate of FPs compared to the RF models.
- It does not have as good of scoring (by any metric of comparison) as the RF models

## ▼ 8.3 Tune Model with GridSearchCV and Run Again

```
In [57]: 1 gs = GridSearchCV(estimator=adaboost_clf,
2                       param_grid={
3                           'n_estimators': [25, 50, 100, 500],
4                           'learning_rate': [.001, .01, .1, 1]
5                       }, scoring=ftwo_scorer, cv=5)
6
7 gs.fit(X_train_resampled, y_train_resampled)
8 gs.best_params_
```

Out[57]: {'learning\_rate': 1, 'n\_estimators': 50}

```
In [58]: 1 adaboost_clf = AdaBoostClassifier(learning_rate=1, n_estimators=50,
2                                           random_state=42)
3 adaboost_clf.fit(X_train_resampled, y_train_resampled)
```

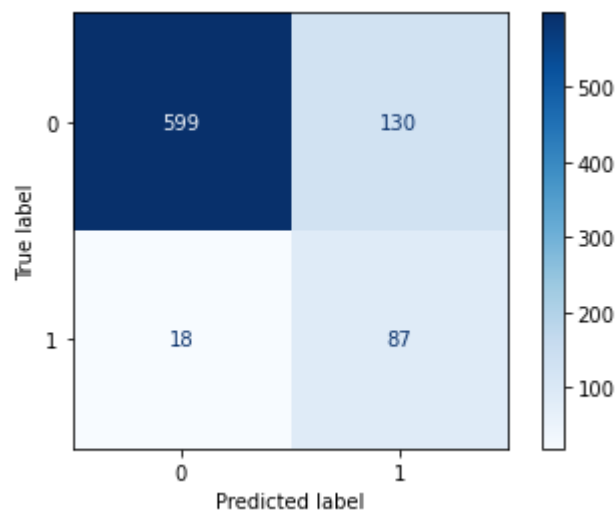
Out[58]: AdaBoostClassifier(learning\_rate=1, random\_state=42)

```
In [59]: 1 # AdaBoost model predictions
2 adaboost_train_preds = adaboost_clf.predict(X_train_resampled)
3 adaboost_test_preds = adaboost_clf.predict(X_test)
4
5 print(classification_report(y_test, adaboost_test_preds))
```

	precision	recall	f1-score	support
0	0.97	0.82	0.89	729
1	0.40	0.83	0.54	105
accuracy			0.82	834
macro avg	0.69	0.83	0.72	834
weighted avg	0.90	0.82	0.85	834

Not much of an improvement from first AdaBoost model.

```
In [60]: 1 plot_confusion_matrix(adaboost_clf, X_test, y_test,
2                               cmap=plt.cm.Blues)
3 plt.show()
```



Still lots of FP's.

```
In [61]: 1 f = fbeta_score(y_test, adaboost_test_preds, beta=2.0)
          2 fbeta2_scores.append(round(f, 3))
          3 fbeta2_scores
```

```
Out[61]: [0.64, 0.736, 0.768, 0.775, 0.683, 0.683]
```

## ▼ 9 RF Without State Column

The 'State' column is probably not particularly useful for SyriaTel, and none of the states seem to have super strong feature importance by themselves.

### ▼ 9.1 Remove State Column and Prepare Dataset to Model

```
In [62]: 1 # define new categorical columns
2 cats = ['international plan', 'intl_calls_bins',
3         'customer service calls', 'voice mail plan']
4
5 # Define X and y
6 y = df['churn']
7
8 # drop state_int column as well this time
9 X = df.drop(['churn', 'total intl calls',
10             'total_intl_calls', 'state', 'total day minutes',
11             'total day calls', 'total day charge', 'total eve minutes',
12             'total eve calls', 'total eve charge', 'total night minutes',
13             'total night calls', 'total night charge', 'state_int'], axis=
14
15 # Split
16 X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=27)
17
18 # Fixing class imbalance with SMOTE
19 smote = SMOTE(random_state=42)
20 X_train_resampled, y_train_resampled = smote.fit_sample(X_train, y_train)
21
22 # Seperate data into categorical for train and test sets
23 X_train_cats = X_train_resampled[cats]
24 X_test_cats = X_test[cats]
25
26 # handle categorical values
27 ohe = OneHotEncoder(handle_unknown="ignore", sparse=False)
28
29 # OHE for training categoricals
30 ohe.fit(X_train_cats)
31 X_train_ohe = pd.DataFrame(
32     ohe.transform(X_train_cats),
33     index=X_train_cats.index,
34     columns=ohe.get_feature_names()
35 )
36
37 # OHE for testing categoricals - transform with train fit
38 X_test_ohe = pd.DataFrame(
39     ohe.transform(X_test_cats),
40     index=X_test_cats.index,
41     columns=ohe.get_feature_names()
42 )
43
44 # Scaling variables to work well with OHE data -scale train and test data
45 X_train_numerics = X_train_resampled.drop(cats, axis=1)
46 X_test_numerics = X_test.drop(cats, axis=1)
47
48 scaler = MinMaxScaler()
49
50 scaler.fit(X_train_numerics)
51 X_train_scaled = pd.DataFrame(
52     scaler.transform(X_train_numerics),
53     index=X_train_numerics.index,
54     columns=X_train_numerics.columns
55 )
56
```

```

57 # Scale variables of test data with the train fit (avoids data Leakage)
58 X_test_scaled = pd.DataFrame(
59     scaler.transform(X_test_numerics),
60     index=X_test_numerics.index,
61     columns=X_test_numerics.columns
62 )
63
64
65 # Concatenate and replace X_train and X_test with OHE+scaled data
66 X_train_resampled = pd.DataFrame()
67 X_train_resampled = pd.concat([X_train_scaled, X_train_ohe], axis=1)
68
69 X_test = pd.DataFrame()
70 X_test = pd.concat([X_test_scaled, X_test_ohe], axis=1)

```

## 9.2 Fit Model

In [63]:

```

1 # random forest model
2 forest = RandomForestClassifier(n_estimators=100, max_depth=13,
3                               random_state=69)
4 forest.fit(X_train_resampled, y_train_resampled)

```

Out[63]: RandomForestClassifier(max\_depth=13, random\_state=69)

## 9.3 Evaluate

In [64]:

```

1 print(classification_report(y_train_resampled, y_hat_train))
2 print(classification_report(y_test, y_hat_test))

```

	precision	recall	f1-score	support
0	0.85	0.93	0.89	2121
1	0.92	0.84	0.88	2121
accuracy			0.88	4242
macro avg	0.88	0.88	0.88	4242
weighted avg	0.88	0.88	0.88	4242

	precision	recall	f1-score	support
0	0.97	0.93	0.95	729
1	0.62	0.83	0.71	105
accuracy			0.91	834
macro avg	0.80	0.88	0.83	834
weighted avg	0.93	0.91	0.92	834

- Metrics look really good, but precision for churn would ideally be higher.

```
In [65]: 1 # Training Accuracy  
2 forest.score(X_train_resampled, y_train_resampled)
```

Out[65]: 0.9464875058934464

```
In [66]: 1 # Testing Accuracy  
2 forest.score(X_test, y_test)
```

Out[66]: 0.9304556354916067

Accuracy is high for both train and test data, and they are similar, indicating lower chance of overfitting.

```
In [67]: 1 final_cols = ['account length', 'total intl charge',  
2                   'total_calls', 'total_charge',  
3                   'no international plan', 'has international plan',  
4                   'low international calls', 'moderate international calls',  
5                   'high international calls', '0 customer service calls',  
6                   '1 customer service call', '2 customer service calls',  
7                   '3 customer service calls',  
8                   '4 customer service calls', '5 customer service calls',  
9                   '6 customer service calls', '7 customer service calls',  
10                  '8 customer service calls', '9 customer service calls',  
11                  'no voice mail plan', 'has voice mail plan']
```

```

In [68]: 1 # Create table of feature importances and values
          2 df_final_feature_importances = pd.DataFrame(index=(final_cols))
          3 # df_final_feature_importances.set_axis([final_cols], axis=1)
          4 df_final_feature_importances['feature importances values'] = np.round(fore
          5 df_final_feature_importances = df_final_feature_importances.sort_values(
          6     by=['feature importances values'], ascending=False)
          7 df_final_feature_importances

```

Out[68]:

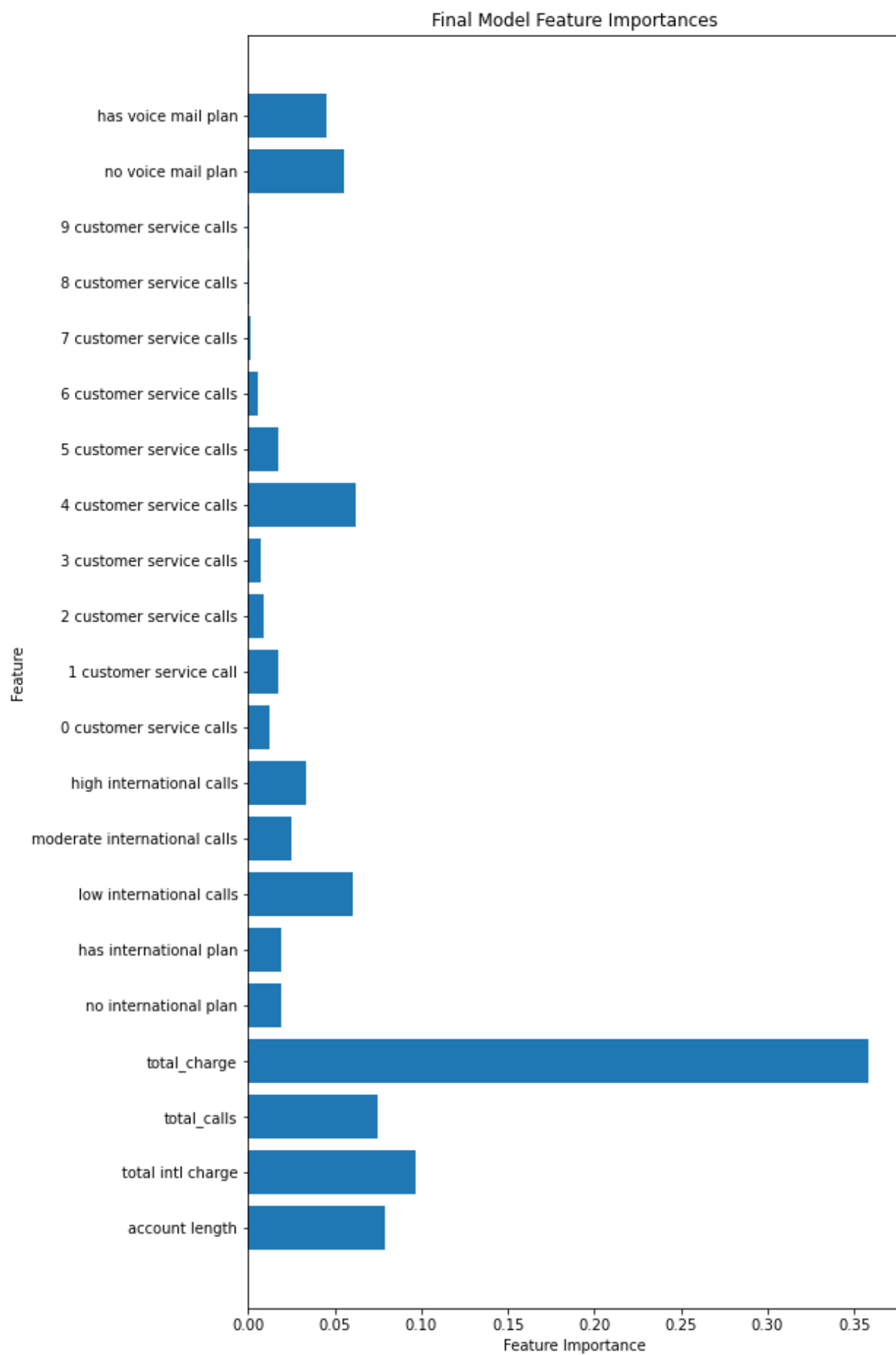
feature importances values	
total_charge	0.359
total intl charge	0.097
account length	0.079
total_calls	0.074
4 customer service calls	0.062
low international calls	0.060
no voice mail plan	0.055
has voice mail plan	0.045
high international calls	0.033
moderate international calls	0.025
has international plan	0.019
no international plan	0.019
1 customer service call	0.018
5 customer service calls	0.017
0 customer service calls	0.013
2 customer service calls	0.009
3 customer service calls	0.007
6 customer service calls	0.006
7 customer service calls	0.001
8 customer service calls	0.001
9 customer service calls	0.000

- Total charge is by far the most important feature.
- International charge is the next important.
- Account length and total calls are about the same and are the third and fourth most important categories.

```
In [69]: 1 # creating new feature importances plot function to allow for more
2 # comprehensability for this plot, which will be used elsewhere as well.
3
4 def plot_feature_importances_final(model):
5     n_features = X_train_resampled.shape[1]
6     plt.figure(figsize=(8,16))
7     plt.barh(range(n_features), model.feature_importances_, align='center')
8     plt.yticks(np.arange(n_features), final_cols)
9     plt.xlabel('Feature Importance')
10    plt.ylabel('Feature')
11    plt.title('Final Model Feature Importances')
```

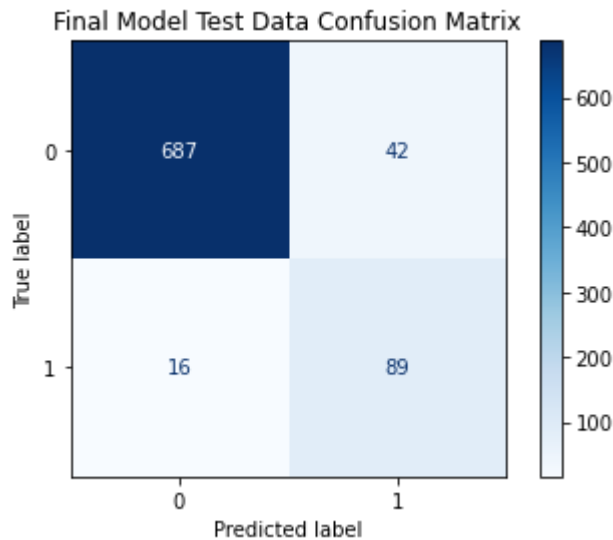


```
In [70]: 1 plot_feature_importances_final(forest)
```



- 4 customer service calls is the most important categorical column.
- total charge is very important in the model
- total international charge, account length, and total calls are the next most important features.

```
In [71]: 1 plot_confusion_matrix(forest, X_test, y_test,  
2                             cmap=plt.cm.Blues)  
3 plt.title("Final Model Test Data Confusion Matrix")  
4 plt.show()
```



- Very low amount of false negatives.
- Pretty low number of false positives as well. This column would ideally be lower, but its penalty is less than FN's.

```
In [72]: 1 y_hat_test = forest.predict(X_test)  
2 f = fbeta_score(y_test, y_hat_test, beta=2.0)  
3 fbeta2_scores.append(round(f, 3))  
4 fbeta2_scores
```

```
Out[72]: [0.64, 0.736, 0.768, 0.775, 0.683, 0.683, 0.808]
```

## ▼ 9.4 Tune Model with GridSearchCV

```
In [73]: 1 forest = RandomForestClassifier(random_state=1)
2
3 gs = GridSearchCV(estimator=forest,
4                   param_grid={
5                       'n_estimators': [10, 100, 1000],
6                       'criterion': ['entropy', 'gini'],
7                       'max_depth': [None, 1, 7, 15, 23, 31],
8                       'min_samples_split': [2, 5, 10],
9                       'min_samples_leaf': [2, 3, 4, 5, 6]
10                  }, scoring='f2_score', cv=5)
11
12 gs.fit(X_train_resampled, y_train_resampled)
```

```
Out[73]: GridSearchCV(cv=5, estimator=RandomForestClassifier(random_state=1),
                  param_grid={'criterion': ['entropy', 'gini'],
                              'max_depth': [None, 1, 7, 15, 23, 31],
                              'min_samples_leaf': [2, 3, 4, 5, 6],
                              'min_samples_split': [2, 5, 10],
                              'n_estimators': [10, 100, 1000]},
                  scoring='f2_score', cv=5)
```

Approximately 42 minute runtime for this GridSearchCV.

```
In [74]: 1 # Best F-2 Score
2 print('Best F2 Score: %.3f' % gs.best_score_)
3
4 # Best params
5 print('\nBest params:\n', gs.best_params_)
```

Best F2 Score: 0.854

Best params:

```
{'criterion': 'gini', 'max_depth': 31, 'min_samples_leaf': 2, 'min_samples_split': 2, 'n_estimators': 1000}
```

```
In [75]: 1 forest = RandomForestClassifier(n_estimators=1000, criterion='gini',
2                                     min_samples_leaf=2, min_samples_split=2,
3                                     max_depth=31, random_state=42)
4 forest.fit(X_train_resampled, y_train_resampled)
```

```
Out[75]: RandomForestClassifier(max_depth=31, min_samples_leaf=2, n_estimators=1000,
                               random_state=42)
```

## ▼ 9.5 Re-Evaluate

```
In [76]: 1 print(classification_report(y_train_resampled, y_hat_train, digits=4))
        2 print(classification_report(y_test, y_hat_test))
```

	precision	recall	f1-score	support
0	0.8499	0.9264	0.8865	2121
1	0.9192	0.8364	0.8758	2121
accuracy			0.8814	4242
macro avg	0.8845	0.8814	0.8812	4242
weighted avg	0.8845	0.8814	0.8812	4242

	precision	recall	f1-score	support
0	0.98	0.94	0.96	729
1	0.68	0.85	0.75	105
accuracy			0.93	834
macro avg	0.83	0.90	0.86	834
weighted avg	0.94	0.93	0.93	834

Fairly similar scores to the previous model.

```
In [77]: 1 # Training Accuracy
        2 forest.score(X_train_resampled, y_train_resampled)
```

Out[77]: 0.9813767090994814

```
In [78]: 1 # Testing Accuracy
        2 forest.score(X_test, y_test)
```

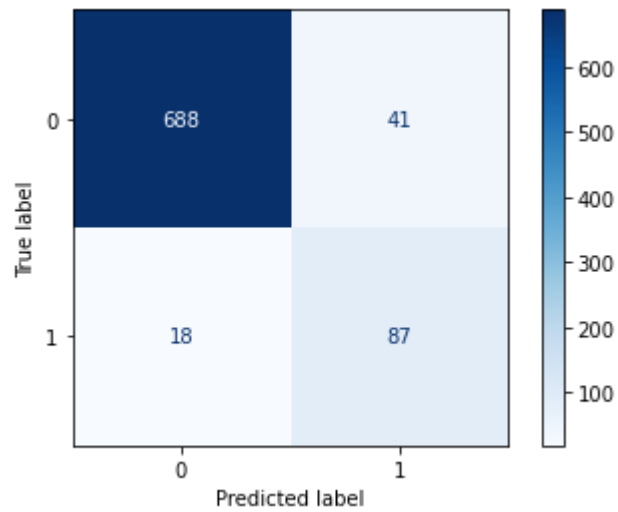
Out[78]: 0.9292565947242206

Difference in testing and training accuracy scores suggests overfitting.

```
In [79]: 1 y_pred = forest.predict(X_test)
        2 conf_matrix(y_test, y_pred)
```

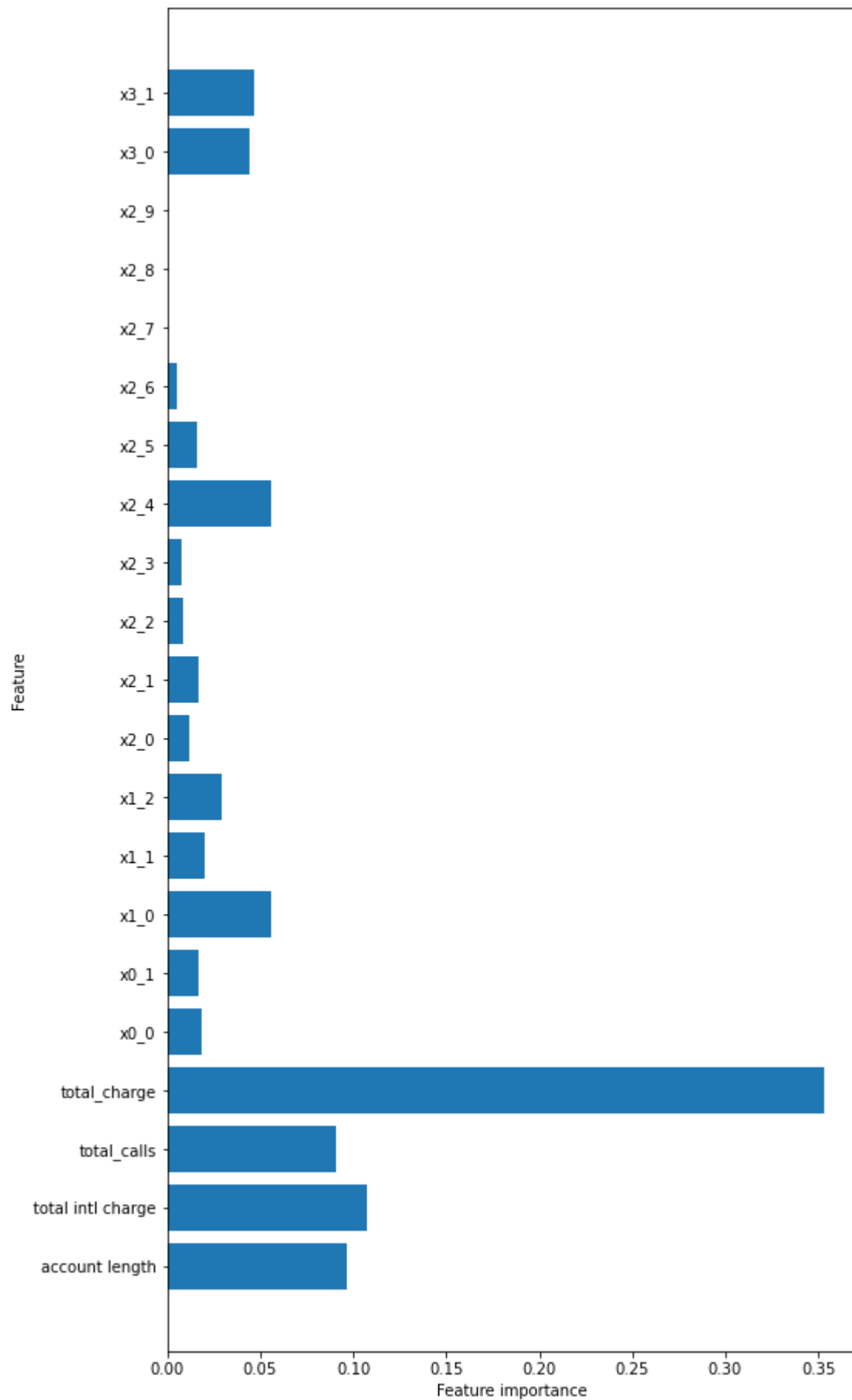
Out[79]: {'TP': 87, 'TN': 688, 'FP': 41, 'FN': 18}

```
In [80]: 1 plot_confusion_matrix(forest, X_test, y_test, cmap=plt.cm.Blues)  
        2 plt.show()
```



FN's and FP's are just slightly above the previous model.

```
In [81]: 1 plot_feature_importances(forest)
```



Pretty similar results as the previous model.

```
In [82]: 1 y_hat_test = forest.predict(X_test)
          2 f = fbeta_score(y_test, y_hat_test, beta=2.0)
          3 fbeta2_scores.append(round(f, 3))
          4 fbeta2_scores
```

```
Out[82]: [0.64, 0.736, 0.768, 0.775, 0.683, 0.683, 0.808, 0.794]
```

Slightly lower F-2 Score than the previous model.

Suppressing the following code because it produces a large file size for the image it makes. The following code creates an image of a single tree in the forest in the section 8 tuned model.



```
In [83]: 1 # try to visualize the tree
2 # adapted from: https://towardsdatascience.com/how-to-visualize-a-decision-
3 # credits to Will Koehresen for the code above (link above)
4
5 # from sklearn.tree import export_graphviz
6 # # Export as dot file
7 # export_graphviz(forest.estimators_[3], out_file='tree.dot',
8 #                 feature_names = X_train_resampled.columns.tolist(),
9 #                 class_names = 'churn',
10 #                 rounded = True, proportion = False,
11 #                 precision = 3, filled = True)
12
13 # # Convert to png using system command (requires Graphviz)
14 # from subprocess import call
15 # call(['dot', '-Tpng', 'tree.dot', '-o', 'tree.png', '-Gdpi=600'])
16
17 # # Display in jupyter notebook
18 # from IPython.display import Image
19 # Image(filename = 'tree.png')
```

## 10 Final Model Selection and Discussion

### 10.1 Overview

The goal of this project is to reduce customer churn rate for SyriaTel. This dataset is based in the US, so it is important for model results to extrapolate well to an application in Syria (where SyriaTel is based). To reduce customer churn rate, it is important to identify customers that are likely to churn, why these customers churn, and deploy a customer retention strategy on these customers.

Main evaluation metrics to use and the basis of justification for utilizing these metrics:

- The main goal is to implement retention strategies such as offering discounted rates, better plans, etc. to customers identified as likely to churn.
- Customers that are classified as churning by the model will be offered a retention strategy that will probably cost money to implement, but it can be justified by retaining the customer in the long term.
- Accuracy is important for minimizing total false negatives and false positives.
- False negatives (customers identified as retained but are not actually retained) are important to minimize because these are customers that "fell through the cracks" and left without any attempt to save their business.
- False positives are important to minimize because these customers will potentially represent an unneeded cost because they will be offered one of the retention strategies, which costs the company money, even though they don't need incentive to stay with the company. This loss could be countered with some type of survey to identify dissatisfied customers, but it is important to decrease the amount of false positives.

- False negatives (customers identified as not churning but actually churned) should be penalized more because these customers do not have the potential to be retained with any particular strategy, they simply fall through the cracks. This cost is greater than a decent retention strategy because the retention strategy will signify an initial cost, but it will allow for the customer to be retained for a longer time.

Other justifications for model selection:

- The model must have potential to extrapolate from this dataset (from a US-based company) to the stakeholder (SyriaTel).
- It is ideal for the model to be deployed somewhat quickly so that customers can be classified as likely to churn or not at the end of various time periods, and so that customer retention strategies can be deployed quickly.

Potential Retention Strategies:

- If the customer has an international plan, then offer a certain amount of months for free. The international plan in this dataset doesn't seem to do anything (check EDA notebook), and it probably wouldn't be of any cost for this particular company to offer international plans free for an extended period. But this idea is one strategy for retaining customers with international plans.
- Offer a discount for call charges. Total charge was an important feature in most models, so it will be useful to offer discounts to retain customers.
- Create better customer service methods and offer solutions quickly to customers that call multiple times.
- Have an effective international plan.

## ▼ 10.2 Customer Retention Strategy and Evaluation Metric

Main Retention Strategy to Consider:

- Offer customers identified as likely to churn the option to have a discount such as 6 months at 50% off.

Items to consider:

- Assumptions have to be made in order to evaluate how effective the retention strategy will be.
- Assumption 1: average cost to acquire a new customer is around 300 dollars. Sources: <https://www.forbes.com/sites/forbestechcouncil/2020/10/30/acquiring-subscribers-is-only-half-the-battle/https://www.entrepreneur.com/article/225415> (<https://www.forbes.com/sites/forbestechcouncil/2020/10/30/acquiring-subscribers-is-only-half-the-battle/https://www.entrepreneur.com/article/225415>)
- Average phone bill cost (in the US) for a single user: 70 dollars/month. Source: <https://www.usmobile.com/blog/cut-cell-phone-bill/> (<https://www.usmobile.com/blog/cut-cell-phone-bill/>)
- Assume that customers that leave the service will cost the company 370 dollars (cost of one month of service + cost to acquire new customer).

- For False Positive customers (model identifies customer as churning, but they do not churn), assume the cost to the company is (35 dollars\* 6 months) = 210 dollars.
- For False Negative customers (model identifies as retained but they churn), assume the cost to the company is 370 dollars.
- For True Positive customers, assume the rate of retention (after offering the discount) is 80%.
- The expected value gained for any True Positive customer is  $.8(-210+370) + .2(-370) = 54$  dollars.

Discussion of Metric to Use (F-Beta2):

- A False Negative costs about 1.762 times as much as a False Positive from the above calculations. The metric will need to implement recognition of this discrepancy.
- Could use expected amount of money gained or lost.
- Use F-Beta with a beta value > 1 to penalize False Negatives more than False Positives.
- Due to the nature of making many assumptions in the previous calculations, F-Beta2 score will be used to compare model performance between each other. This scoring method ensures that FN's are penalized more than FP's. F-Beta2 will not necessarily penalize FN's exactly 1.762 times more than FP's, but it will give a good comparison between models that will offer a good estimation that doesn't have to rely on the assumptions made.

## ▼ 10.3 Final Model and Inferences from the Model

In [84]:

```
1 fbeta2_scores
```

Out[84]: [0.64, 0.736, 0.768, 0.775, 0.683, 0.683, 0.808, 0.794]

The first RF model without the state column performed the best in terms of F-2 Score.

### ▼ 10.3.1 Review the top feature importances of the model

In [85]: 1 df\_final\_feature\_importances.head(10)

Out[85]:

feature importances values	
total_charge	0.359
total intl charge	0.097
account length	0.079
total_calls	0.074
4 customer service calls	0.062
low international calls	0.060
no voice mail plan	0.055
has voice mail plan	0.045
high international calls	0.033
moderate international calls	0.025

By far the most important feature is total charge. It will be important to focus on discounting costs of customers that are predicted to churn. It would be worthwhile to investigate better overall methods to decrease costs for customers. The second most important feature is total international charge, and the third most important feature is total calls. It seems like there may be a common theme with either charging too much or not providing comprehensive enough packages for calling. There are various business strategies that could reduce churn in these cases such as rewarding customers for utilizing the service more or offering more comprehensive plans for customers such as unlimited calling.

### ▼ 10.3.2 Estimate money saved with this model

Without a retention strategy, how much money is lost? SyriaTel has around 8 million customers (according to its LinkedIn page), so let's assume that and account for lost money without a strategy. With the assumptions made in section 10.2, 370 dollars is lost for customers that are not identified as churning. 14.5% of customers in the original dataset churned. With eight million total customers, that would represent a loss of business from 1.16 million customers and \$429.2M.

In the final model, the frequency of FN's (customers that churn but are predicted as retained) is about 1.80%. The frequency of TP's (TP's / total observations) is about 10.80%. The frequency of FP's is about 5.16%. Therefore, FN's will account for 144,000 customers, FP's will account for 412,800 customers, and TP's will account for 864,000 customers. The expected loss (derived from the calculations in section 10.2) equates to 53,280,000 dollars from FN's and 86,688,000 dollars from FP's. The expected gain from TP's is 46,656,000 dollars. In total, SyriaTel is expected to lose about 93,312,000 dollars with this retention strategy, but SyriaTel will save 335,888,000 dollars by implementing this strategy in comparison to having no strategy.

### ▼ 10.3.3 Conclusions

With the retention strategy covered in this notebook, SyriaTel is estimated to save just under 336 million dollars. Obviously SyriaTel is based in Syria and probably does not use dollars, and the estimations for these calculations is based somewhat on US telecommunications data (reference section 10.2). However, the point of this notebook is that the model created allows for SyriaTel to save an enormous amount of money. There are also other considerations that can be inferred from this model that SyriaTel can make use of. For instance, this dataset highlights the need for good customer service, especially after the customer has contacted customer service previously. The dataset also highlights the need for a comprehensive international plan (if that plan is offered) because many customers with an international plan ended up churning in this dataset.

Further considerations:

- Tune estimation parameters to account for telecommunication data in Syria.
- Use a larger telecommunications dataset. There were 3,333 datapoints in this set, which is by no means a large dataset. A larger dataset will allow for more confidence in the predictive models covered in this notebook.
- Tune the hyperparameters in the model and re-adjust the model after more data has been collected.
- Ideally this model and various parameters would be adjusted and predictions would be created again after another time period has surpassed and customer churn rate can be re-evaluated.
- There are other macroeconomic variables in Syria that will affect customer retention. Syria is not the most stable of areas, and the EU has even imposed sanctions on SyriaTel in the past (reference: <https://www.reuters.com/article/us-syria-eu-sanctions-idUSTRE78N1DY20110924> (<https://www.reuters.com/article/us-syria-eu-sanctions-idUSTRE78N1DY20110924>))