

# Table of Contents

<a href="#">1 Abstract</a>		
<a href="#">2 Business Problem For a Telecommunication Company</a>		
<b>Contents ↗ ⚙</b>		
<a href="#">3 Import Libaries</a>		
10.2 Results		
▼ 11 XGBoost	<a href="#">4 Load Data</a>	
11.1 Results	<a href="#">5 Data Exploration &amp; Visualization</a>	
▼ 12 Gaussian NB	▼ <a href="#">6 Data Engineering</a>	
12.1 Results	▼ <a href="#">6.1 Group state column into south and north</a>	
▼ 13 AdaBoostC	6.1.1 Comments	
13.1 Results	6.2 Drop unnecessary columns	
▼ 14 Gradient Bo	6.3 Drop higly correlated features	
14.1 Results	6.4 One Hot Encoder Categorical Columns	
▼ 15 SVM Classi	6.5 Train test split	
▼ 15.1 Tuning	6.6 Standardize Features	
▼ 15.1.1 Hyp	6.7 Building analysis functions	
15.1.1.1	▼ <a href="#">7 Logistic Regression</a>	
▼ 15.2 Results	7.1 Optimization and Tuning	
15.2.1 SVI	7.2 Best performing parameters for Logistic Regression Model	
15.2.2 SV	7.3 Results Comments	
15.3 Results	▼ <a href="#">8 Descion Tree Classifier</a>	
▼ 16 Stacking	8.1 Hyperparameter Tuning for decision tree using Gridsearch	
16.1 Results	8.1.1 Select Best Parameters for Descion tree classifier	
▼ 17 SGDClassif	8.2 Results Comments	
17.1 Results	▼ <a href="#">9 Random Forest Classifier</a>	
▼ 18 Logistic Re	9.1 Hyperparameter Tuning for Random Forest using Gridsearch	
18.1 Results	9.2 Best Paramerts for Random Forest	
▼ 19 SGDClassif	9.3 Results & Comments	
19.1 Tuning	▼ <a href="#">10 KNN Classifier</a>	
19.2 Results	▼ <a href="#">10.1 Tuning &amp; Optimization (Iterate over n values)</a>	
▼ 20 Data Interp	▼ <a href="#">10.1.1 Select Best K for KNN Machine Learning Model</a>	
▼ 20.1 Most In	10.1.1.1 K=3	
20.1.1 Top	▼ <a href="#">10.1.2 Iterate over 1 &lt; n &lt; 26</a>	
▼ 20.2 Top Po	10.1.2.1 K=1	
20.2.1 To	▼ <a href="#">10.2 Results &amp; Comments</a>	
20.2.2 Toj	▼ <a href="#">11 XGBoost</a>	
▼ 21 Conclusion	21.1 Total	11.1 Results & comments
▼ 21.1 Increas	21.1.2 Cus	▼ <a href="#">12 Gaussian NB</a>
21.1.1 Tot	21.1.3 Tot	12.1 Results & Comments
21.1.2 Tot	21.1.4 Tot	▼ <a href="#">13 AdaBoostClassifier</a>
21.1.3 Tot	13.1 Results & Comments	
21.1.4 Tot	▼ <a href="#">14 Gradient Boosting Classifier</a>	
▼ 21.2 Decre	21.2.1 Voi	14.1 Results & comments
21.2.1 Voi	21.2.2 Inte	▼ <a href="#">15 SVM Classifier</a>
21.2.2 Inte	15.1 Tuning & Optimization	
22 Recomend	▼ <a href="#">15.1.1 Hypertuning of SVM Classifier throgh GridSearhCV</a>	
23 Future Wor	15.1.1.1 Best paramets for SVM Classifier	

▼ [15.2 Results & Comments](#)

[15.2.1 SVM with RFE](#)

[15.2.2 SVM with SelectKBest](#)

[15.3 Results & Comments](#)

▼ [16 Stacking](#)

[16.1 Results & Comments](#)

**Contents** ⚙ [17 SGDClassifier \(Lasso & Ridge regularizations\)](#)

10.2 Results

▼ 11 XGBoost

11.1 Results

▼ 12 Gaussian N

12.1 Results

▼ 13 AdaBoostC

13.1 Results

▼ 14 Gradient Bo

14.1 Results

▼ 15 SVM Classi

15.1 Tuning

15.1.1 Hyp

15.1.1.1

▼ 15.2 Results

15.2.1 SVI

15.2.2 SV

15.3 Results

▼ 16 Stacking

16.1 Results

▼ 17 SGDClassif

17.1 Results

▼ 18 Logistic Re

18.1 Results

▼ 19 SGDClassif

19.1 Tuning

19.2 Results

▼ 20 Data Interp

20.1 Most In

20.1.1 Top

20.2 Top Po

20.2.1 To

20.2.2 To

▼ 21 Conclusion

21.1 Increasing chun

[21.1.1 Total day charge](#)

[21.1.2 Customer service calls](#)

[21.1.3 Total evening charge](#)

[21.1.4 Total internatioanl charge](#)

[21.2 Decreasing chun factors](#)

[21.2.1 Voice mail plan](#)

[21.2.2 International plan](#)

[22 Recomendations](#)

[23 Future Work](#)

# 1 Abstract

Customer churn in the telecommunication industry In this project, machine learning models are made to predict possible customer churn in a telecommunication industry. Our models are powerful

which can determine churn factors Therefore, we can give a company the advantage of taking

increas action before losing customers.

21.1.1 Tot

21.1.2 Cus

21.1.3 Tot

21.1.4 Tot

# 2 Business Problem For a Telecommunication Company

21.2 Decrea

21.2.1 Voi

21.2.2 Inte

1) Determine important factors affection churn.

22 Recomend

23 Future Wor

2) Determine factors increasing churn rate.

3) Determine factors decresing churn rate.

### 3 Import Libaries

#### Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
    - 20.1.2 Toj
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [1]:

```

1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4
5 from sklearn.preprocessing import LabelEncoder
6 from sklearn.pipeline import Pipeline
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.model_selection import train_test_split
9
10
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.neighbors import KNeighborsClassifier
13 from sklearn.ensemble import RandomForestClassifier
14
15
16 from sklearn import tree
17 from sklearn import svm
18 from sklearn import ensemble
19 from sklearn import neighbors
20 from sklearn import linear_model
21 from sklearn import metrics
22 from sklearn import preprocessing
23
24 from sklearn import metrics
25 from sklearn.linear_model import LinearRegression
26 from sklearn.model_selection import cross_val_score
27 from sklearn.model_selection import KFold
28 from sklearn.preprocessing import OneHotEncoder
29 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
30 from sklearn.tree import DecisionTreeClassifier
31 from sklearn.ensemble import BaggingClassifier, RandomForestClassifier
32 from sklearn.model_selection import cross_val_score
33
34
35 from sklearn.model_selection import cross_validate
36 %matplotlib inline
37
38 from IPython.display import Image
39 import matplotlib as mlp
40 import matplotlib.pyplot as plt
41 import numpy as np
42 import os
43 import pandas as pd
44 import sklearn
45
46
47 from sklearn.model_selection import ShuffleSplit
48
49 from sklearn.model_selection import StratifiedKFold
50
51 from sklearn.metrics import classification_report
52
53 import warnings
54 warnings.filterwarnings('ignore')
55
56

```

## Contents

10.2 Results
▼ 11 XGBoost
11.1 Results
▼ 12 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 Gradient Bo
14.1 Results
▼ 15 SVM Classi
▼ 15.1 Tuning
15.1.1 Hyp
15.1.1.1
▼ 15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results
▼ 18 Logistic Re
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
▼ 20.1 Most In
20.1.1 Top
▼ 20.2 Top Po
20.2.1 To
20.2.2 Toj
▼ 21 Conclusion
▼ 21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
▼ 21.2 Decreas
21.2.1 Voic
21.2.2 Inte
22 Recomend
23 Future Wor

```

57 import seaborn as sns
58 sns.set_style('darkgrid', {'axes.facecolor': '0.9'})
59
60 from sklearn.preprocessing import LabelEncoder, OneHotEncoder, Standard
61 from sklearn.model_selection import train_test_split, cross_val_score,
62 from sklearn.metrics import make_scorer
63 from sklearn.metrics import classification_report, confusion_matrix, ac
64 from sklearn.metrics import f1_score, fbeta_score, r2_score, roc_auc_s
65
66 from sklearn.linear_model import LogisticRegression, SGDClassifier
67 from sklearn.neighbors import KNeighborsClassifier
68 from sklearn.svm import SVC
69 from sklearn.svm import SVC
70 from sklearn.naive_bayes import GaussianNB
71
72 from sklearn.metrics import precision_recall_curve
73 from sklearn.metrics import plot_precision_recall_curve, average_precis
74
75 from sklearn.metrics import plot_confusion_matrix
76

```

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
    - 15.3 Results
- in [2]: ▶ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

## 4 Load Data

```

1 df = pd.read_csv('bigml_59c28831336c6604c800002a.csv')
2
3 print (df.shape)
4

```

(3333, 21)

```
In [231]: # Load data
          2
          3 df.head(10)
```

**Out[231]:**

ce ail an	number vmail messages	total day minutes	total day calls	total day charge	...	total eve calls	total eve charge	total night minutes	total night calls	c
yes	25	265.1	110	45.07	...	99	16.78	244.7	91	
yes	26	161.6	123	27.47	...	103	16.62	254.4	103	
no	0	243.4	114	41.38	...	110	10.30	162.6	104	
no	0	299.4	71	50.90	...	88	5.26	196.9	89	
no	0	166.7	113	28.34	...	122	12.61	186.9	121	
no	0	223.4	98	37.98	...	101	18.75	203.9	118	
yes	24	218.2	88	37.09	...	108	29.62	212.6	118	
no	0	157.0	79	26.69	...	94	8.76	211.8	96	
no	0	184.5	97	31.37	...	80	29.89	215.8	90	
yes	37	258.6	84	43.96	...	111	18.87	326.4	97	

**Contents** ↗ ↘

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [232]: 1 df.dtypes

```
Out[232]: state          object
account length      int64
area code           int64
phone number        object
international plan  object
voice mail plan    object
number vmail messages  int64
total day minutes   float64
total day calls     int64
total day charge    float64
total eve minutes   float64
total eve calls     int64
total eve charge    float64
total night minutes  float64
total night calls   int64
total night charge  float64
total intl minutes  float64
total intl calls    int64
total intl charge   float64
customer service calls  int64
churn                bool
dtype: object
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

## 5 Data Exploration & Visualization

In [233]: #Data Visualization

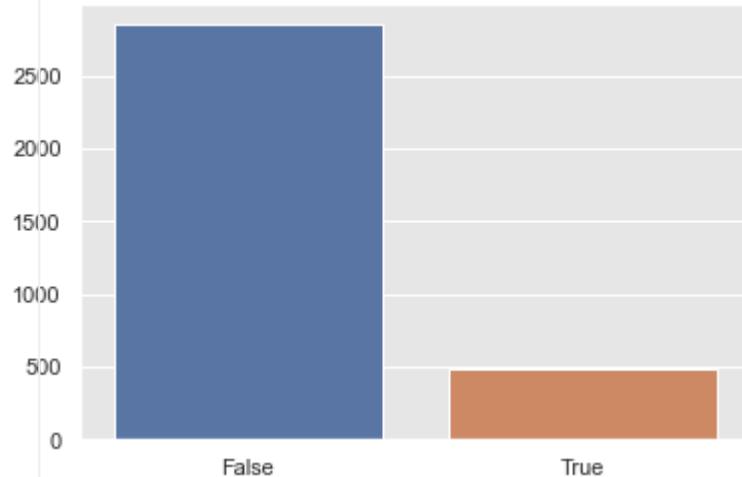
```

1
2
3 y = df["churn"].value_counts()
4 #print (y)
5 sns.barplot(y.index, y.values)
6
7 #sns.countplot(x ='churn', data = churn_data)
8 #plt.show()

```

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
  - ▼ 15.1.1 Hyp
  - 15.1.1.1
  - ▼ 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results



- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning

```

1 y_True = df["churn"][df["churn"] == True]
2 print ("Churn Percentage = "+str( (y_True.shape[0] / df["churn"].shape[0]) * 100 )

```

Churn Percentage = 14.491449144914492

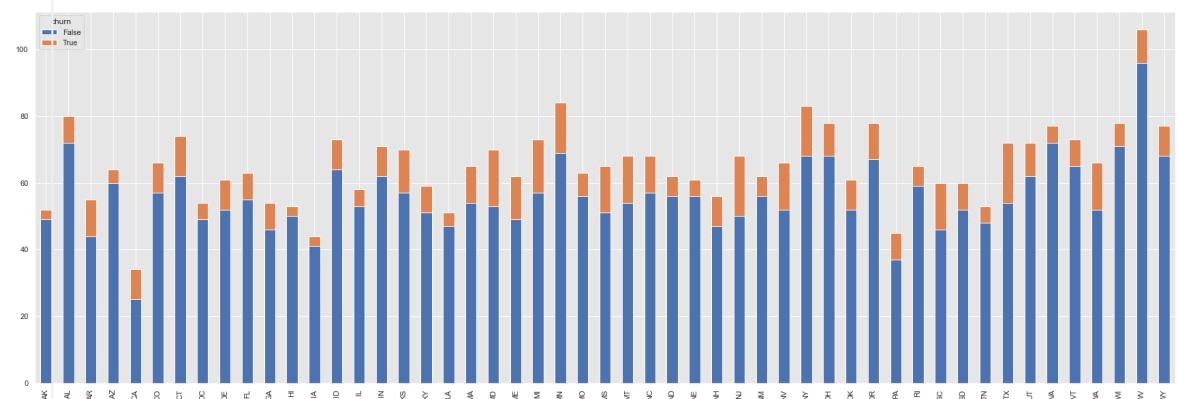
- ▼ 20 Data Interp
  - 20.1 Most In
  - 20.1.1 Top
  - 20.2 Top 25

```

1 #Churn by state
2 df.groupby(["state", "churn"]).size().unstack().plot(kind='bar', stacked=True)

```

Out[225]: <AxesSubplot:xlabel='state'>

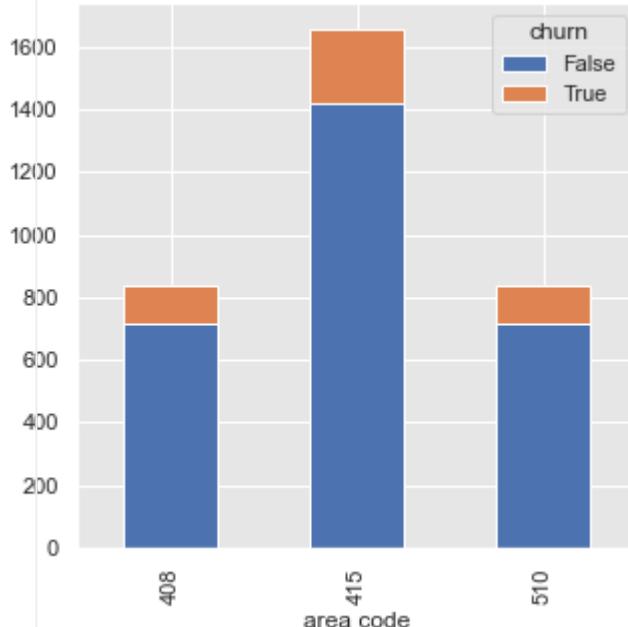


```
In [236]: #Churn by area code
           df.groupby(["area code", "churn"]).size().unstack().plot(kind='bar', st
```

Out[236]: <AxesSubplot:xlabel='area code'>

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



In [237]: 1 df.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
```

## Contents ⚙️

10.2	Results
▼ 11	XGBoost
11.1	Results
▼ 12	Gaussian N
12.1	Results
▼ 13	AdaBoostC
13.1	Results
▼ 14	Gradient Bo
14.1	Results
▼ 15	SVM Classi
▼ 15.1	Tuning
15.1.1	Hyp
15.1.1.1	
▼ 15.2	Results
15.2.1	SVI
15.2.2	SV
15.3	Results
▼ 16	Stacking
16.1	Results
▼ 17	SGDClassif
17.1	Results
▼ 18	Logistic Re
18.1	Results
▼ 19	SGDClassif
19.1	Tuning
19.2	Results
▼ 20	Data Interp
▼ 20.1	Most In
20.1.1	Top
▼ 20.2	Top Po
20.2.1	To
20.2.2	Toj
▼ 21	Conclusion
▼ 21.1	Increas
21.1.1	Total
21.1.2	Cus
21.1.3	Total
21.1.4	Total
▼ 21.2	Decrea
21.2.1	Voi
21.2.2	Inte
22	Recomend
23	Future Wor

In [238]: 1 df.dtypes

```
Out[238]: state          object
account length      int64
area code           int64
phone number        object
international plan  object
voice mail plan    object
number vmail messages  int64
total day minutes   float64
total day calls     int64
total day charge    float64
total eve minutes   float64
total eve calls     int64
total eve charge    float64
total night minutes  float64
total night calls   int64
total night charge  float64
total intl minutes  float64
total intl calls    int64
total intl charge   float64
customer service calls  int64
churn                bool
dtype: object
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results

In [239]: 1 #check for missing values  
2 df.isnull().sum()

```
Out[239]: 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
```

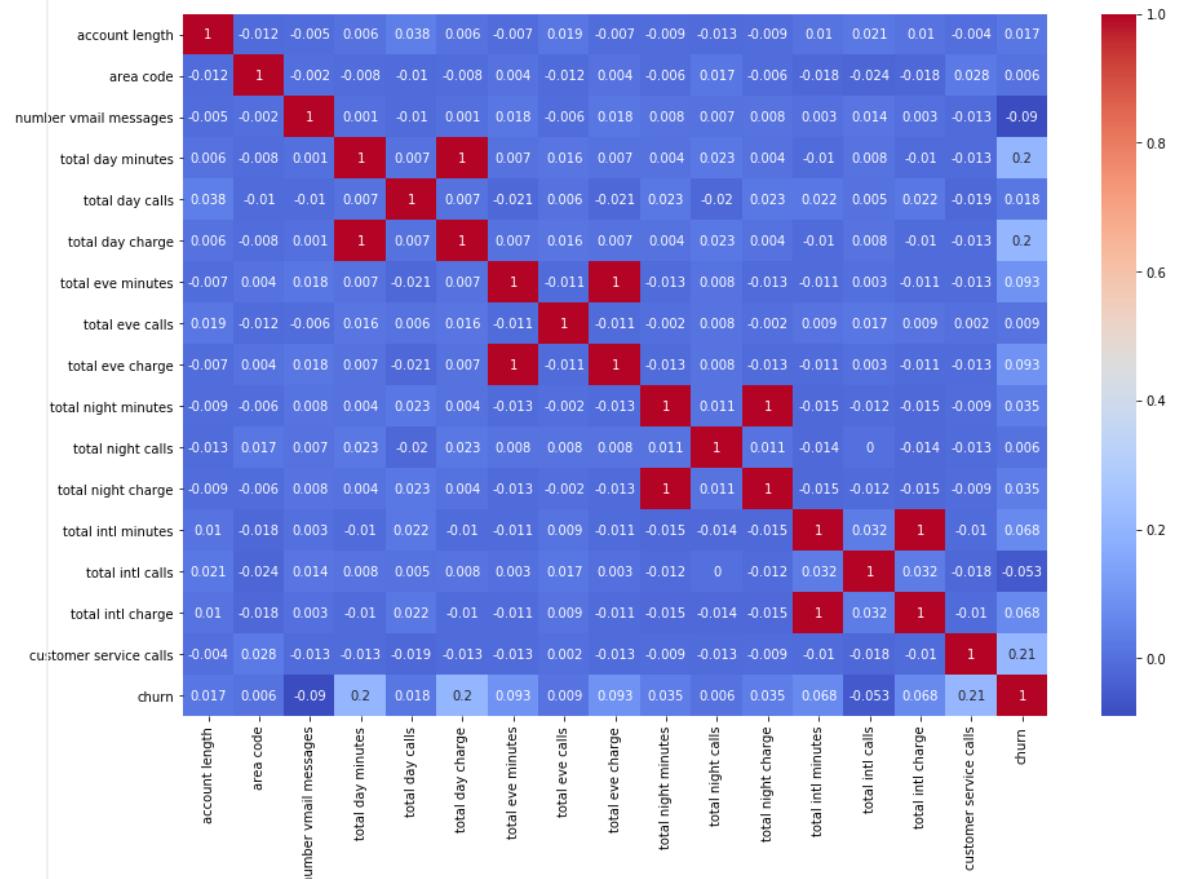
In [240]: # check for duplicated rows  
df.duplicated().sum()  
3

Out[240]: 0

## Contents

- 10.2 Results
- 11 XGBoost
- 11.1 Results
- 12 Gaussian N
- 12.1 Results
- 13 AdaBoostC
- 13.1 Results
- 14 Gradient Bc
- 14.1 Results
- 15 SVM Classi
- 15.1 Tuning
- 15.1.1 Hyp
- 15.1.1.1
- 15.1.2 Results
- 15.2.1 SVI
- 15.2.2 SV
- 15.3 Results
- 16 Stacking
- 16.1 Results
- 17 SGDClassif
- 17.1 Results
- 18 Logistic Re
- 18.1 Results
- 19 SGDClassif
- 19.1 Tuning
- 19.2 Results
- 20 Data Interp
- 20.1 Most In
- 20.1.1 Top
- 20.2 Top Po
- 20.2.1 To
- 20.2.2 To
- 21 Conclusion
- 21.1 Increas
- 21.1.1 Total
- 21.1.2 Cus
- 21.1.3 Total
- 21.1.4 Total
- 21.2 Decreas
- 21.2.1 Voic
- 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

#Correlation Analysis  
cor = df.corr()  
plt.figure(figsize=(15,10))  
sns.heatmap(cor.round(3), annot=True, cmap='coolwarm')  
plt.show()



In [267]: 1 df.dtypes

```
Out[267]: state          object
account length      int64
area code           int64
phone number        object
international plan  object
voice mail plan    object
number vmail messages  int64
total day minutes   float64
total day calls     int64
total day charge    float64
total eve minutes   float64
total eve calls     int64
total eve charge    float64
total night minutes  float64
total night calls   int64
total night charge  float64
total intl minutes  float64
total intl calls    int64
total intl charge   float64
customer service calls  int64
churn                bool
dtype: object
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [260]: 1 df.corr()

Out[260]:

	account length	area code	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes	total charge
account length	1.000000	-0.012463	-0.004628	0.006216	0.038470	0.006214	-0.006757	0.0
area code	-0.012463	1.000000	-0.001994	-0.008264	-0.009646	-0.008264	0.003580	-0.0
number vmail messages	-0.004628	-0.001994	1.000000	0.000778	-0.009548	0.000776	0.017562	-0.0
total day minutes	0.006216	-0.008264	0.000778	1.000000	0.006750	1.000000	0.007043	0.0
total day calls	0.038470	-0.009646	-0.009548	0.006750	1.000000	0.006753	-0.021451	0.0
total day charge	0.006214	-0.008264	0.000776	1.000000	0.006753	1.000000	0.007050	0.0
total eve minutes	-0.006757	0.003580	0.017562	0.007043	-0.021451	0.007050	1.000000	-0.0
total eve calls	0.019260	-0.011886	-0.005864	0.015769	0.006462	0.015769	-0.011430	1.0
total eve charge	-0.006745	0.003607	0.017578	0.007029	-0.021449	0.007036	1.000000	-0.0
total night minutes	-0.008955	-0.005825	0.007681	0.004323	0.022938	0.004324	-0.012584	-0.0
total night calls	-0.013176	0.016522	0.007123	0.022972	-0.019557	0.022972	0.007586	0.0
total night charge	-0.008960	-0.005845	0.007663	0.004300	0.022927	0.004301	-0.012593	-0.0
total intl minutes	0.009514	-0.018288	0.002856	-0.010155	0.021565	-0.010157	-0.011035	0.0
total intl calls	0.020661	-0.024179	0.013957	0.008033	0.004574	0.008032	0.002541	0.0
total intl charge	0.009546	-0.018395	0.002884	-0.010092	0.021666	-0.010094	-0.011067	0.0
customer service calls	-0.003796	0.027572	-0.013263	-0.013423	-0.018942	-0.013427	-0.012985	0.0
churn	0.016541	0.006174	-0.089728	0.205151	0.018459	0.205151	0.092796	0.0

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voi
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

# 6 Data Engineering

## Contents ↗ ⚙

In [268]:	1 df.groupby(['state', 'churn']).sum().head()
11 XGBoost 11.1 Results	
12 Gaussian N 12.1 Results	
13 AdaBoostC 13.1 Results	
14 Gradient Bo 14.1 Results	
15 SVM Classi 15.1 Tuning 15.1.1 Hyp 15.1.1.1	
15.2 Results 15.2.1 SVI 15.2.2 SV 15.3 Results	
16 Stacking 16.1 Results	
17 SGDClassif 17.1 Results	
18 Logistic Re 18.1 Results	
19 SGDClassif 19.1 Tuning 19.2 Results	
20 Data Interp 20.1 Most In 20.1.1 Top 20.2 Top Po 20.2.1 To 20.2.2 To	
21 Conclusion 21.1 Increas 21.1.1 Total 21.1.2 Cus 21.1.3 Total 21.1.4 Total 21.2 Decreas 21.2.1 Voi 21.2.2 Inte	
22 Recomend 23 Future Wor	
In [269]:	1 df['churn'] = df['churn'].astype('str')
In [179]:	1 south=["MD", "DE", "VA", "WV", "KY", "TN", "NC", "SC", "FL", "GA", "AL" 2 df['state_south']=df['state'].isin(south)
In [180]:	1 df['state_south'].value_counts()
Out[180]:	False 2226 True 1107 Name: state_south, dtype: int64
In [181]:	1 north=[ "AZ", "CA", "CO", "CT", "DC", "HI", "ID", "IL", "IN", "IA", "KS" 2 df['state_north']=df['state'].isin(north)

```
In [182]: 1 df['state_north'] = df['state_north'].astype('str')
2 df['state_south'] = df['state_south'].astype('str')
```

```
In [183]: 1 df.head()
```

## Contents ⚙️

10.2 Results
▼ 11 XGBoost
11.1 Results
▼ 12 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 Gradient Bo
14.1 Results
▼ 15 SVM Classi
▼ 15.1 Tuning
▼ 15.1.1 Hyp
15.1.1.1
▼ 15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results I tried to include state_north and state_south in the models but both features are insignificant. I
▼ 18 Logistic Re decided to remove state feature.
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
▼ 20.1 Most In
20.1.1 Top
▼ 20.2 Top Po
20.2.1 To
20.2.2 To
▼ 21 Conclusion
▼ 21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
▼ 21.2 Decreas
21.2.1 Voic
21.2.2 Inte
22 Recomend
23 Future Wor

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

5 rows × 23 columns

### 6.1.1 Comments

In [332]: 1 df.dtypes

```
Out[332]: state          object
account length      int64
area code           int64
phone number        object
international plan  object
voice mail plan    object
number vmail messages  int64
total day minutes   float64
total day calls     int64
total day charge    float64
total eve minutes   float64
total eve calls     int64
total eve charge    float64
total night minutes  float64
total night calls   int64
total night charge  float64
total intl minutes   float64
total intl calls    int64
total intl charge   float64
customer service calls  int64
churn              bool
state_south         bool
state_north         bool
dtype: object
```

## Contents ↗ ⚙

10.2	Results
▼ 11	XGBoost
11.1	Results
▼ 12	Gaussian N
12.1	Results
▼ 13	AdaBoostC
13.1	Results
▼ 14	Gradient Bo
14.1	Results
▼ 15	SVM Classi
▼ 15.1	Tuning
15.1.1	Hyp
15.1.1.1	
▼ 15.2	Results
15.2.1	SVI
15.2.2	SV
15.3	Results
▼ 16	Stacking
16.1	Results
▼ 17	SGDClassif
17.1	Results
▼ 18	Logistic Re
18.1	Results
▼ 19	In [241]: SGDClassif
19.1	Tuning
19.2	Results
▼ 20	Data Interp
▼ 20.1	Most In
In [241]: Top	
▼ 20.2	Top Po
20.2.1	To
20.2.2	To
▼ 21	Conclusion
21.1	Increas
21.1.1	Total
21.1.2	Cus
21.1.3	Total
21.1.4	Total
▼ 21.2	Decrea
In [241]: Voic	
21.2.2	Inte
22	Recomend
23	Future Wor

## 6.2 Drop unnecessary columns

In [241]: 1 #Remove phonenumer (is not important) ; redundant column  
2  
3 df=df.drop(['phone number'], axis=1)

In [242]: 1 df=df.drop(['state'], axis=1)

In [243]: 1 df=df.drop(['account length'], axis=1)

## 6.3 Drop higly correlated features

In [244]: 1 # We kept total intl charge  
2 df=df.drop(['total intl minutes'], axis=1)

```
In [245]: 1 # We kept total night charge
2 df=df.drop(['total night minutes'], axis=1)
```

```
In [246]: 1 #We kept total eve charge
2 df=df.drop(['total eve minutes'], axis=1)
3
```

## Contents ⚙

10.2 Results
▼ 11 XGBoost
11.1 Results
11.2 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 GradientBo
14.1 Results
▼ 15 SVM Classi
15.1 Tuning
15.1.1 Hyp
15.1.1.1
15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results
▼ 18 Logistic Re
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
20.1 Most In
20.1.1 Top
20.2 Top Po
20.2.1 To
20.2.2 Toj
▼ 21 Conclusion
21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
21.2 Decreas
21.2.1 Voi
21.2.2 Inte
22 Recomend
23 Future Wor

```
1 df.dtypes
```

state	object
account length	int64
area code	int64
phone number	object
international plan	object
voice mail plan	object
number vmail messages	int64
total day minutes	float64
total day calls	int64
total day charge	float64
total eve minutes	float64
total eve calls	int64
total eve charge	float64
total night minutes	float64
total night calls	int64
total night charge	float64
total intl minutes	float64
total intl calls	int64
total intl charge	float64
customer service calls	int64
churn	bool
	dtype: object

In [41]: 1 df.head()

Out[41]:

	state	account length	area code	international plan	voice mail plan	number vmail messages	total day minutes	total day calls	t ch
0	KS	128	415	no	yes	25	265.1	110	
1	OH	107	415	no	yes	26	161.6	123	
2	NJ	137	415	no	no	0	243.4	114	
3	OH	84	408	yes	no	0	299.4	71	
4	OK	75	415	yes	no	0	166.7	113	

## Contents ↗

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

```
In [145]: 1 for c in cat_col_names:
            2     print(df[c].value_counts())
```

0	3010
1	323

## Contents ⚙️

10.2 Results	23	84
▼ 11 XGBoost	34	83
11.1 Results	1	80
▼ 12 Gaussian N	35	78
12.1 Results	48	78
▼ 13 AdaBoostC	37	78
13.1 Results	50	77
▼ 14 Gradient Bo	45	77
14.1 Results	6	74
▼ 15 SVM Classi	22	73
▼ 15.1 Tuning	13	73
▼ 15.1.1 Hyp	46	73
15.1.1.1	43	72
▼ 15.2 Results	44	72
15.2.1 SVI	15	71
15.2.2 SV	16	70
15.3 Results	20	70
15.3 Results	27	68
▼ 16 Stacking	26	68
16.1 Results	31	68
▼ 17 SGDClassif	47	66
17.1 Results	33	66
▼ 18 Logistic Re	5	66
18.1 Results	19	65
▼ 19 SGDClassif	39	65
19.1 Tuning	25	65
19.2 Results	3	64
▼ 20 Data Interp	9	63
20.1 Most In	24	63
20.1.1 Top	32	62
20.1.2 Top	28	62
▼ 20.2 Top Po	21	62
20.2.1 To	36	61
20.2.2 To	8	61
▼ 21 Conclusion	29	61
▼ 21.1 Increas	40	60
21.1.1 Total	41	60
21.1.2 Cus	17	59
21.1.3 Tot	14	58
21.1.4 Tot	30	56
▼ 21.2 Decrease	2	55
21.2.1 Voic	7	54
21.2.2 Inte	10	54
22 Recomendat	42	53
23 Future Wor	11	53
	0	52
	18	51
	38	45

```

12      44
4       34
Name: state, dtype: int64
0      2411
1      922
Name: voice mail plan, dtype: int64

```

## Contents ⚙

<ul style="list-style-type: none"> <li>10.2 Results</li> <li>▼ 11 XGBoost           <ul style="list-style-type: none"> <li>In.1 [Results] ►</li> </ul> </li> <li>▼ 12 Gaussian N           <ul style="list-style-type: none"> <li>12.1 Results</li> </ul> </li> <li>▼ 13 AdaBoostC           <ul style="list-style-type: none"> <li>13.1 Results</li> </ul> </li> <li>▼ 14 Gradient Bc           <ul style="list-style-type: none"> <li>14.1 Results</li> </ul> </li> <li>▼ 15 SVM Classi           <ul style="list-style-type: none"> <li>▼ 15.1 Tuning               <ul style="list-style-type: none"> <li>▼ 15.1.1 Hyp                   <ul style="list-style-type: none"> <li>15.1.1.1</li> </ul> </li> </ul> </li> <li>▼ 15.2 Results               <ul style="list-style-type: none"> <li>15.2.1 SVI</li> <li>15.2.2 SV</li> <li>15.3 Results</li> </ul> </li> </ul> </li> <li>▼ 16 Stacking           <ul style="list-style-type: none"> <li>16.1 Results</li> </ul> </li> <li>▼ 17 SGDClassif           <ul style="list-style-type: none"> <li>17.1 Results</li> </ul> </li> <li>▼ 18 Logistic Re           <ul style="list-style-type: none"> <li>18.1 Results</li> </ul> </li> <li>▼ 19 SGDClassif           <ul style="list-style-type: none"> <li>19.1 Tuning</li> <li>19.2 Results</li> </ul> </li> <li>▼ 20 Data Interp           <ul style="list-style-type: none"> <li>▼ 20.1 Most In               <ul style="list-style-type: none"> <li>20.1.1 Top</li> </ul> </li> <li>▼ 20.2 Top Po               <ul style="list-style-type: none"> <li>20.2.1 To</li> <li>20.2.2 Toj</li> </ul> </li> </ul> </li> <li>▼ 21 Conclusion           <ul style="list-style-type: none"> <li>▼ 21.1 Increas               <ul style="list-style-type: none"> <li>21.1.1 Total</li> <li>21.1.2 Cus</li> <li>21.1.3 Total</li> <li>21.1.4 Total</li> </ul> </li> <li>▼ 21.2 Decrea               <ul style="list-style-type: none"> <li>21.2.1 Voic</li> <li>21.2.2 Inte</li> </ul> </li> </ul> </li> <li>22 Recomend</li> <li>23 Future Wor</li> </ul>	<h2>6.4 One Hot Encoder Categorical Columns</h2> <pre> 1 total_col_names=df.columns 2 3 #find numeric columns 4 5 num_cols=df._get_numeric_data().columns 6 7 # Remove "object"-type features from X 8 cont_features = [col for col in df.columns if df[col].dtype in [np.float64, np.int64]] 9 10 # Remove "object"-type features from X_train and X_test 11 df_cont =df.loc[:, cont_features] 12 13 #Category column 14 15 cat_col_names=list(set(total_col_names)-set(num_cols)) 16 17 # Create X_cat which contains only the categorical variables 18 features_cat = [col for col in df.columns if df[col].dtype in [np.object]] 19 df_cat = df.loc[:, features_cat] 20 21 # OneHotEncode categorical variables 22 ohe = OneHotEncoder(handle_unknown='ignore') 23 24 df_cat_ohe = ohe.fit_transform(df_cat) 25 26 # Convert these columns into a DataFrame 27 columns = ohe.get_feature_names(input_features=df_cat.columns) 28 df_cat = pd.DataFrame(df_cat_ohe.todense(), columns=columns) 29 </pre>
--	---

In [140]: 1 df\_cont.head()

Out[140]:

	area code	number vmail messages	total day calls	total day charge	total eve calls	total eve charge	total night calls	total night charge	total intl calls	churn
0	415	25	110	45.07	99	16.78	91	11.01	3	
1	415	26	123	27.47	103	16.62	103	11.45	3	
2	415	0	114	41.38	110	10.30	104	7.32	5	
3	408	0	71	50.90	88	5.26	89	8.86	7	
4	415	0	113	28.34	122	12.61	121	8.41	3	

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
- ▼ 19.2 19.1
  - 19.2.1 Results
- ▼ 20 Data Interp
  - 20.1 Most In
    - 20.1.1 Top
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [49]: 1 df\_cont.columns

Out[49]: Index(['area code', 'number vmail messages', 'total day calls', 'total day charge', 'total eve calls', 'total eve charge', 'total night calls', 'total night charge', 'total intl calls', 'customer service calls', 'churn'], dtype='object')

In [141]: 1 df = pd.concat([pd.DataFrame(df\_cont), df\_cat], axis=1)

In [196]: 1 df.head()

Out[196]:

	account length	area code	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes	total eve calls	total eve charge
0	128	415	25	265.1	110	45.07	197.4	99	16
1	107	415	26	161.6	123	27.47	195.5	103	16
2	137	415	0	243.4	114	41.38	121.2	110	10
3	84	408	0	299.4	71	50.90	61.9	88	5
4	75	415	0	166.7	113	28.34	148.3	122	12

5 rows × 72 columns

In [349]: 1 df.columns

Out[349]:

```
Index(['area code', 'number vmail messages', 'total day calls',
       'total day charge', 'total eve calls', 'total eve charge',
       'total night calls', 'total night charge', 'total intl calls',
       'total intl charge', 'customer service calls', 'churn',
       'international plan_no', 'international plan_yes', 'voice mail plan_no',
       'voice mail plan_yes', 'state_south_False', 'state_south_True',
       'state_north_False', 'state_north_True'],
      dtype='object')
```

In [350]: 1 df.dtypes

Out[350]:

account length	int64
area code	int64
number vmail messages	int64
total day minutes	float64
total day calls	int64
...	
state_WY	float64
international plan_no	float64
international plan_yes	float64
voice mail plan_no	float64
voice mail plan_yes	float64

Length: 72, dtype: object

## Contents ↗

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp

- In [349]: 1 df.columns
- 15.2 Results
- 15.2.1 SVI
- 15.2.2 SV
- 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - In [350]: 1 df.dtypes
  - 19.2 Results
- ▼ 20 Data Interp
  - 20.1 Most In
  - 20.1.1 Top 3
  - 20.2 Top Po
  - 20.2.1 To
  - 20.2.2 Toj
- ▼ 21 Conclusion
  - 21.1 Increases
  - 21.1.1 Total
  - 21.1.2 Cus
  - 21.1.3 Total
  - 21.1.4 Total
  - 21.2 Decreases
  - 21.2.1 Voic
  - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [50]: 1 df.shape  
2

Out[50]: (3333, 16)

## Contents

10.2 Results

▼ 11 XGBoost

Out[51]: (3333, 16)

11.1 Results

▼ 12 Gaussian N

12.1 Results

▼ 13 AdaBoostC

13.1 Results

▼ 14 GradientBo

14.1 Results

▼ 15 SVM Classi

15.1 Tuning

15.1.1 Hyp

15.1.1.1

▼ In [142]:

1 df.shape

```
1 ## separate dependent and independent variables
2 X = df.drop(['churn'], axis=1)
3 y = df['churn']
```

15.2 Results

15.2.1 SVI

15.2.2 SV

15.3 Results

▼ 16 Stacking

16.1 Results

▼ 17 SGDClassif

17.1 Results

▼ 18 Logistic Re

18.1 Results

▼ 19 SGDClassif

## 6.5 Train test split

```
1 ## splitting whole dataset into train and test dataset
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
3 print(f"Size Of The Train Dataset :- {len(X_train)}")
4 print(f"Size Of The Test Dataset :- {len(X_test)}")
```

Size Of The Train Dataset :- 2666

Size Of The Test Dataset :- 667

## 6.6 Standarize Features

```
1 scaler = StandardScaler()
2 scaled_x_train = scaler.fit_transform(X_train)
3 scaled_X_test = scaler.transform(X_test)
4 scaled_X_train = pd.DataFrame(scaled_x_train, columns=X_train.columns)
```

▼ In [144]:

## 6.7 Building analysis functions

▼ 21 Conclusion

21.1 Increases

21.1.1 Total

21.1.2 Cus

21.1.3 Total

21.1.4 Total

▼ 21.2 Decrease

21.2.1 Voic

21.2.2 Inter

22 Recomend

23 Future Wor

In [145]:

```

1 #Function to draw the ROC curve and to calculate the area under the cur
2 from sklearn.metrics import roc_curve, auc
3 from sklearn.svm import SVC
4
5 def roc_curve_and_auc(clf, X_train, X_test, y_train, y_test):
6
7     # Calculate the probability scores of each point in the training se
8     y_train_score = clf.fit(X_train, y_train).decision_function(X_train)
9
10    # Calculate the fpr, tpr, and thresholds for the training set
11    train_fpr, train_tpr, thresholds = roc_curve(y_train, y_train_score)
12
13    # Calculate the probability scores of each point in the test set
14    y_test_score = clf.decision_function(X_test)
15
16    # Calculate the fpr, tpr, and thresholds for the test set
17    test_fpr, test_tpr, test_thresholds = roc_curve(y_test, y_test_score)
18
19    # ROC curve for training set
20    plt.figure(figsize=(10, 8))
21    lw = 2
22    plt.plot(train_fpr, train_tpr, color='darkorange',
23              lw=lw, label='Train ROC curve')
24    plt.plot(test_fpr, test_tpr, color='blue',
25              lw=lw, label='Test 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(
34        'Receiver operating characteristic (ROC) Curve for Training and'
35        'Testing Data')
36    plt.legend(loc='lower right')
37    plt.show()
38    # Print the area under the roc curve
39    print('Training AUC: {}'.format(round(auc(train_fpr, train_tpr), 5)))
40    print('Testing AUC: {}'.format(round(auc(test_fpr, test_tpr), 5)))

```

## Contents ⚙️

10.2 Results	
11 XGBoost	
11.1 Results	
12 Gaussian N	
12.1 Results	
13 AdaBoostC	
13.1 Results	
14 Gradient Bo	
14.1 Results	
15 SVM Classi	
15.1 Tuning	
15.1.1 Hyp	
15.1.1.1	
15.2 Results	
15.2.1 SVI	
15.2.2 SV	
15.3 Results	
16 Stacking	
16.1 Results	
17 SGDClassif	
17.1 Results	
18 Logistic Re	
18.1 Results	
19 SGDClassif	
19.1 Tuning	
19.2 Results	
20 Data Interp	
20.1 Most In	
20.1.1 Top	
20.2 Top Po	
20.2.1 To	
20.2.2 Toj	
21 Conclusion	
21.1 Increases	
21.1.1 Total	
21.1.2 Cus	
21.1.3 Total	
21.1.4 Total	
21.2 Decrease	
21.2.1 Voic	
21.2.2 Inte	
22 Recomend	
23 Future Wor	

In [146]:

```

1 #Plot coefficients
2 def plot_coefficients(clf):
3     weights_clf = pd.Series(clf.coef_[0], index=X.columns.values)
4     weights_clf.sort_values(inplace=True)
5     plt.figure(figsize=(15, 6))
6     plt.xticks(rotation=90)
7     features = plt.bar(weights_clf.index, weights_clf.values)

```

In [147]:

```

1 #Function to draw precision recall curve
2 def plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train):
3     fig = plt.figure(figsize = (8, 8))
4     lw = 3
5     no_skill = len(y[y==1]) / len(y)
6     # plot the no skill precision-recall curve
7     plt.plot([0, 1], [no_skill, no_skill], linestyle='--', lw = lw, label='No Skill')
8     # calculate model precision-recall curve
9     precision, recall, _ = precision_recall_curve(y_test, y_pred_test)
10    train_precision, train_recall, _ = precision_recall_curve(y_train, y_pred_train)
11    # plot the model precision-recall curve
12    plt.plot(recall, precision, marker='.', lw = lw, label='Test Set')
13    plt.plot(train_recall,train_precision, marker='.',lw=lw, label='Training Set')
14    # axis labels
15    plt.xlabel('Recall')
16    plt.ylabel('Precision')
17    # show the legend
18    plt.legend()
19
20    # show the plot
21    plt.show()

```

## Contents ↗

10.2	Results
▼ 11	XGBoost
11.1	Results
▼ 12	Gaussian NB
12.1	Results
▼ 13	AdaBoostC
13.1	Results
▼ 14	Gradient Bo
14.1	Results
▼ 15	SVM Classi
▼ 15.1	Tuning
15.1.1	Hyp
15.1.1.1	
▼ 15.2	Results
15.2.1	SVI
15.2.2	SV
15.3	Results
▼ 16	Stacking
16.1	Results
▼ 17	SGDClassif
17.1	Results
▼ 18	Logistic Re
18.1	Results
▼ 19	SGDClassif
19.1	Tuning
19.2	Results
▼ In [148]:	Data Interp
▼ 20.1	Most Inf
20.1.1	Top
20.1.2	Top Po
20.2.1	To
20.2.2	To
▼ 21	Conclusion
▼ 21.1	Increas
21.1.1	Total
21.1.2	Cus
21.1.3	Total
21.1.4	Total
▼ 21.2	Decreas
21.2.1	Voi
21.2.2	Inte
22	Recomend
23	Future Wor

In [148]:

```

1 #To plot coefficients
2 def plot_feature_importances(model):
3     n_features = X_train.shape[1]
4     plt.figure(figsize=(8, 8))
5     plt.barh(range(n_features), model.feature_importances_, align='center')
6     plt.yticks(np.arange(n_features), X_train.columns.values)
7     plt.xlabel('Feature importance')
8     plt.ylabel('Feature')

```

In [149]:

```

1
2 #Plot confusion matrix
3 def confusion_matrix_df(y_test, y_pred_test):
4     confusion_mat = confusion_matrix(y_test, y_pred_test, labels=[0, 1])
5     _row = confusion_mat.sum(axis=0)
6     _col = [np.nan] + list(confusion_mat.sum(axis=1)) + [_row]
7     con_df = pd.DataFrame({})
8     con_df["Predicted"] = ["Actual"] + ["No Churn", "Churn"] + ["All"]
9     for label, idx in {"No Churn": 0, "Churn": 1}.items():
10        temp = [np.nan] + list(confusion_mat[:, idx]) + [_row[idx]]
11        con_df[label] = temp
12
13    con_df["All"] = _col
14    return con_df

```

```
In [150]: #Function to print the metrics: Confusion Matrix, Classification Report
def model_evaluation(X_train, X_test, y_train, y_test, y_pred_train, y_
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - 20.1 Most In
    - 20.1.1 Top
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

# 7 Logistic Regression

## 7.1 Optimization and Tuning

In [61]:

```

1 # Now let's compare a few different regularization performances on the
2 weights = [None, 'balanced', {1:2, 0:1}, {1:10, 0:1}, {1:100, 0:1}, {1:
3 names = ['None', 'Balanced', '2 to 1', '10 to 1', '100 to 1', '1000 to
4 colors = sns.color_palette('Set2')
5
6 plt.figure(figsize=(10,8))
7
8 for n, weight in enumerate(weights):
9     # Fit a model
10    logreg = LogisticRegression(fit_intercept=False, C=1e20, class_weight=
11        weight)
12    model_log = logreg.fit(scaled_X_train, y_train)
13    print(model_log)

14    # Predict
15    y_hat_test = logreg.predict(scaled_X_test)

16    y_score = logreg.fit(scaled_X_train, y_train).decision_function(sca
17
18    fpr, tpr, thresholds = roc_curve(y_test, y_score)

19    print('AUC for {}: {}'.format(names[n], auc(fpr, tpr)))
20    print('-----')
21    lw = 2
22
23    plt.plot(fpr, tpr, color=colors[n],
24              lw=lw, label='ROC curve {}'.format(names[n]))
25
26
27    plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
28    plt.xlim([0.0, 1.0])
29    plt.ylim([0.0, 1.05])
30
31    plt.yticks([i/20.0 for i in range(21)])
32    plt.xticks([i/20.0 for i in range(21)])
33    plt.xlabel('False Positive Rate')
34    plt.ylabel('True Positive Rate')
35    plt.title('Receiver operating characteristic (ROC) Curve')
36    plt.legend(loc='lower right')
37    plt.show()

```

```
LogisticRegression(C=1e+20, fit_intercept=False)
```

```
AUC for None: 0.81423403870501
```

---

```
LogisticRegression(C=1e+20, class_weight='balanced', fit_intercept=False)
```

```
AUC for Balanced: 0.814794718755652
```

---

```
LogisticRegression(C=1e+20, class_weight={0: 1, 1: 2}, fit_intercept=False)
```

```
AUC for 2 to 1: 0.8150479291011034
```

---

```
LogisticRegression(C=1e+20, class_weight={0: 1, 1: 10}, fit_intercept=False)
```

```
AUC for 10 to 1: 0.8137457044673538
```

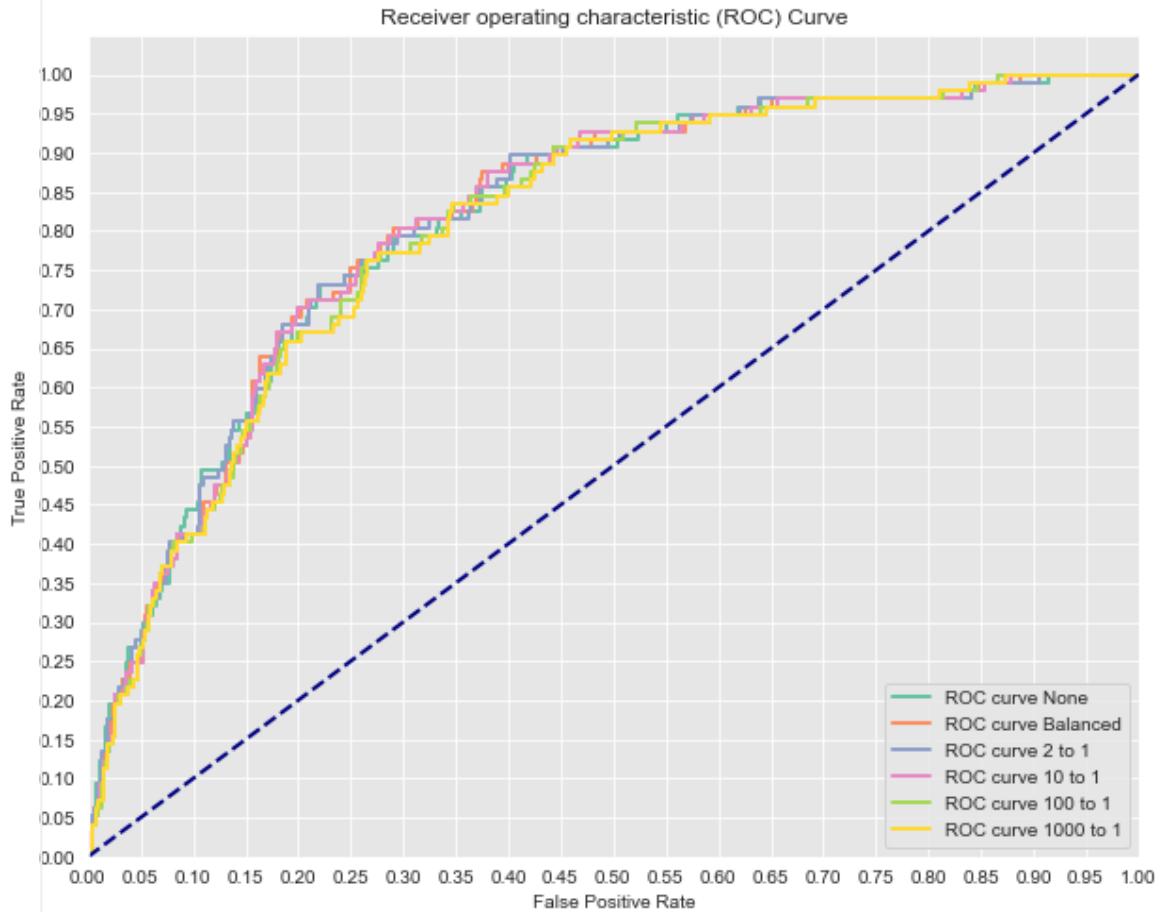
## Contents

10.2	Results
▼ 11	XGBoost
11.1	Results
▼ 12	Gaussian N
12.1	Results
▼ 13	AdaBoostC
13.1	Results
▼ 14	Gradient Bo
14.1	Results
▼ 15	SVM Classi
▼ 15.1	Tuning
15.1.1	Hyp
15.1.1.1	
▼ 15.2	Results
15.2.1	SVI
15.2.2	SV
15.3	Results
▼ 16	Stacking
16.1	Results
▼ 17	SGDClassif
17.1	Results
▼ 18	Logistic Re
18.1	Results
▼ 19	SGDClassif
19.1	Tuning
19.2	Results
▼ 20	Data Interp
▼ 20.1	Most In
20.1.1	Top
▼ 20.2	Top Po
20.2.1	To
20.2.2	Toj
▼ 21	Conclusion
▼ 21.1	Increas
21.1.1	Total
21.1.2	Cus
21.1.3	Total
21.1.4	Total
▼ 21.2	Decreas
21.2.1	Voi
21.2.2	Inte
22	Recomend
23	Future Wor

```
LogisticRegression(C=1e+20, class_weight={0: 1, 1: 100}, fit_intercept=False)
AUC for 100 to 1: 0.8077590884427563
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
      - 15.3 Results
  - ▼ 16 Stacking
    - 16.1 Results
  - ▼ 17 SGDClassif
    - 17.1 Results
  - ▼ 18 Logistic Re
    - 18.1 Results
  - ▼ 19 SGDClassif
    - 19.1 Tuning
    - 19.2 Results
  - ▼ 20 Data Interp
    - ▼ 20.1 Most In
      - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 Toj
  - ▼ 21 Conclusion
    - ▼ 21.1 Increas
      - 21.1.1 Total
      - 21.1.2 Cus
      - 21.1.3 Total
      - 21.1.4 Total
    - ▼ 21.2 Decreas
      - 21.2.1 Voic
      - 21.2.2 Inte
  - 22 Recomend
  - 23 Future Wor



In [362]: 1 ! pip install imbalanced-learn

```
Requirement already satisfied: imbalanced-learn in c:\users\mirnamamaranda\anaconda3\anaconda\lib\site-packages (0.8.0)
Requirement already satisfied: numpy>=1.13.3 in c:\users\mirnamamaranda\anaconda3\anaconda\lib\site-packages (from imbalanced-learn) (1.19.2)
Requirement already satisfied: scipy>=0.19.1 in c:\users\mirnamamaranda\anaconda3\anaconda\lib\site-packages (from imbalanced-learn) (1.5.2)
Requirement already satisfied: joblib>=0.11 in c:\users\mirnamamaranda\anaconda3\anaconda\lib\site-packages (from imbalanced-learn) (0.17.0)
Requirement already satisfied: scikit-learn>=0.24 in c:\users\mirnamamaranda\anaconda3\anaconda\lib\site-packages (from imbalanced-learn) (0.24.1)
Requirement already satisfied: threadpoolctl>=2.0.0 in c:\users\mirnamamaranda\anaconda3\anaconda\lib\site-packages (from scikit-learn>=0.24->imbalanced-learn) (2.1.0)
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [63]:

```

1 # Now let's compare a few different regularization performances on the
2 C_param_range = [0.001, 0.01, 0.1, 1, 10, 100]
3 names = [0.001, 0.01, 0.1, 1, 10, 100]
4 colors = sns.color_palette('Set2')
5
6 plt.figure(figsize=(10, 8))
7
8 for n, c in enumerate(C_param_range):
9     # Fit a model
10    logreg = LogisticRegression(fit_intercept=False, C=c, solver='liblinear')
11    model_log = logreg.fit(scaled_X_train, y_train)
12    print(model_log)
13    # Preview model params
14
15    # Predict
16    y_hat_test = logreg.predict(scaled_X_test)
17
18    y_score = logreg.fit(scaled_X_train, y_train).decision_function(sca
19
20    fpr, tpr, thresholds = roc_curve(y_test, y_score)
21
22    print('AUC for {}: {}'.format(names[n], auc(fpr, tpr)))
23    print('-----')
24    lw = 2
25    plt.plot(fpr, tpr, color=colors[n],
26              lw=lw, label='ROC curve Normalization Weight: {}'.format(r
27
28    plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
29    plt.xlim([0.0, 1.0])
30    plt.ylim([0.0, 1.05])
31
32    plt.yticks([i/20.0 for i in range(21)])
33    plt.xticks([i/20.0 for i in range(21)])
34    plt.xlabel('False Positive Rate')
35    plt.ylabel('True Positive Rate')
36    plt.title('Receiver operating characteristic (ROC) Curve')
37    plt.legend(loc='lower right')
38    plt.show()

```

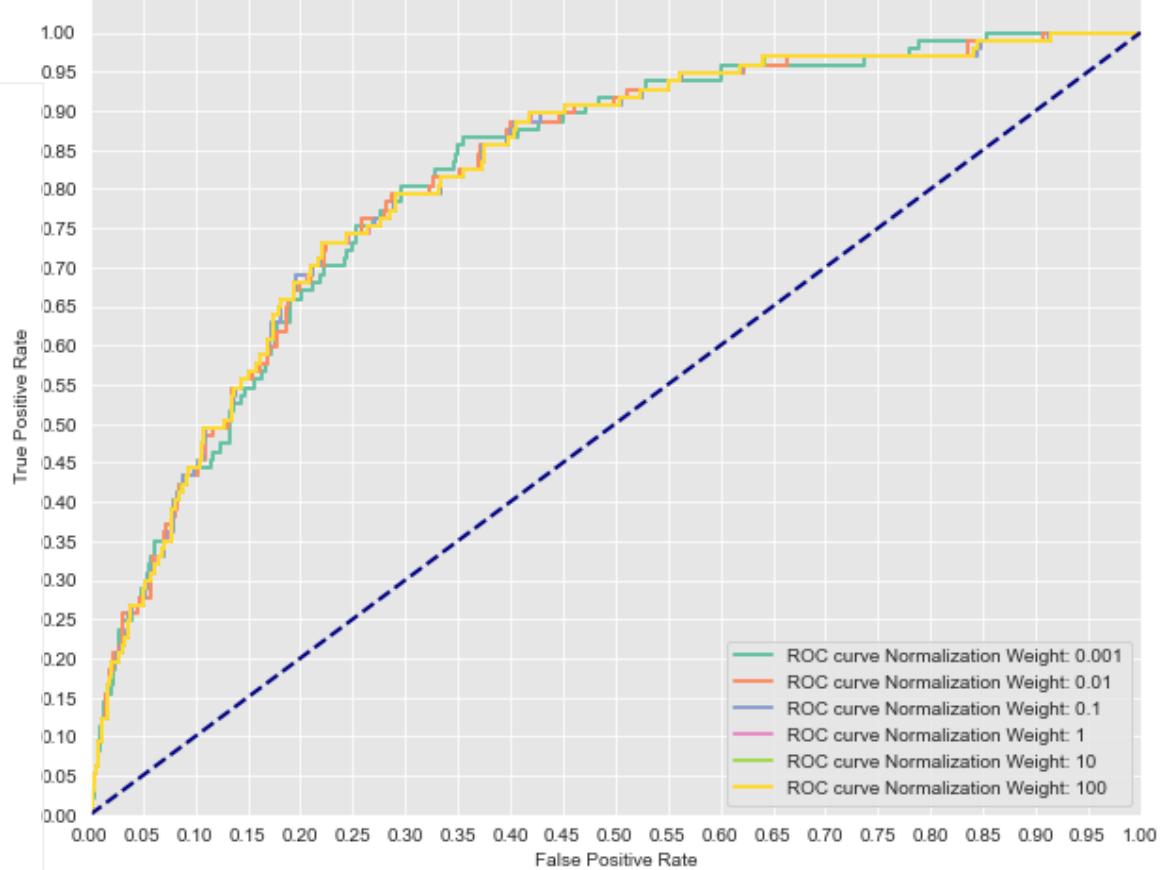
```

LogisticRegression(C=0.001, fit_intercept=False, solver='liblinear')
AUC for 0.001: 0.8137999638270935
-----
LogisticRegression(C=0.01, fit_intercept=False, solver='liblinear')
AUC for 0.01: 0.8139989148128052
-----
LogisticRegression(C=0.1, fit_intercept=False, solver='liblinear')
AUC for 0.1: 0.8140893470790378
-----
LogisticRegression(C=1, fit_intercept=False, solver='liblinear')
AUC for 1: 0.8142159522517635
-----
LogisticRegression(C=10, fit_intercept=False, solver='liblinear')
AUC for 10: 0.8142340387050099
-----
LogisticRegression(C=100, fit_intercept=False, solver='liblinear')

```

AUC for 100: 0.81423403870501

Receiver operating characteristic (ROC) Curve



## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif Model
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

## 7.2 Best performing parameters for Logistic Regression Model

In [64]:

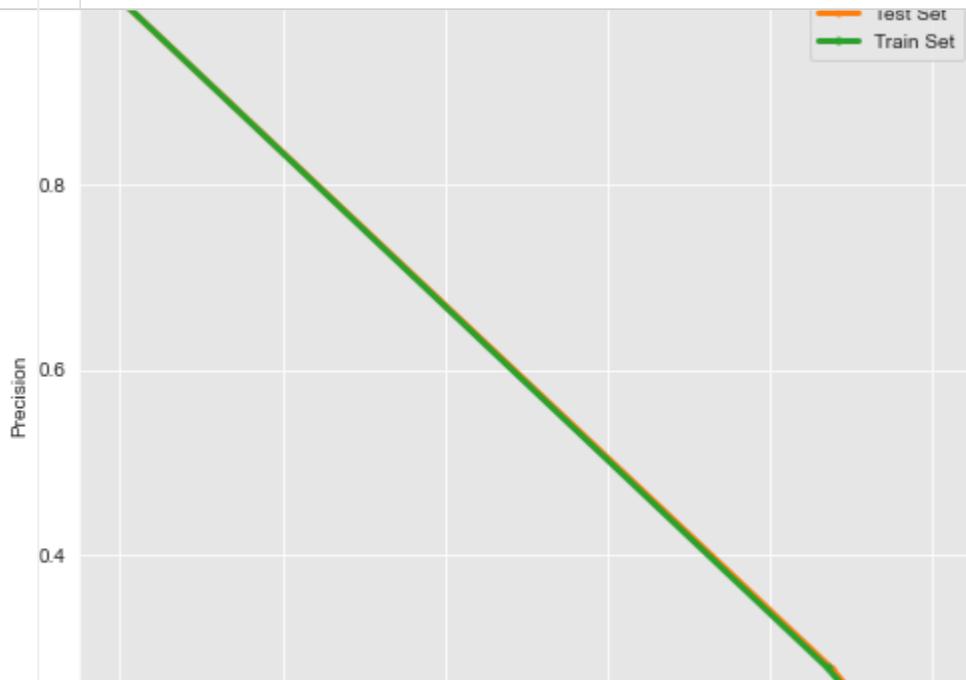
```

1 clf = LogisticRegression(fit_intercept=False, C=100,
                           solver='liblinear', class_weight='balanced')
2
3 clf.fit(scaled_X_train, y_train)
4 y_pred_train = clf.predict(scaled_X_train)
5 y_pred_test = clf.predict(scaled_X_test)
6
7
8 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred)
9
10 roc_curve_and_auc(clf, scaled_X_train, scaled_X_test, y_train, y_test)
11
12 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)

```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
      - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning



In [65]:

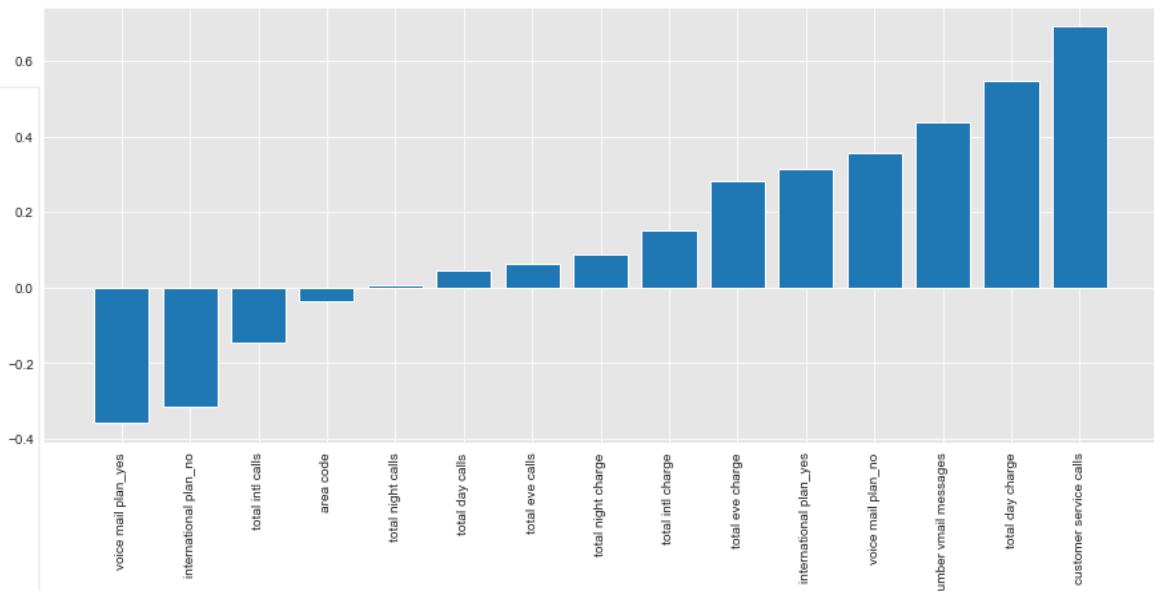
```
1 clf.coef_.round(2)
```

Out[65]:

```
array([[-0.04,  0.44,  0.05,  0.55,  0.06,  0.28,  0.01,  0.09, -0.15,
       0.15,  0.69, -0.32,  0.32,  0.36, -0.36]])
```

- ▼ 20 Data Interp
  - 20.1 Most In
  - 20.1.1 Top
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [66]: 1 | plot\_coefficients(clf)



## Contents ⚙️

10.2	Results
▼ 11	XGBoost
11.1	Results
▼ 12	Gaussian N
12.1	Results
▼ 13	AdaBoostC
13.1	Results
▼ 14	Gradient Bo
14.1	Results
▼ 15	SVM Classi
▼ 15.1	Tuning
15.1.1	Hyp
15.1.1.1	
▼ 15.2	Results
15.2.1	SVI
15.2.2	SV
15.3	Results
▼ 16	Stacking
16.1	Results
▼ 17	SGDClassif
17.1	Results
▼ 18	Logistic Re
18.1	Results
▼ 19	SGDClassif
19.1	Tuning
19.2	Results
▼ 20	Data Interp
▼ 20.1	Most In
20.1.1	Top
20.1.2	Total
20.1.3	Top
▼ 20.2	Top Po
20.2.1	To
20.2.2	Total
▼ 21	Conclusion
21.1	Increas
21.1.1	Total
21.1.2	Cus
21.1.3	Total
21.1.4	Total
▼ 21.2	Decrea
21.2.1	Voi
21.2.2	Inte
22	Recomend
23	Future Wor

## 7.3 Results Comments

- The mean cross validation is 76.87%. We have a precision of 0.28 and recall of 0.88 for the test set. We don't have overfitting. It means we are able to detect the churn customers 88% but 28% of the customers we say 'churn' was true.
- Most important positive churn factors:
  - 1)Customer service calls
  - 2)Total day charge
  - 3)Number voicemail messages
  - 4)Voice mail plan \_no
  - 5)International plan\_yes

6)Total evening charge

Most important negative churn features:

1)Voice mail plan \_yes

2)International plan\_no

## Contents ↻ ⚙

10.2 Results 3)Total International call

▼ 11 XGBoost 4)Area code

11.1 Results

▼ 12 Gaussian N

12.1 Results

▼ 13 AdaBoostC

13.1 Results

▼ 14 Gradient Bo

In [67]: 14.1 Results

▼ 15 SVM Classi

15.1 Tuning

▼ 15.1.1 Hyp

15.1.1.1

▼ 15.2 Results

15.2.1 SVI

15.2.2 SV

15.3 Results

▼ 16 Stacking

16.1 Results

▼ 17 SGDClassif

17.1 Results

▼ 18 Logistic Re

18.1 Results

▼ 19 SGDClassif

19.1 Tuning

19.2 Results

▼ 20 Data Interp

▼ 20.1 [68]: In

20.1.1 Top

▼ 20.2 Top Po

20.2.1 To

20.2.2 Toj

▼ 21 Conclusion

▼ 21.1 Increas

21.1.1 Total

21.1.2 Cus

21.1.3 Total

21.1.4 Total

▼ 21.2 Decrea

21.2.1 Voi

21.2.2 Inte

22 Recomend

23 Future Wor

## 8 Descion Tree Classifier

```
1 from sklearn.tree import DecisionTreeClassifier
2
3
4 tree_clf = DecisionTreeClassifier(criterion='gini', max_depth=5)
5 tree_clf.fit(scaled_X_train, y_train)
6
7 # Model Performance
8 # Test set predictions
9 pred = tree_clf.predict(scaled_X_test)
10
11
12
13 tree_clf_score = cross_val_score(tree_clf, scaled_X_train, y_train, cv=5)
14 mean_tree_clf_score = np.mean(tree_clf_score)
15
16 print(f"Mean Cross Validation Score: {mean_tree_clf_score :.2%}")
```

Mean Cross Validation Score: 93.96%

```
1 print('Training r^2:', tree_clf.score(scaled_X_train, y_train))
2 print('Test r^2:', tree_clf.score(scaled_X_test, y_test))
3
4 print(classification_report(y_test,pred))
```

Training r^2: 0.9538634658664666

Test r^2: 0.9175412293853074

	precision	recall	f1-score	support
False	0.94	0.96	0.95	570
True	0.76	0.64	0.69	97
accuracy			0.92	667
macro avg	0.85	0.80	0.82	667
weighted avg	0.91	0.92	0.91	667

## 8.1 Hyperparameter Tuning for decision tree using Gridsearch

In [69]:

```

1 dt_param_grid = {
2     'criterion': ['gini', 'entropy'],
3     'max_depth': [None, 2, 3, 4, 5, 6],
4     'min_samples_split': [2, 5, 10],
5     'min_samples_leaf': [1, 2, 3, 4, 5, 6]
6 }
7
8 num_decision_trees = 3 * 2 * 6 * 3 * 6
9 print(f"Grid Search will have to search through {num_decision_trees} di

```

### Contents ↗

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
      - 15.3 Results
  - ▼ 16 Stacking
    - 16.1 Results
  - ▼ 17 SGDClassif
    - 17.1 Results
  - ▼ 18 Logistic Re
    - 18.1 Results
  - ▼ 19 SGDClassif
    - 19.1 Tuning
    - 19.2 Results
  - ▼ 20 Data Interp
    - ▼ 20.1 Most In
      - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 Toj
  - ▼ 21 Conclusion[70]:
    - ▼ 21.1 Increas
      - 21.1.1 Total
      - 21.1.2 Cus
      - 21.1.3 Total
      - 21.1.4 Total
    - ▼ 21.2 Decreas
      - 21.2.1 Voi
      - 21.2.2 Inte
  - 22 Recomend:
  - 23 Future Wor

Grid Search will have to search through 648 different permutations.

In [70]:

```

1 from sklearn.model_selection import train_test_split, GridSearchCV, cross_val_score
2 # Instantiate GridSearchCV
3 dt_grid_search = GridSearchCV(tree_clf, dt_param_grid, cv=10, return_train_score=True)
4
5 # Fit to the data
6 dt_grid_search.fit(scaled_X_train, y_train)
7
8 # Examine best parameters
9 dt_grid_search.best_params_
10 # Mean training score
11 dt_gs_training_score = np.mean(dt_grid_search.cv_results_['mean_train_score'])
12
13 # Mean test score
14 dt_gs_testing_score = dt_grid_search.score(scaled_X_test, y_test)
15
16 print(f"Mean Training Score: {dt_gs_training_score :.2%}")
17 print(f"Mean Test Score: {dt_gs_testing_score :.2%}")
18 print("Best Parameter Combination Found During Grid Search:")
19 dt_grid_search.best_params_

```

Mean Training Score: 93.67%

Mean Test Score: 92.80%

Best Parameter Combination Found During Grid Search:

```
{'criterion': 'gini',
 'max_depth': 6,
 'min_samples_leaf': 4,
 'min_samples_split': 10}
```

### 8.1.1 Select Best Parameters for Descion tree classifier

```
In [71]: N
1 dt_param_grid2 = {
2     'criterion': ['gini'],
3     'max_depth': [ 6],
4     'min_samples_split': [4],
5     'min_samples_leaf': [10]
6 }
7
8 dt_grid_search2 = GridSearchCV(tree_clf, dt_param_grid2, cv=10, return_
9
10 # Fit to the data
11
12 clf=dt_grid_search2 .fit(scaled_X_train, y_train)
13
14
15 #Examine best parameters
16
17 # Mean training score
18 dt_gs_training_score2 = np.mean(clf.cv_results_['mean_train_score'])
19
20 # Mean test score
21 dt_gs_testing_score2 = clf.score(scaled_X_test, y_test)
22
23
24 print(f"Mean Training Score: {dt_gs_training_score2 :.2%}")
25 print(f"Mean Test Score: {dt_gs_testing_score2 :.2%}")
26
27 # Model Performance
28 # Test set predictions
29
30 y_pred_test =clf.predict(scaled_X_test)
31 y_pred_train =clf.predict(scaled_X_train)
32
33 print(classification_report(y_test,y_pred_test))
34
35 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred_
36
37
38
39 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
40
41 plot_feature_importances(tree_clf)
```

Mean Training Score: 95.79%

Mean Test Score: 92.05%

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

False	0.94	0.96	0.95	570
True	0.76	0.66	0.71	97

▼ 21.2 Decrease

21.2.1 Void

21.2.2 Inter

	accuracy			
--	----------	--	--	--

macro avg	0.85	0.81	0.83	667
-----------	------	------	------	-----

weighted avg	0.92	0.92	0.92	667
--------------	------	------	------	-----

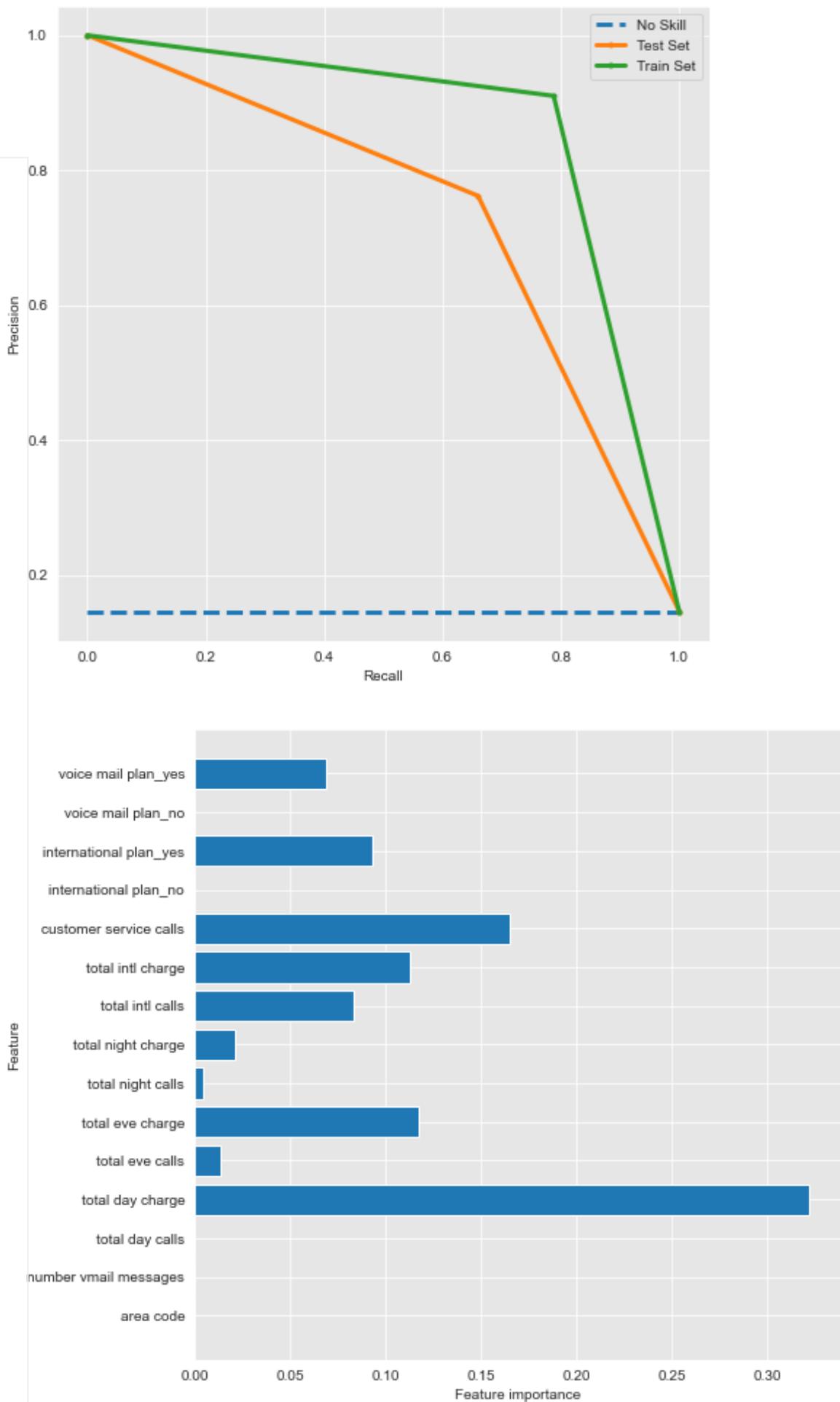
MODEL EVALUATION METRICS:

## Confusion Matrix for train &amp; test set:

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	2250.0	30.0	2280.0
2	Churn	82.0	304.0	386.0
3	All	2332.0	334.0	2666.0

**Contents** ⚙️

10.2 Results	Predicted	No Churn	Churn	All			
▼ 11 XGBoost	0 Actual	NaN	NaN	NaN			
11.1 Results	1 No Churn	550.0	20.0	570.0			
▼ 12 Gaussian N	2 Churn	33.0	64.0	97.0			
12.1 Results	3 All	583.0	84.0	667.0			
▼ 13 AdaBoostC	-----						
13.1 Results	Classification Report for train & test set						
▼ 14 Gradient Bo	Train set						
14.1 Results		precision	recall	f1-score			
▼ 15 SVM Classi				support			
15.1 Tuning	False	0.96	0.99	0.98			
15.1.1	True	0.91	0.79	0.84			
▼ 15.2 Results	accuracy			0.96			
15.2.1 SVI	macro avg	0.94	0.89	0.91			
15.2.2 SV	weighted avg	0.96	0.96	0.96			
15.3 Results				2666			
▼ 16 Stacking	Test set						
16.1 Results		precision	recall	f1-score			
▼ 17 SGDClassif				support			
17.1 Results	False	0.94	0.96	0.95			
▼ 18 Logistic Re	True	0.76	0.66	0.71			
18.1 Results	accuracy			0.92			
▼ 19 SGDClassif	macro avg	0.85	0.81	0.83			
19.1 Tuning	weighted avg	0.92	0.92	0.92			
19.2 Results				667			
▼ 20 Data Interp	-----						
20.1 Most In	accuracy			667			
20.1.1 Top	macro avg	0.85	0.81	0.83			
20.2 Top Po	weighted avg	0.92	0.92	0.92			
20.2.1 To				667			
20.2.2 Toj	-----						
▼ 21 Conclusion	Cohen's Kappa for train and test set:						
21.1 Increases	0.8203 0.6615						
21.1.1 Total	f2 score for train and test set:						
21.1.2 Cus	0.8094 0.678						
21.1.3 Total	roc auc score for train and test set:						
21.1.4 Total	0.8872 0.8124						
▼ 21.2 Decrease	Mean Cross Validation Score:						
21.2.1 Voic	0.9373						
21.2.2 Inte	-----						
22 Recomend	-----						
23 Future Wor	-----						



In [72]:

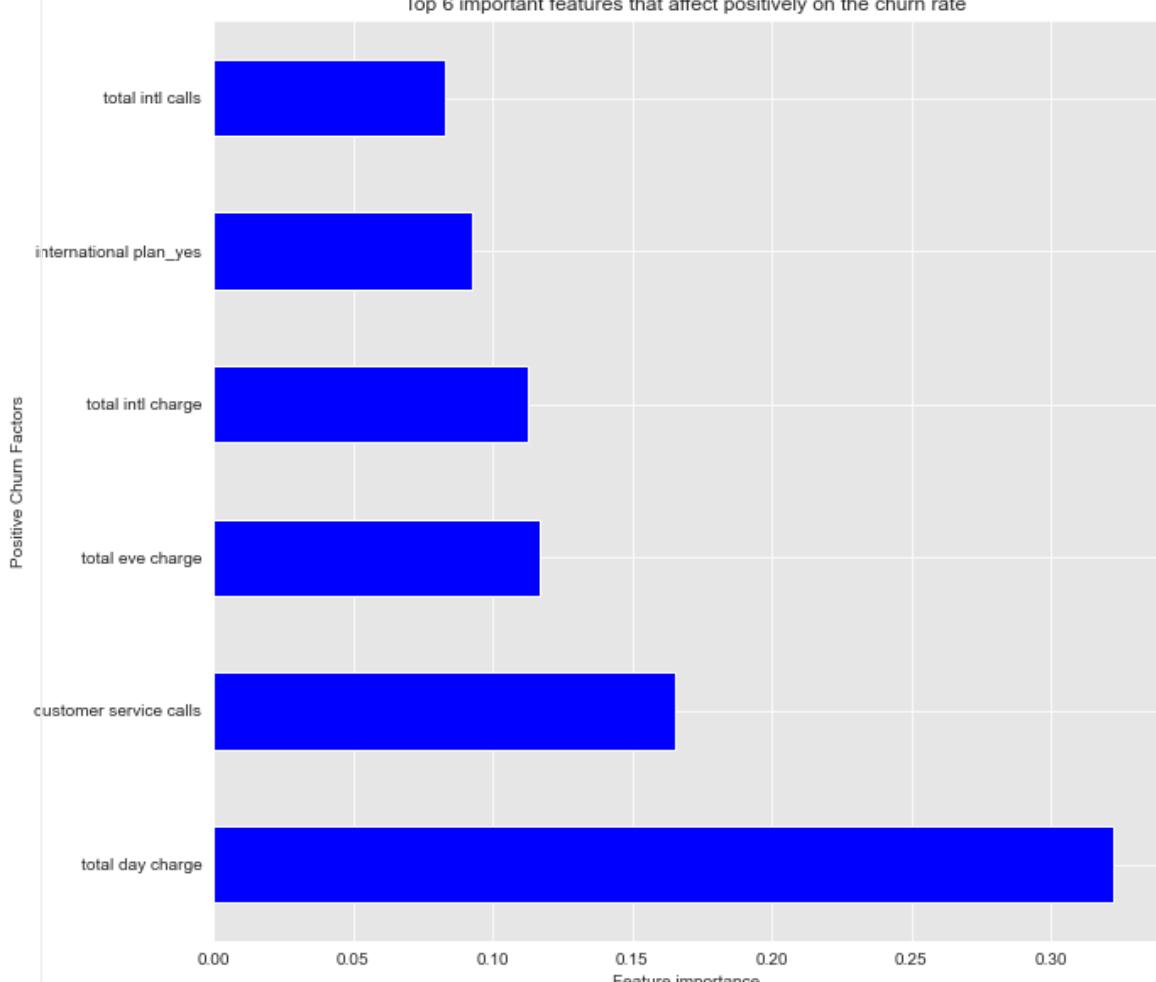
```

1 # Get Feature Importance from the classifier
2
3 feature_importance = tree_clf.feature_importances_
4 print (tree_clf.feature_importances_)
5 feat_importances = pd.Series(tree_clf.feature_importances_, index=X.columns)
6 feat_importances = feat_importances.nlargest(6)
7 feat_importances.plot(kind='barh', figsize=(10,10), color="b")
8
9 plt.xlabel('Feature importance')
10 plt.ylabel('Positive Churn Factors')
11 plt.title('Top 6 important features that affect positively on the churn rate')

```

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boos
  - 14.1 Results
- ▼ 15 SVM Classifi
  - 15.1 Tuning
  - ▼ 15.1.1 Hyp
  - 15.1.1.1
  - ▼ 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassifi
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most Infl
    - 20.1.1 Top 1
    - 20.1.2 Top 2
  - ▼ 20.2 Top Po
    - 20.2.1 Total
    - 20.2.2 Total
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases
    - 21.2.1 Voicemail
    - 21.2.2 Internet
- 22 Recomendation
- 23 Future Work



## 8.2 Results Comments

We have a precision of 0.87 and recall of 0.65 for the test set. We don't have overfitting. It means we are able to detect the churn customers 65% but 87% of the customers we say 'churn' was true.

The top 5 factors affecting te churn rate:

- 1) Total day charge
- 2) Customer service calls
- 3) Total evening charge
- 4) Total international charge
- 3) Total International plan\_yes
- 4) Total international calls

## Contents ↻ 9 Random Forest Classifier

10.2 Results <b>11 XGBBoost</b> 11.1 Results <b>12 Gaussian N</b> 12.1 Results <b>13 AdaBoostC</b> 13.1 Results <b>14 Gradient Bo</b> 14.1 Results <b>15 SVM Classi</b> ▼ 15.1 Tuning ▼ 15.1.1 Hyp 15.1.1.1 ▼ 15.2 Results 15.2.1 SVI 15.2.2 SV 15.3 Results <b>16 Stacking</b> In [77]: 16.1 Results <b>17 SGDClassif</b> 17.1 Results <b>18 Logistic Re</b> 18.1 Results <b>19 SGDClassif</b> 19.1 Tuning 19.2 Results <b>20 Data Interp</b> ▼ 20.1 Most In 20.1.1 Top ▼ 20.2 Top Po 20.2.1 To 20.2.2 To <b>21 Conclusion</b> ▼ 21.1 Increas 21.1.1 Total 21.1.2 Cus 21.1.3 Total 21.1.4 Total ▼ 21.2 Decrea 21.2.1 Voi 21.2.2 Inte <b>22 Recomend</b> <b>23 Future Wor</b>	<pre> 1 rf_clf = RandomForestClassifier() 2 rf_clf .fit(scaled_X_train, y_train) 3 4 # Model Performance 5 # Test set predictions 6 pred = rf_clf .predict(scaled_X_test) 7 8 9 10 rf_clf_score = cross_val_score(rf_clf, scaled_X_train, y_train, cv=10) 11 mean_rf_clf_score = np.mean(tree_clf_score) 12 13 print(f"Mean Cross Validation Score: {mean_rf_clf_score :.2%}") </pre> <p>Mean Cross Validation Score: 93.96%</p> <pre> 1 print('Training r^2:', tree_clf.score(scaled_X_train, y_train)) 2 print('Test r^2:', tree_clf.score(scaled_X_test, y_test)) 3 4 print(classification_report(y_test,pred)) </pre> <p>Training r^2: 0.9538634658664666    Test r^2: 0.9175412293853074</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th></th> <th style="text-align: center;">precision</th> <th style="text-align: center;">recall</th> <th style="text-align: center;">f1-score</th> <th style="text-align: center;">support</th> </tr> </thead> <tbody> <tr> <td>False</td> <td style="text-align: center;">0.94</td> <td style="text-align: center;">0.98</td> <td style="text-align: center;">0.96</td> <td style="text-align: center;">570</td> </tr> <tr> <td>True</td> <td style="text-align: center;">0.87</td> <td style="text-align: center;">0.60</td> <td style="text-align: center;">0.71</td> <td style="text-align: center;">97</td> </tr> <tr> <td>accuracy</td> <td></td> <td></td> <td style="text-align: center;">0.93</td> <td style="text-align: center;">667</td> </tr> <tr> <td>macro avg</td> <td style="text-align: center;">0.90</td> <td style="text-align: center;">0.79</td> <td style="text-align: center;">0.83</td> <td style="text-align: center;">667</td> </tr> <tr> <td>weighted avg</td> <td style="text-align: center;">0.92</td> <td style="text-align: center;">0.93</td> <td style="text-align: center;">0.92</td> <td style="text-align: center;">667</td> </tr> </tbody> </table>		precision	recall	f1-score	support	False	0.94	0.98	0.96	570	True	0.87	0.60	0.71	97	accuracy			0.93	667	macro avg	0.90	0.79	0.83	667	weighted avg	0.92	0.93	0.92	667
	precision	recall	f1-score	support																											
False	0.94	0.98	0.96	570																											
True	0.87	0.60	0.71	97																											
accuracy			0.93	667																											
macro avg	0.90	0.79	0.83	667																											
weighted avg	0.92	0.93	0.92	667																											

## 9.1 Hyperparameter Tuning for Random Forest using Gridsearch

```
In [78]: N
1 rf_param_grid = {
2     'n_estimators': [10, 30, 100],
3     'criterion': ['gini', 'entropy'],
4     'max_depth': [None, 2, 6, 10],
5     'min_samples_split': [5, 10],
6     'min_samples_leaf': [3, 6]
7 }
8
9 rf_grid_search = GridSearchCV(rf_clf, dt_param_grid, cv=10, return_trainin
10
11 # Fit to the data
12 rf_grid_search.fit(scaled_X_train, y_train)
13
14 #Examine best parameters
15
16 # Mean training score
17 rf_gs_training_score = np.mean(rf_grid_search.cv_results_['mean_trainin
18
19 # Mean test score
20 rf_gs_testing_score = rf_grid_search.score(scaled_X_test, y_test)
21
22 print(f"Mean Training Score: {rf_gs_training_score :.2%}")
23 print(f"Mean Test Score: {rf_gs_testing_score :.2%}")
24 print("Best Parameter Combination Found During Grid Search:")
25 rf_grid_search.best_params_
```

Mean Training Score: 91.03%  
 Mean Test Score: 93.85%  
 Best Parameter Combination Found During Grid Search:

```
Out[78]: {'criterion': 'entropy',
          'max_depth': None,
          'min_samples_leaf': 1,
          'min_samples_split': 5}
```

## 9.2 Best Paramerts for Random Forest

### Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - 15.1 Tuning
    - 15.1.1 Hyp
    - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassifi
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - 20.1 Most In
    - 20.1.1 Top
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

```
In [79]: N
1 dt_param_grid2 = {
2     'criterion': ['entropy'],
3     'max_depth': [None],
4     'min_samples_split': [5],
5     'min_samples_leaf': [1]
6 }
7
8 dt_grid_search2 = GridSearchCV(rf_clf, dt_param_grid2, cv=10, return_tr
9
10 # Fit to the data
11
12 clf=dt_grid_search2.fit(scaled_X_train, y_train)
13
14
15 #Examine best parameters
16
17 # Mean training score
18 dt_gs_training_score2 = np.mean(clf.cv_results_['mean_train_score'])
19
20 # Mean test score
21 dt_gs_testing_score2 = clf.score(scaled_X_test, y_test)
22
23
24 print(f"Mean Training Score: {dt_gs_training_score2 :.2%}")
25 print(f"Mean Test Score: {dt_gs_testing_score2 :.2%}")
26
27 # Model Performance
28 # Test set predictions
29
30 y_pred_test =clf.predict(scaled_X_test)
31 y_pred_train =clf.predict(scaled_X_train)
```

Mean Training Score: 98.82%

Mean Test Score: 93.25%

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - 15.1 Tuning
    - 15.1.1 Hyp
    - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - 20.1 Most In
    - 20.1.1 Top
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [80]:

```

1 print(classification_report(y_test,y_pred_test))
2
3 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred)
4
5
6
7 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
8
9 plot_feature_importances(rf_clf)

```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boos
  - 14.1 Results
- ▼ 15 SVM Classifi
  - 15.1 Tuning
    - 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassifi
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - 20.1 Most Infl
    - 20.1.1 Top
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreases
    - 21.2.1 Voic
    - 21.2.2 Inter
- 22 Recomendati
- 23 Future Wor

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

False	0.94	0.99	0.96	570
True	0.89	0.61	0.72	97
accuracy			0.93	667
macro avg	0.92	0.80	0.84	667
weighted avg	0.93	0.93	0.93	667

MODEL EVALUATION METRICS:

Confusion Matrix for train &amp; test set:

		Predicted	No Churn	Churn	All
		Actual	NaN	NaN	NaN
0	No Churn	2280.0	0.0	2280.0	
	Churn	31.0	355.0	386.0	
1	All	2311.0	355.0	2666.0	
	Actual	NaN	NaN	NaN	
2	No Churn	563.0	7.0	570.0	
	Churn	38.0	59.0	97.0	
3	All	601.0	66.0	667.0	

Classification Report for train &amp; test set

Train set

	precision	recall	f1-score	support
False	0.99	1.00	0.99	2280
	1.00	0.92	0.96	386
accuracy			0.99	2666
	0.99	0.96	0.98	2666
macro avg				
weighted avg	0.99	0.99	0.99	2666

Test set

	precision	recall	f1-score	support
False	0.94	0.99	0.96	570
	0.89	0.61	0.72	97

accuracy		0.93	667
macro avg	0.92	0.80	0.84
weighted avg	0.93	0.93	0.93

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

Cohen's Kappa for train and test set:

0.9514 0.6871

f2 score for train and test set:

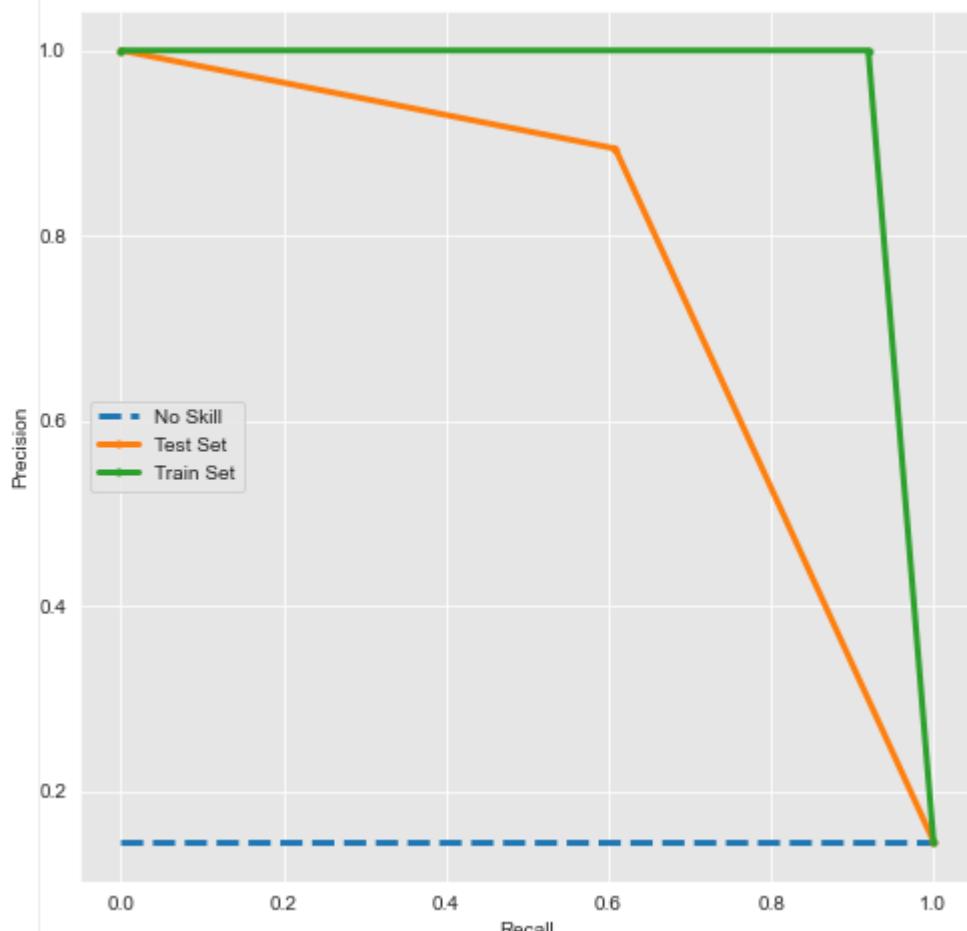
0.9347 0.6498

roc auc score for train and test set:

0.9598 0.798

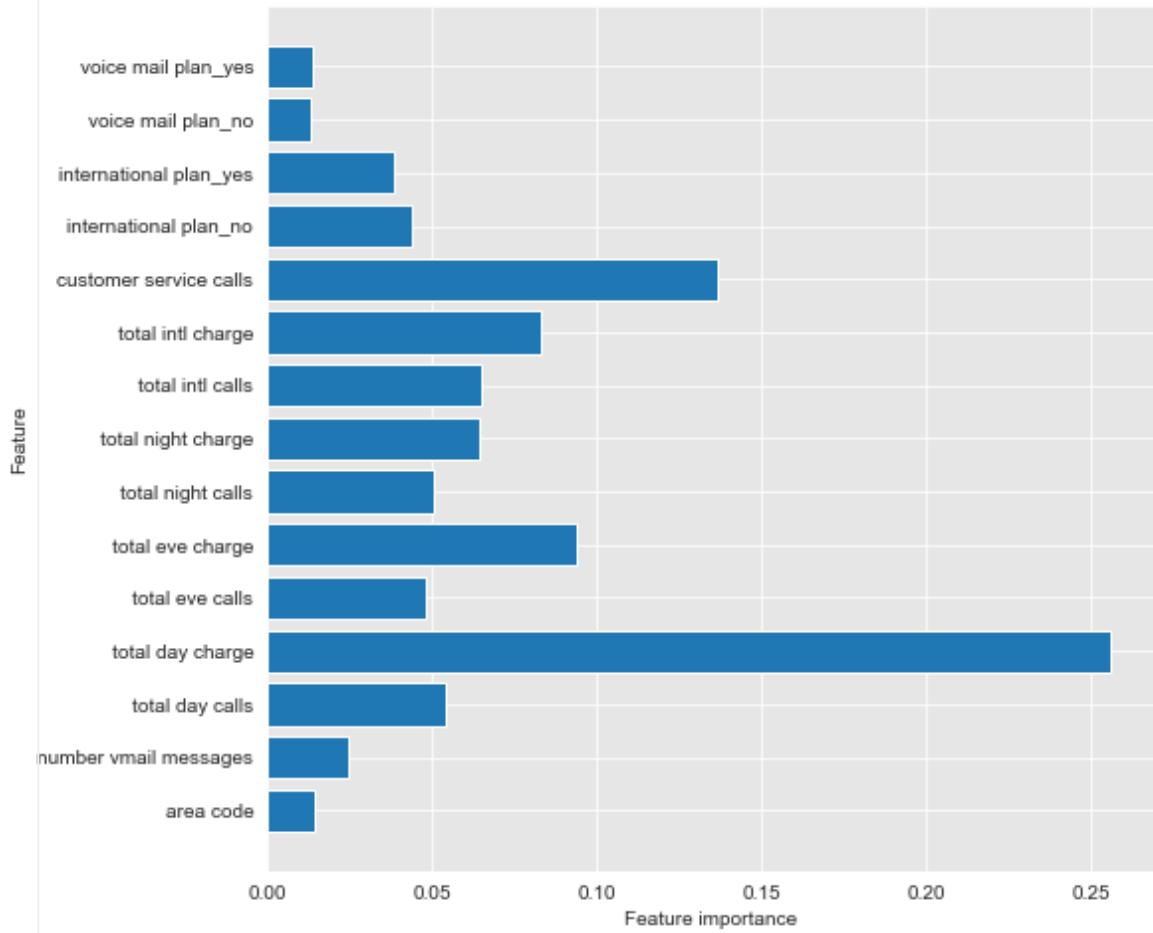
Mean Cross Validation Score:

0.9466



## Contents ⚙

10.2 Results
▼ 11 XGBoost
11.1 Results
▼ 12 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 Gradient Bo
14.1 Results
▼ 15 SVM Classi
▼ 15.1 Tuning
▼ 15.1.1 Hyp
15.1.1.1
▼ 15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results
▼ 18 Logistic Re
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
▼ 20.1 Most In
20.1.1 Top
▼ 20.2 Top Po
20.2.1 To
20.2.2 Toj
▼ 21 Conclusion
▼ 21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
▼ 21.2 Decreas
21.2.1 Voi
21.2.2 Inte
22 Recomend
23 Future Wor



In [81]:

```

1 feature_importance = rf_clf.feature_importances_
2 print (rf_clf.feature_importances_)
3 feat_importances = pd.Series(rf_clf.feature_importances_, index=X.columns)
4 feat_importances = feat_importances.nlargest(6)
5 feat_importances.plot(kind='barh' , figsize=(10,10), color="b")
6
7 plt.xlabel('Feature importance')
8 plt.ylabel('Positive Churn Factors')
9 plt.title('Top 6 important features that affect positively on the churn rate')

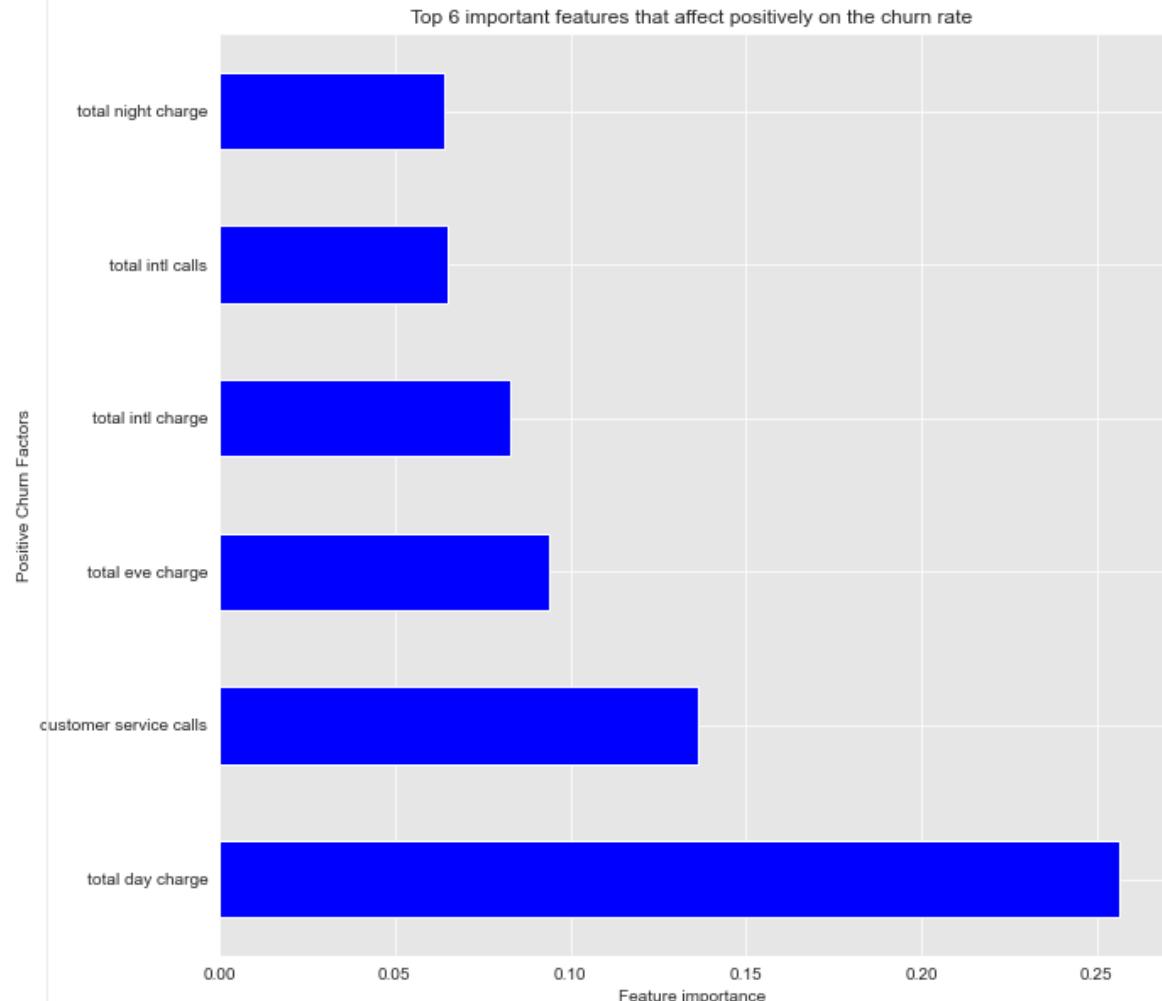
```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boosting
  - 14.1 Results
- ▼ 15 SVM Classification
  - ▼ 15.1 Tuning
    - 15.1.1 Hypothesis
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassification
  - 17.1 Results
- ▼ 18 Logistic Regression
  - 18.1 Results
- ▼ 19 SGDClassifier
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interpolation
  - ▼ 20.1 Most Influential Features
    - 20.1.1 Top 5
    - 20.1.2 Top 10
  - ▼ 20.2 Top Percentage
    - 20.2.1 Top 10%
    - 20.2.2 Top 20%
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total Churn
    - 21.1.2 Customer Service Calls
    - 21.1.3 Total Day Charge
    - 21.1.4 Total Night Charge
  - ▼ 21.2 Decreases
    - 21.2.1 Voicemail Requests
    - 21.2.2 Internet Usage
- 22 Recomendations
- 23 Future Work

[0.01450957 0.02484146 0.05424423 0.25617375 0.04835243 0.09392077  
 0.05054068 0.06424145 0.06487218 0.08277946 0.13643413 0.04379005  
 0.03817828 0.01335585 0.01376572]

**Out[81]:** Text(0.5, 1.0, 'Top 6 important features that affect positively on the churn rate')



## 9.3 Results & Comments

We have a precision of 0.61 and recall of 0.89 for the test set. We have some overfitting. It means we are able to detect the churn customers 61% but 89% of the customers we say 'churn' was true.

### Contents ⚙

- 10.1 The top 5 factors affecting te churn rate:
  - 10.2 Results
- 11 XGBoost
  - 11.1 Results
  - 1) Total day charge 2) Customer service calls 3) Total evening charge 4) Total international charge
  - 5) Total international calls 6) Total night charge
- 12 Gaussian N
  - 12.1 Results
- 13 AdaBoostC
  - 13.1 Results
- 14 Gradient Bc
  - 14.1 Results
- 15 SVM Classi
  - 15.1 Tuning
    - 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- 16 Stacking
  - 16.1 Results
- 17 SGDClassif
  - 17.1 Results
- 18 Logistic Re
  - 18.1 Results
- 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- 20 Data Interp
  - 20.1 Most In
    - 20.1.1 Top
    - 20.1.2 Toj
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- 21 Conclusion
  - 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreases
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [82]: N ▾ 1 #Instantiate KNeighborsClassifier

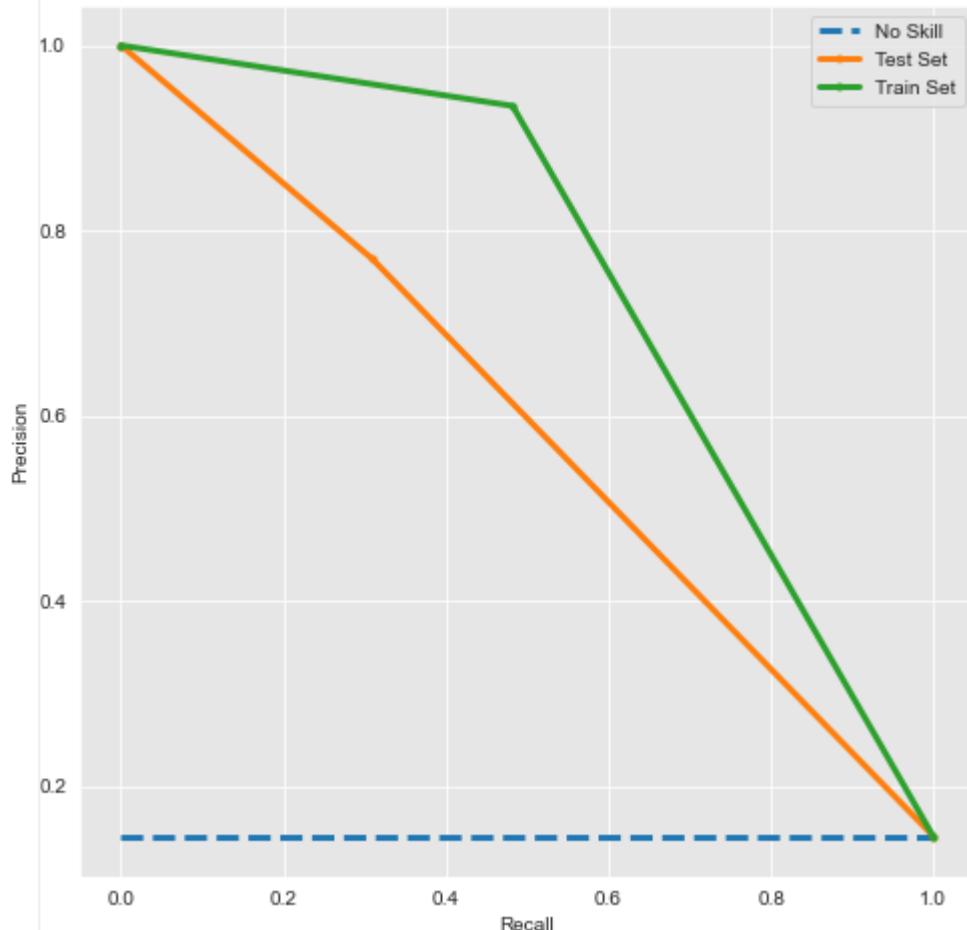
```

2 clf= KNeighborsClassifier()
3
4 # Fit the classifier
5 clf.fit(scaled_X_train, y_train)
6
7 # Predict on the test set
8 test_preds = clf.predict(scaled_X_test)
9 train_preds=clf.predict(scaled_X_train)
10
11
12
13
14
15
16 plot_pr_rc_curve(y_test, test_preds, y_train, train_preds)
17
18
19
20

```

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



In [83]: 1 model\_evaluation(scaled\_X\_train, scaled\_X\_test, y\_train, y\_test, train\_

#### MODEL EVALUATION METRICS:

##### Confusion Matrix for train & test set:

### Contents ⚙️

		Predicted	No Churn	Churn	All
	0	Actual	NaN	NaN	NaN
10.2	Results	0	2267.0	13.0	2280.0
11	XGBoost	1	No Churn	200.0	186.0
11.1	Results	2	Churn	186.0	386.0
12	Gaussian N	3	All	2467.0	199.0
12.1	Results				2666.0
13	AdaBoostC	Predicted	No Churn	Churn	All
13.1	Results	0	Actual	NaN	NaN
14	Gradient Bo	1	No Churn	561.0	9.0
14.1	Results	2	Churn	67.0	30.0
15	SVM Classi	3	All	628.0	39.0
15.1	Tuning				667.0
15.1.1	Hyp				
15.1.1.1					
15.2	Results				
15.2.1	SVI				
15.2.2	SV				
15.3	Results				
16	Stacking	False	0.92	0.99	0.96
16.1	Results	True	0.93	0.48	0.64
17	SGDClassif	accuracy			0.92
17.1	Results	macro avg	0.93	0.74	0.80
18	Logistic Re	weighted avg	0.92	0.92	0.91
18.1	Results				2666
19	SGDClassif				
19.1	Tuning				
19.2	Results				
20	Data Interp				
20.1	Most In	precision	recall	f1-score	support
20.1.1	Top	False	0.89	0.98	0.94
20.2	Top Po	True	0.77	0.31	0.44
20.2.1	To	accuracy			570
20.2.2	Toj	macro avg	0.83	0.65	0.69
21	Conclusion	weighted avg	0.88	0.89	0.86
21.1	Increas				667
21.1.1	Total				
21.1.2	Cus				
21.1.3	Total				
21.1.4	Total				
21.2	Decreas				
21.2.1	Voi				
21.2.2	Inte				
22	Recomend				
23	Future Wor				

Cohen's Kappa for train and test set:

0.5961 0.3903

f2 score for train and test set:

0.5336 0.3513

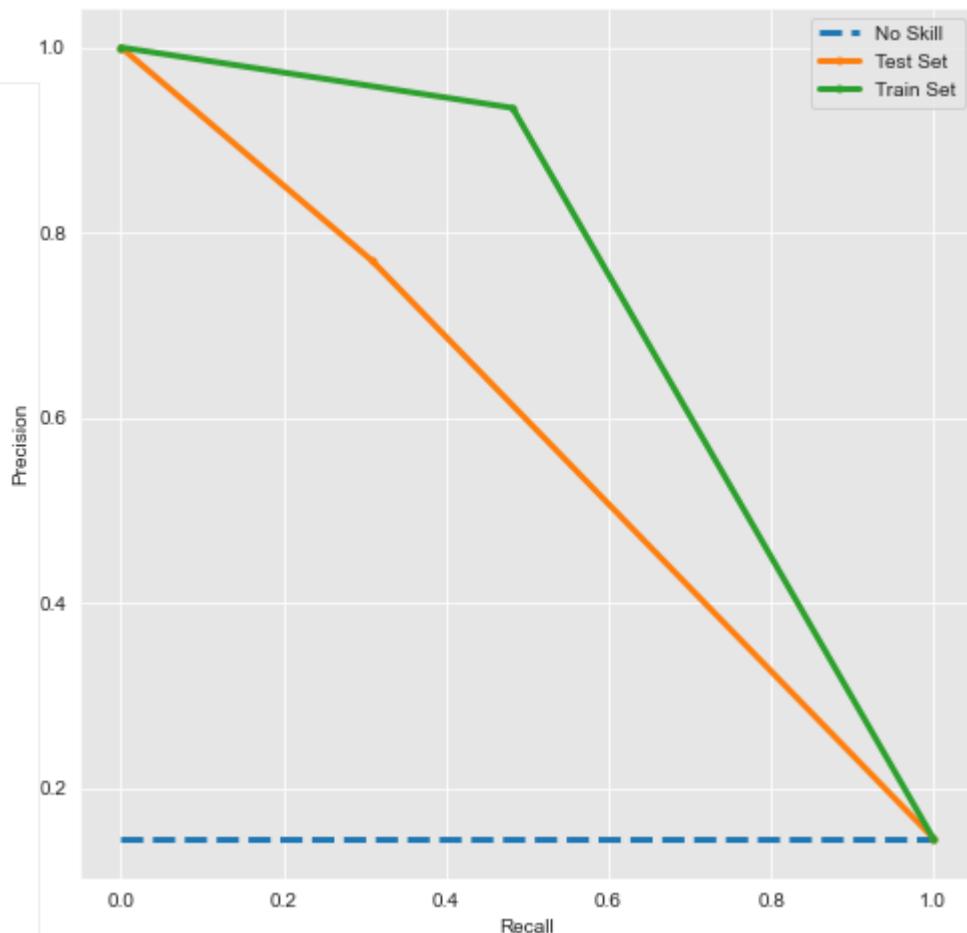
roc auc score for train and test set:

0.7381 0.6467

Mean Cross Validation Score:

0.8614

In [84]: 1 plot\_pr\_rc\_curve(y\_test, test\_preds, y\_train, train\_preds)



## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - 20.1 Most In
  - 20.1.1 Top
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj

## 10.1 Tuning & Optimization (Iterate over n values)

```
In [85]: 1 def find_best_k(scaled_X_train, y_train, scaled_X_test, y_test, min_k=1)
          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(scaled_X_train, y_train)
          7         preds = knn.predict(scaled_X_test)
          8         f1 = f1_score(y_test, preds)
          9         if f1 > best_score:
          10             best_k = k
          11             best_score = f1
          12
          13     print("Best Value for k: {}".format(best_k))
          14     print("F1-Score: {}".format(best_score))
```

```
In [86]: 1 find_best_k(scaled_X_train, y_train, scaled_X_test, y_test)
```

```
Best Value for k: 3  
F1-Score: 0.4583333333333334
```

## Contents ⚙

### 10.1.1 Select Best K for KNN Machine Learning Model

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
    - 20.1.2 Toj
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [87]:

```

1 clf = KNeighborsClassifier(n_neighbors=3)
2 clf.fit(X_train, y_train)
3 train_preds = clf.predict(scaled_X_train)
4 test_preds = clf.predict(scaled_X_test)
5 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, train_
6 plot_pr_rc_curve(y_test, test_preds, y_train, train_preds)
7

```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

**MODEL EVALUATION METRICS:****Confusion Matrix for train & test set:**

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	2280.0	0.0	2280.0
2	Churn	386.0	0.0	386.0
3	All	2666.0	0.0	2666.0

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	570.0	0.0	570.0
2	Churn	97.0	0.0	97.0
3	All	667.0	0.0	667.0

**Classification Report for train & test set****Train set**

		precision	recall	f1-score	support
▼ 19 SGDClassif	False	0.86	1.00	0.92	2280
	True	0.00	0.00	0.00	386
▼ 20 Data Interp	accuracy			0.86	2666
▼ 20.1 Most In	macro avg	0.43	0.50	0.46	2666
20.1.1 Top	weighted avg	0.73	0.86	0.79	2666

**Test set**

		precision	recall	f1-score	support
▼ 21 Conclusion	False	0.85	1.00	0.92	570
▼ 21.1 Increas	True	0.00	0.00	0.00	97
21.1.1 Total	accuracy			0.85	667
21.1.2 Cus	macro avg	0.43	0.50	0.46	667
21.1.3 Total	weighted avg	0.73	0.85	0.79	667

**Cohen's Kappa for train and test set:**

0.0 0.0

**f2 score for train and test set:**

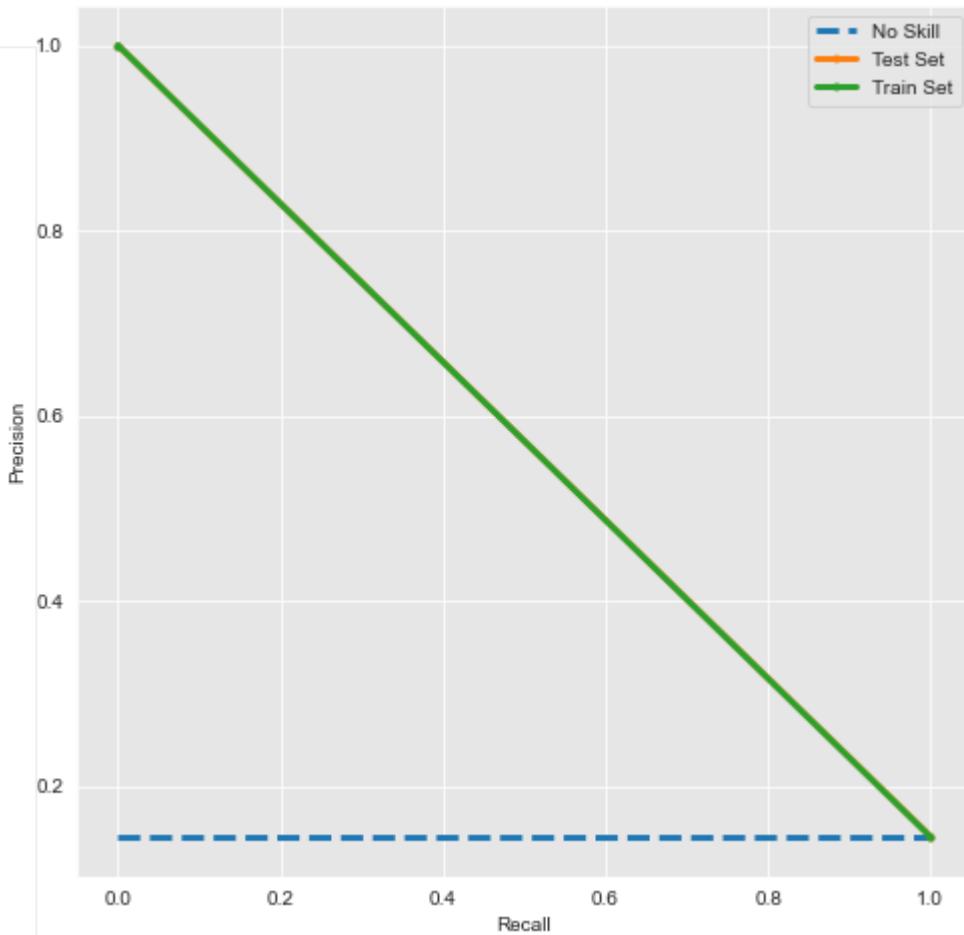
0.0 0.0

roc auc score for train and test set:

0.5 0.5

Mean Cross Validation Score:

0.8437



```

1
2 KNN_clf_score = cross_val_score(clf, scaled_X_train, y_train, cv=10)
3 mean_KNN_clf_score = np.mean(KNN_clf_score)
4
5
6 print(f"Mean Cross Validation Score: {mean_KNN_clf_score :.2%}")
7
8

```

Mean Cross Validation Score: 89.05%

## 10.1.2 Iterate over $1 < n < 26$

- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
- ▼ 21.2 Decreases
  - 21.2.1 Void
  - 21.2.2 Inter
- 22 Recomend
- 23 Future Wor

```
In [89]: 1 k_range = range(1, 26)
2
3 for k in k_range:
4     clf = KNeighborsClassifier(n_neighbors=k)
5     clf.fit(scaled_X_train, y_train)
6     y_pred_train = clf.predict(scaled_X_train)
7     y_pred_test = clf.predict(scaled_X_test)
8
9     print('-----')
10    print(k)
11    print('\nClassification Report for train & test set\n',
12          '\nTrain set\n',
13          classification_report(y_train, y_pred_train),
14          '\n\nTest set\n',
15          classification_report(y_test, y_pred_test))
16    print('-----\n')
17
18    print("f2 score for train and test set: \n",
19          round(fbeta_score(y_train, y_pred_train, 2.0), 4),
20          round(fbeta_score(y_test, y_pred_test, 2.0), 4))
```

## Contents ⚙️

10.2 Results
11 XGBoost
11.1 Results
12 Gaussian N
12.1 Results
13 AdaBoostC
13.1 Results
14 Gradient Bo
14.1 Results
15 SVM Classi
15.1 Tuning
15.1.1 Hyp
15.1.1.1
15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
16 Stacking
16.1 Results
17 SGDClassif
17.1 Results
18 Logistic Re
18.1 Results
19 SGDClassif
19.1 Tuning
19.2 Results
20 Data Interp
20.1 Most In
20.1.1 Top
20.2 Top Po
20.2.1 To
20.2.2 To
21 Conclusion
21.1 Increas
21.1.1 Tot
21.1.2 Cus
21.1.3 Tot
21.1.4 Tot
21.2 Decreas
21.2.1 Voic
21.2.2 Inte
22 Recomend
23 Future Wor

### 10.1.2.1 K=1

In [90]:

```

1 clf = KNeighborsClassifier(n_neighbors=1)
2 clf.fit(X_train, y_train)
3 train_preds = clf.predict(scaled_X_train)
4 test_preds = clf.predict(scaled_X_test)
5 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, train_
6 plot_pr_rc_curve(y_test, test_preds, y_train, train_preds)
7

```

**Contents**

10.2 Results
▼ 11 XGBoost
11.1 Results
▼ 12 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 Gradient Bo
14.1 Results
▼ 15 SVM Classi
▼ 15.1 Tuning
▼ 15.1.1 Hyp
15.1.1.1
▼ 15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results
▼ 18 Logistic Re
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
▼ 20.1 Most In
20.1.1 Top
▼ 20.2 Top Po
20.2.1 To
20.2.2 To
▼ 21 Conclusion
▼ 21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
▼ 21.2 Decreas
21.2.1 Voic
21.2.2 Inte
22 Recomend
23 Future Wor

**MODEL EVALUATION METRICS:****Confusion Matrix for train & test set:**

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	2280.0	0.0	2280.0
2	Churn	386.0	0.0	386.0
3	All	2666.0	0.0	2666.0

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	570.0	0.0	570.0
2	Churn	97.0	0.0	97.0
3	All	667.0	0.0	667.0

**Classification Report for train & test set****Train set**

	precision	recall	f1-score	support
False	0.86	1.00	0.92	2280
True	0.00	0.00	0.00	386
accuracy			0.86	2666
macro avg	0.43	0.50	0.46	2666
weighted avg	0.73	0.86	0.79	2666

**Test set**

	precision	recall	f1-score	support
False	0.85	1.00	0.92	570
True	0.00	0.00	0.00	97
accuracy			0.85	667
macro avg	0.43	0.50	0.46	667
weighted avg	0.73	0.85	0.79	667

**Cohen's Kappa for train and test set:**

0.0 0.0

**f2 score for train and test set:**

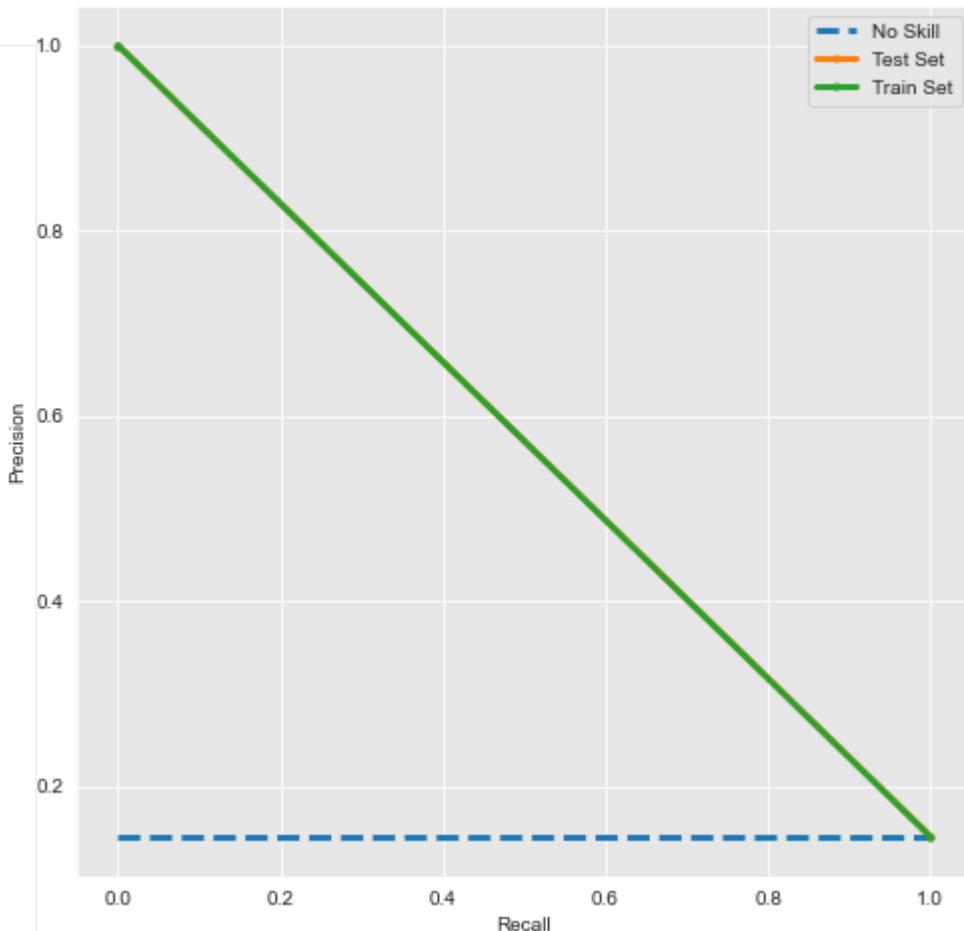
0.0 0.0

roc auc score for train and test set:

0.5 0.5

Mean Cross Validation Score:

0.7942



## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
  - ▼ 15.1.1 Hyp
  - 15.1.1.1
  - ▼ 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
  - 20.1.1 Top
  - ▼ 20.2 Top Po
  - 20.2.1 To
  - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
  - The mean cross validation is around 85%. This model performed poorly even for the optimal selected K (K=1 or 3) with poor recall snd precision for True Churn.
  - 21.1.1 Tot
  - 21.1.2 Cus
  - In [93]:
  - 21.1.3 Tot
  - 21.1.4 Tot
- ▼ 21.2 Decreas
- 21.2.1 Voi
- 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

## 10.2 Results & Comments

### 21 Conclusion

- ▼ 21.1 Increas
  - The mean cross validation is around 85%. This model performed poorly even for the optimal selected K (K=1 or 3) with poor recall snd precision for True Churn.

21.1.2 Cus

In [93]:

```
1 ! pip install xgboost
```

```
Requirement already satisfied: xgboost in c:\users\mirnamamaranda\anaconda3\lib\site-packages (1.3.3)
Requirement already satisfied: scipy in c:\users\mirnamamaranda\anaconda3\lib\site-packages (from xgboost) (1.4.1)
Requirement already satisfied: numpy in c:\users\mirnamamaranda\anaconda3\lib\site-packages (from xgboost) (1.19.2)
```

# 11 XGBoost

In [ ]:

```

1 from xgboost import XGBClassifier
2
3
4
5 clf = XGBClassifier(objective='binary:logistic',
6                      nthread=4,
7                      scale_pos_weight=3,
8                      seed=42,
9                      max_depth=3,
10                     n_estimators=140,
11                     learning_rate=0.005)
12
13 clf.fit(scaled_X_train, y_train)
14 y_pred_train = clf.predict(scaled_X_train)
15 y_pred_test = clf.predict(scaled_X_test)

```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

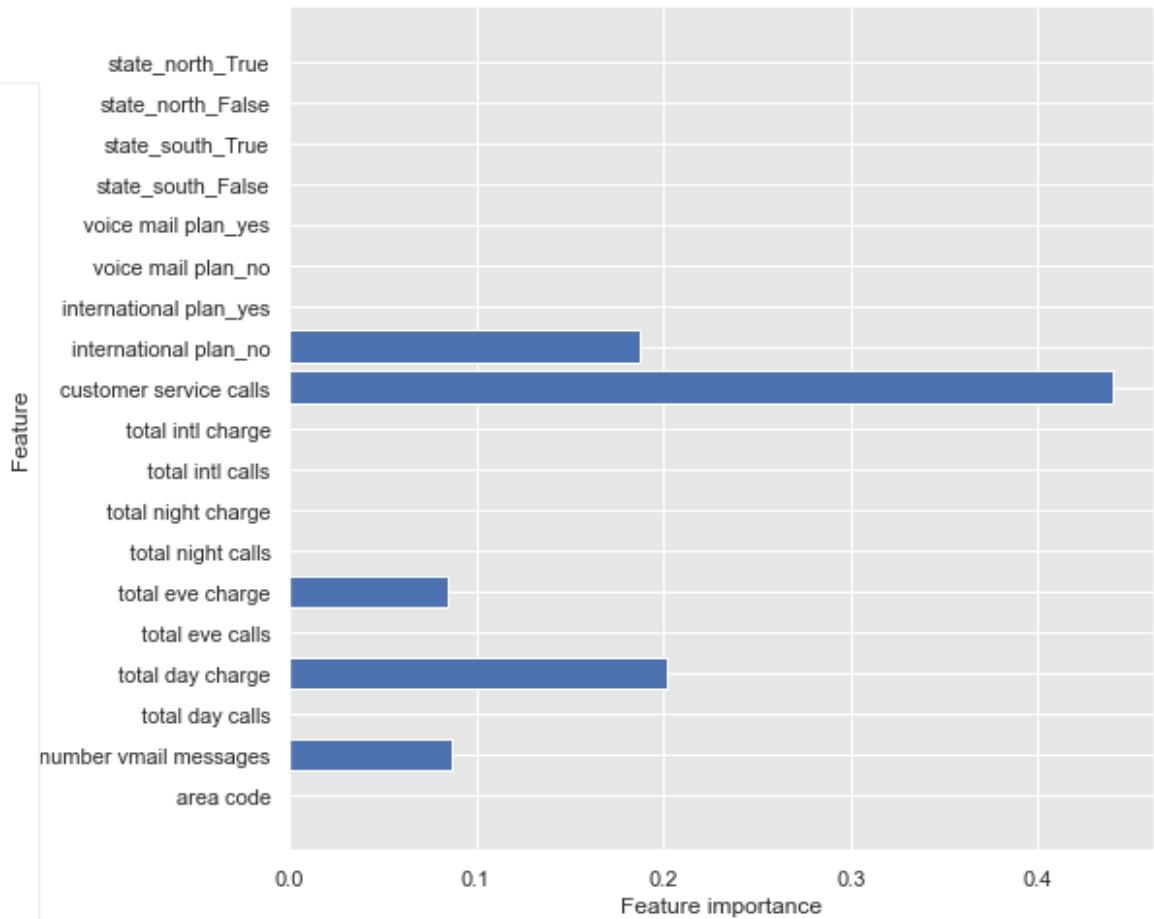
In [ ]:

```

1 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred_train, y_pred_test)
2 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
3
4

```

In [399]: 1 plot\_feature\_importances(clf)



## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voi
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [401]:

```

1 feature_importance = clf.feature_importances_
2 print (clf.feature_importances_)
3 feat_importances = pd.Series(clf.feature_importances_, index=X.columns)
4 feat_importances = feat_importances.nlargest(5)
5 feat_importances.plot(kind='barh' , figsize=(10,10), color="b")
6
7 plt.xlabel('Feature importance')
8 plt.ylabel('Positive Churn Factors')
9 plt.title('Top 5 important features that affect positively on the churn rate')

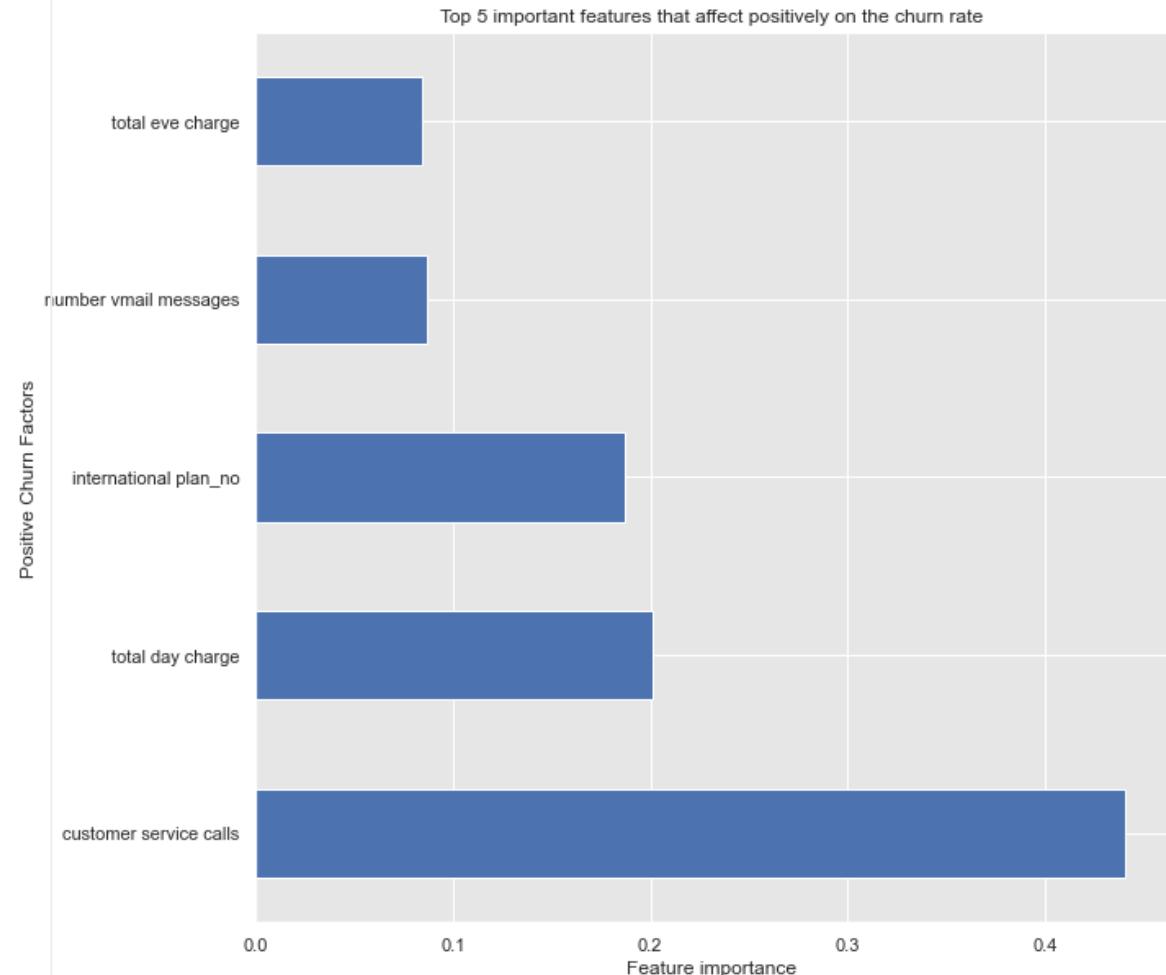
```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 GradientBoostingClassifier
  - 14.1 Results
- ▼ 15 SVM Classification
  - 15.1 Tuning
    - 15.1.1 Hypothesis
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassification
  - 17.1 Results
- ▼ 18 Logistic Regression
  - 18.1 Results
- ▼ 19 SGDClassification
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interpolation
  - 20.1 Most Influential Features
    - 20.1.1 Top 5
  - 20.2 Top Percentage
    - 20.2.1 Total
    - 20.2.2 Total
- ▼ 21 Conclusion
  - 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreases
    - 21.2.1 Voicemail
    - 21.2.2 International
- 22 Recomendations
- 23 Future Work

```
[0.        0.08669046 0.        0.20160858 0.        0.08463372
 0.        0.        0.        0.        0.4402261 0.18684117
 0.        0.        0.        0.        0.        0.]
```

`Out[401]:` Text(0.5, 1.0, 'Top 5 important features that affect positively on the churn rate')



## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boosting
  - 14.1 Results
- ▼ 15 SVM Classification
  - 15.1 Tuning
    - 15.1.1 Hypothesis: The mean cross validation is 88.69%. We have a precision of 0.62 and recall of 0.74 for the test set. We don't have overfitting. It means we are able to detect the 74% of churn customers but 62% of the customers we say 'churn' was true.
    - 15.2.1 SVI
    - 15.2.2 SV: The top 5 factors affecting the churn rate:
    - 15.3 Results
  - 16 Stacking
    - 1) Customer service calls
    - 16.1 Results
  - 17 SGDClassification
    - 2) Total day charge
    - 17.1 Results
  - 18 Logistic Regression
    - 3) International plan\_no
    - 18.1 Results
  - 19 SGDClassification
    - 4) Number voice messages
    - 19.1 Tuning
    - 19.2 Results
  - 20 Data Interpolation
    - 20.1 Most Influential Factors
      - 20.1.1 Top 5 factors
    - 20.2 Top Percentage
  - 21 Conclusion
    - 21.1 Increases in Churn
      - 21.1.1 Total duration
      - 21.1.2 Customer service calls
      - 21.1.3 Total number of calls
      - 21.1.4 Total number of voice messages
    - 21.2 Decreases in Churn
      - 21.2.1 Voice messages
      - 21.2.2 International plan
  - 22 Recomendations
  - 23 Future Work

## 12 Gaussian NB

```
In [95]: nb = GaussianNB()
          nb_model=nb.fit(scaled_X_train, y_train)
          y_pred_train = nb_model.predict(scaled_X_train)
          y_pred_test = nb_model.predict(scaled_X_test)
```

```
In [96]: 1 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_train)
          2
          3 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
```

## MODEL EVALUATION METRICS:

**Contents ↗**

## 10.2 Results

## ▼ 11 XGBoost

## 11.1 Results

## ▼ 12 Gaussian N

## 12.1 Results

## ▼ 13 AdaBoostC

## 13.1 Results

## ▼ 14 Gradient Bo

## 14.1 Results

## ▼ 15 SVM Classi

## ▼ 15.1 Tuning

## ▼ 15.1.1 Hyp

## 15.1.1.1

## ▼ 15.2 Results

## 15.2.1 SVI

## 15.2.2 SV

## 15.3 Results

## ▼ 16 Stacking

## 16.1 Results

## ▼ 17 SGDClassif

## 17.1 Results

## ▼ 18 Logistic Re

## 18.1 Results

## ▼ 19 SGDClassif

## 19.1 Tuning

## 19.2 Results

## ▼ 20 Data Interp

## ▼ 20.1 Most In

## 20.1.1 Top

## ▼ 20.2 Top Po

## 20.2.1 To

## 20.2.2 To

## ▼ 21 Conclusion

## ▼ 21.1 Increas

## 21.1.1 Total

## 21.1.2 Cus

## 21.1.3 Total

## 21.1.4 Total

## ▼ 21.2 Decreas

## 21.2.1 Voic

## 21.2.2 Inte

## 22 Recomend

## 23 Future Wor

## Confusion Matrix for train &amp; test set:

	Predicted	No Churn	Churn	All
0	Actual	Nan	Nan	Nan
1	No Churn	2280.0	0.0	2280.0
2	Churn	0.0	386.0	386.0
3	All	2280.0	386.0	2666.0

	Predicted	No Churn	Churn	All
0	Actual	Nan	Nan	Nan
1	No Churn	570.0	0.0	570.0
2	Churn	0.0	97.0	97.0
3	All	570.0	97.0	667.0

## Classification Report for train &amp; test set

## Train set

	precision	recall	f1-score	support
False	1.00	1.00	1.00	2280
True	1.00	1.00	1.00	386
accuracy			1.00	2666
macro avg	1.00	1.00	1.00	2666
weighted avg	1.00	1.00	1.00	2666

## Test set

	precision	recall	f1-score	support
False	1.00	1.00	1.00	570
True	1.00	1.00	1.00	97
accuracy			1.00	667
macro avg	1.00	1.00	1.00	667
weighted avg	1.00	1.00	1.00	667

## Cohen's Kappa for train and test set:

1.0 1.0

## f2 score for train and test set:

1.0 1.0

## roc auc score for train and test set:

1.0 1.0

[14:57:27] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was cha

nged from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:57:27] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:57:27] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:57:27] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:57:27] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

Mean Cross Validation Score:

0.8869

## Contents ⚙️

10.2 Results

▼ 11 XGBoost

11.1 Results

▼ 12 Gaussian N

12.1 Results

▼ 13 AdaBoostC

13.1 Results

▼ 14 Gradient Bo

14.1 Results

▼ 15 SVM Classi

▼ 15.1 Tuning

▼ 15.1.1 Hyp

15.1.1.1

▼ 15.2 Results

15.2.1 SVI

15.2.2 SV

15.3 Results

▼ 16 Stacking

16.1 Results

▼ 17 SGDClassif

17.1 Results

▼ 18 Logistic Re

18.1 Results

▼ 19 SGDClassif

19.1 Tuning

19.2 Results

▼ 20 Data Interp

▼ 20.1 Most In

20.1.1 Top

▼ 20.2 Top Po

20.2.1 To

20.2.2 Toj

▼ 21 Conclusion

▼ 21.1 Increas

21.1.1 Total

21.1.2 Cus

21.1.3 Total

21.1.4 Total

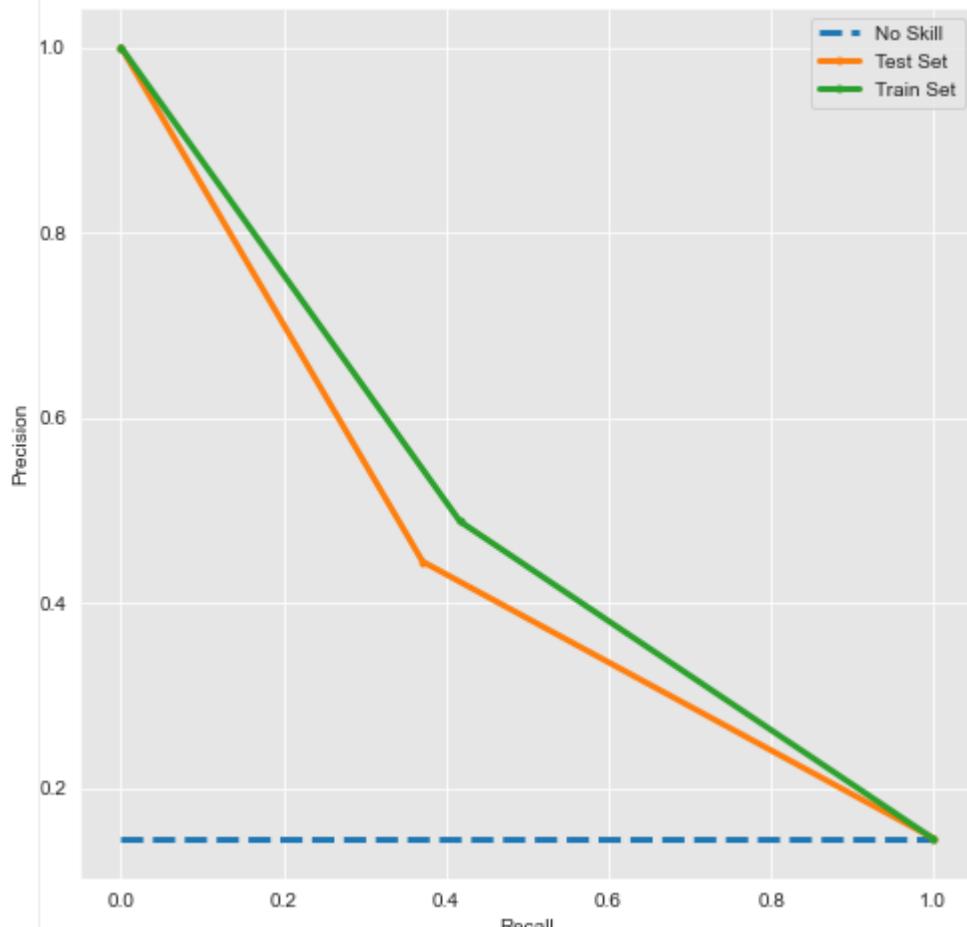
▼ 21.2 Decreas

21.2.1 Voic

21.2.2 Inte

22 Recomend

23 Future Wor



```
In [82]: nb_model.get_params
```

```
Out[82]: <bound method BaseEstimator.get_params of GaussianNB()>
```

## Contents ⚙

## 12.1 Results & Comments

- 10.2 Results
- 11 XGBoost
  - 11.1 Results The mean cross validation is 88.7%. We have a precision of 1 and recall of 1 for the test set. We
- 12 Gaussian NB
  - 12.1 Results don't have overfitting. It means we are able to detect the 100% of churn customers and 100% of the customers we say 'churn' was true.
- 13 AdaBoostC
  - 13.1 Results
- 14 Gradient BoC
  - 14.1 Results
- 15 SVM Classi
  - 15.1 Tuning
    - 15.1.1 Hyp
      - 15.1.1.1
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- 16 Stacking
  - 16.1 Results
- 17 SGDClassif
  - 17.1 Results
- 18 Logistic Re
  - 18.1 Results
- 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- 20 Data Interp
  - 20.1 Most In
    - 20.1.1 Top
    - 20.1.2 To
  - 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- 21 Conclusion
  - 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreases
    - 21.2.1 Void
    - 21.2.2 Inter
- 22 Recomend
- 23 Future Wor

In [97]:

```

1 from sklearn.ensemble import AdaBoostClassifier
2
3
4 ada = AdaBoostClassifier(n_estimators=120,random_state=5,learning_rate=0.01)
5 ada_model=ada.fit(scaled_X_train,y_train)
6
7 y_pred_train = ada_model.predict(scaled_X_train)
8 y_pred_test = ada_model.predict(scaled_X_test)
9
10 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_train,
11                      y_pred_train, y_pred_test)
12 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)

```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voi
    - 21.2.2 Inte
  - 22 Recomend
  - 23 Future Wor

MODEL EVALUATION METRICS:

Confusion Matrix for train &amp; test set:

		Predicted	No Churn	Churn	All
		Actual	NaN	NaN	NaN
0	No Churn	2280.0	0.0	2280.0	
	Churn	0.0	386.0	386.0	
3	All	2280.0	386.0	2666.0	

Classification Report for train &amp; test set

Train set

		precision	recall	f1-score	support
▼ 20.1 Most In	False	1.00	1.00	1.00	2280
	True	1.00	1.00	1.00	386
▼ 20.2 Top Po	accuracy			1.00	2666
	macro avg	1.00	1.00	1.00	2666
▼ 21 Conclusion	weighted avg	1.00	1.00	1.00	2666

Test set

		precision	recall	f1-score	support
▼ 21.2 Decreas	False	1.00	1.00	1.00	570
	True	1.00	1.00	1.00	97
22 Recomend	accuracy			1.00	667
	macro avg	1.00	1.00	1.00	667
23 Future Wor	weighted avg	1.00	1.00	1.00	667

Cohen's Kappa for train and test set:

1.0 1.0

f2 score for train and test set:

1.0 1.0

roc auc score for train and test set:

1.0 1.0

[14:58:28] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:58:28] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:58:28] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:58:29] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[14:58:29] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

Mean Cross Validation Score:

0.8869

## Contents ⚙️

10.2 Results

▼ 11 XGBoost

11.1 Results

▼ 12 Gaussian N

12.1 Results

▼ 13 AdaBoostC

13.1 Results

▼ 14 Gradient Bo

14.1 Results

▼ 15 SVM Classi

▼ 15.1 Tuning

15.1.1 Hyp

15.1.1.1

▼ 15.2 Results

15.2.1 SVI

15.2.2 SV

15.3 Results

▼ 16 Stacking

16.1 Results

▼ 17 SGDClassif

17.1 Results

▼ 18 Logistic Re

18.1 Results

▼ 19 SGDClassif

19.1 Tuning

19.2 Results

▼ 20 Data Interp

▼ 20.1 Most In

20.1.1 Top

▼ 20.2 Top Po

20.2.1 To

20.2.2 Toj

▼ 21 Conclusion

▼ 21.1 Increas

21.1.1 Total

21.1.2 Cus

21.1.3 Total

21.1.4 Total

▼ 21.2 Decrea

21.2.1 Voi

21.2.2 Inte

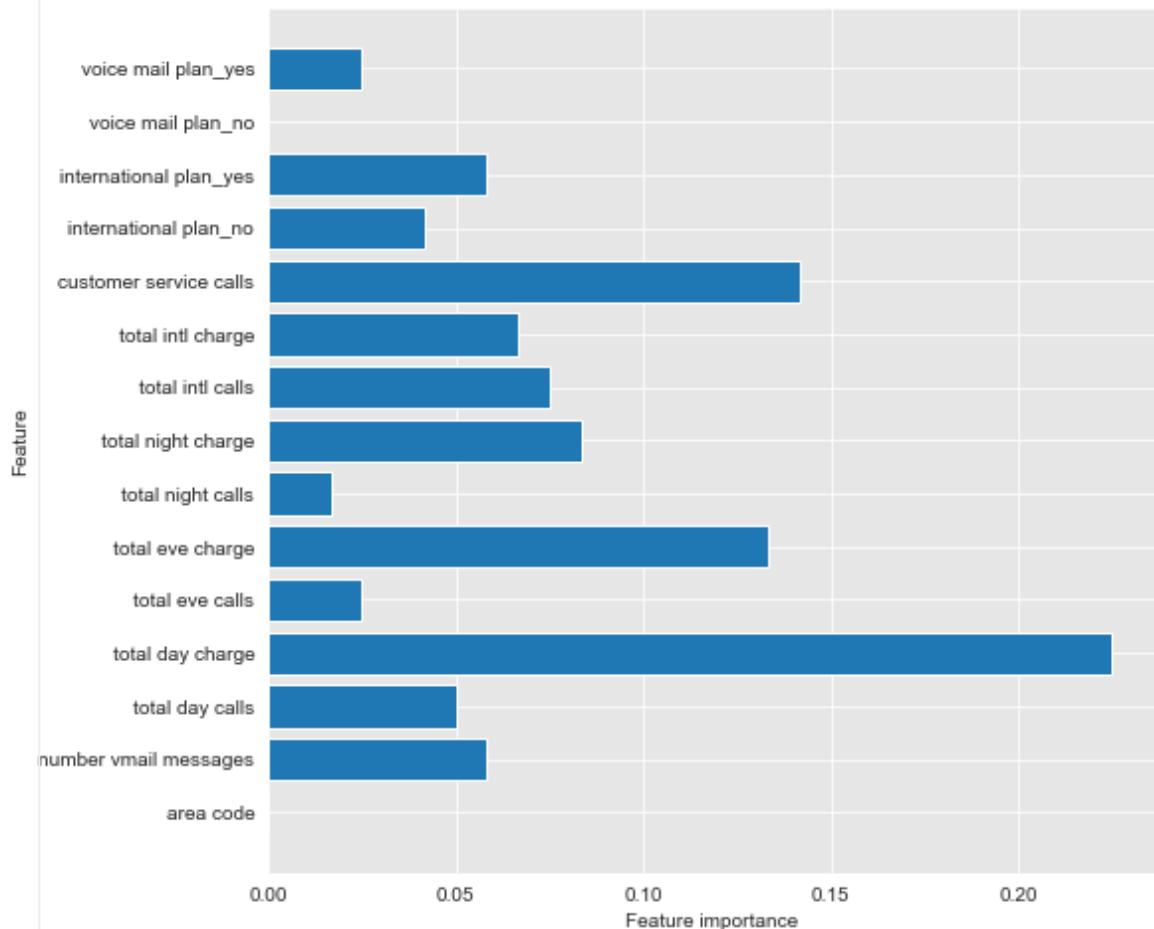
22 Recomend

23 Future Wor

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - In [98]:
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
  - ▼ 15.1.1 Hyp
  - 15.1.1.1
  - ▼ 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
  - 20.1.1 Top
  - ▼ 20.2 Top Po
  - 20.2.1 To
  - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
  - 21.1.1 Total
  - 21.1.2 Cus
  - 21.1.3 Total
  - 21.1.4 Total
  - ▼ 21.2 Decreas
  - 21.2.1 Voic
  - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

```
1 | plot_feature_importances(ada_model)
```



In [99]:

```

1 # Get Feature Importance from the classifier
2 feature_importance = ada_model.feature_importances_
3 print (ada_model.feature_importances_)
4 feat_importances = pd.Series(ada_model.feature_importances_, index=X.columns)
5 feat_importances = feat_importances.nlargest(13)
6 feat_importances.plot(kind='barh' , figsize=(10,10))

```

## Contents ⚙

10.2 Results  
 ▾ 11 XGBoost  
   11.1 Results  
 ▾ 12 GaussianNB  
   12.1 Results

13 AdaBoostC

  13.1 Results

14 Gradient Bo

  14.1 Results

15 SVM Classi

  15.1 Tuning

    15.1.1 Hyp

      15.1.1.1

    15.1.2 Results

      15.2.1 SVI

      15.2.2 SV

    15.3 Results

16 Stacking

  16.1 Results

17 SGDClassif

  17.1 Results

18 Logistic Re

  18.1 Results

19 SGDClassif

  19.1 Tuning

  19.2 Results

20 Data Interp

  20.1 Most In

    20.1.1 Top

    20.2 Top Po

      20.2.1 To

      20.2.2 Toj

21 Conclusion

  21.1 Increas

    21.1.1 Total

    21.1.2 Cus

    21.1.3 Total

    21.1.4 Total

  21.2 Decreas

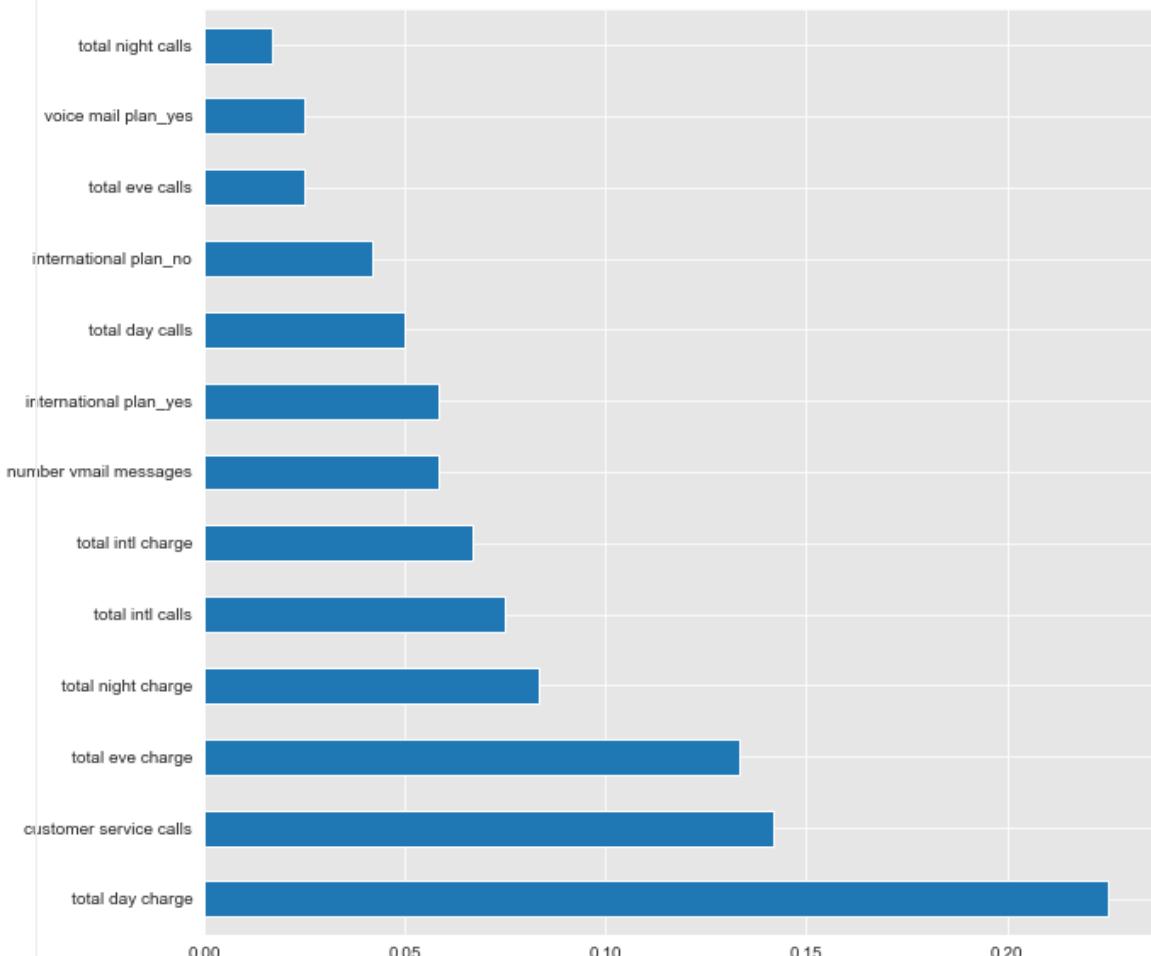
    21.2.1 Voi

    21.2.2 Inte

22 Recomend

23 Future Wor

Out[99]: <AxesSubplot:>



## 13.1 Results & Comments

The mean cross validation is 88.7%. We have a precision of 1 and recall of 1 for the test set. We don't have overfitting. It means we are able to detect the 100% of churn customers and 100% of non-churn customers we say 'churn' was true.

### Contents

- 10.2 Results The most important churn factors:
- ▼ 11 XGBoost
  - 11.1 Results 1) Total day charge
- ▼ 12 Gaussian N
  - 12.1 Results 2) Customer service calls
- ▼ 13 AdaBoostC
  - 13.1 Results 3) Total evening charge
- ▼ 14 Gradient Bc
  - 14.1 Results 4) Total International calls
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning 5) Total night charge
  - ▼ 15.1.1 Hyp
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [187]:

```

1 # Gradient Boosting
2 from sklearn.ensemble import GradientBoostingClassifier
3
4 gb = GradientBoostingClassifier(learning_rate=0.1,n_estimators=150,random_state=42)
5 gb_model=gb.fit(scaled_X_train,y_train)
6
7
8
9 y_pred_train = gb_model.predict(scaled_X_train)
10 y_pred_test = gb_model.predict(scaled_X_test)
11
12 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_train,
13 roc_curve_and_auc(gb_model, scaled_X_train, scaled_X_test, y_train, y_test))
14
15 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)

```

## Contents

10.2 Results				
▼ 11 XGBoost				
11.1 Results				
▼ 12 Gaussian N				
12.1 Results				
▼ 13 AdaBoostC				
13.1 Results				
▼ 14 Gradient Bo				
14.1 Results				
▼ 15 SVM Classif				
▼ 15.1 Tuning				
▼ 15.1.1 Hyp				
15.1.1.1				
▼ 15.2 Results				
15.2.1 SVI				
15.2.2 SV				
15.3 Results				
▼ 16 Stacking				
16.1 Results				
▼ 17 SGDClassif				
17.1 Results				
▼ 18 Logistic Re				
18.1 Results				
▼ 19 SGDClassif				
19.1 Tuning				
19.2 Results				
▼ 20 Data Interp				
▼ 20.1 Most In				
20.1.1 Top				
▼ 20.2 Top Po				
20.2.1 To				
20.2.2 Toj				
▼ 21 Conclusion				
▼ 21.1 Increas				
21.1.1 Total				
21.1.2 Cus				
21.1.3 Total				
21.1.4 Total				
▼ 21.2 Decreas				
21.2.1 Voic				
21.2.2 Inte				
22 Recomendat				
23 Future Wor				

weighted avg        1.00        1.00        1.00        667

---

Cohen's Kappa for train and test set:

1.0 1.0

f2 score for train and test set:

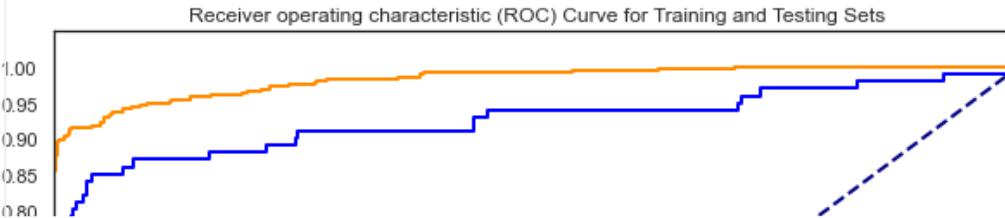
1.0 1.0

roc auc score for train and test set:

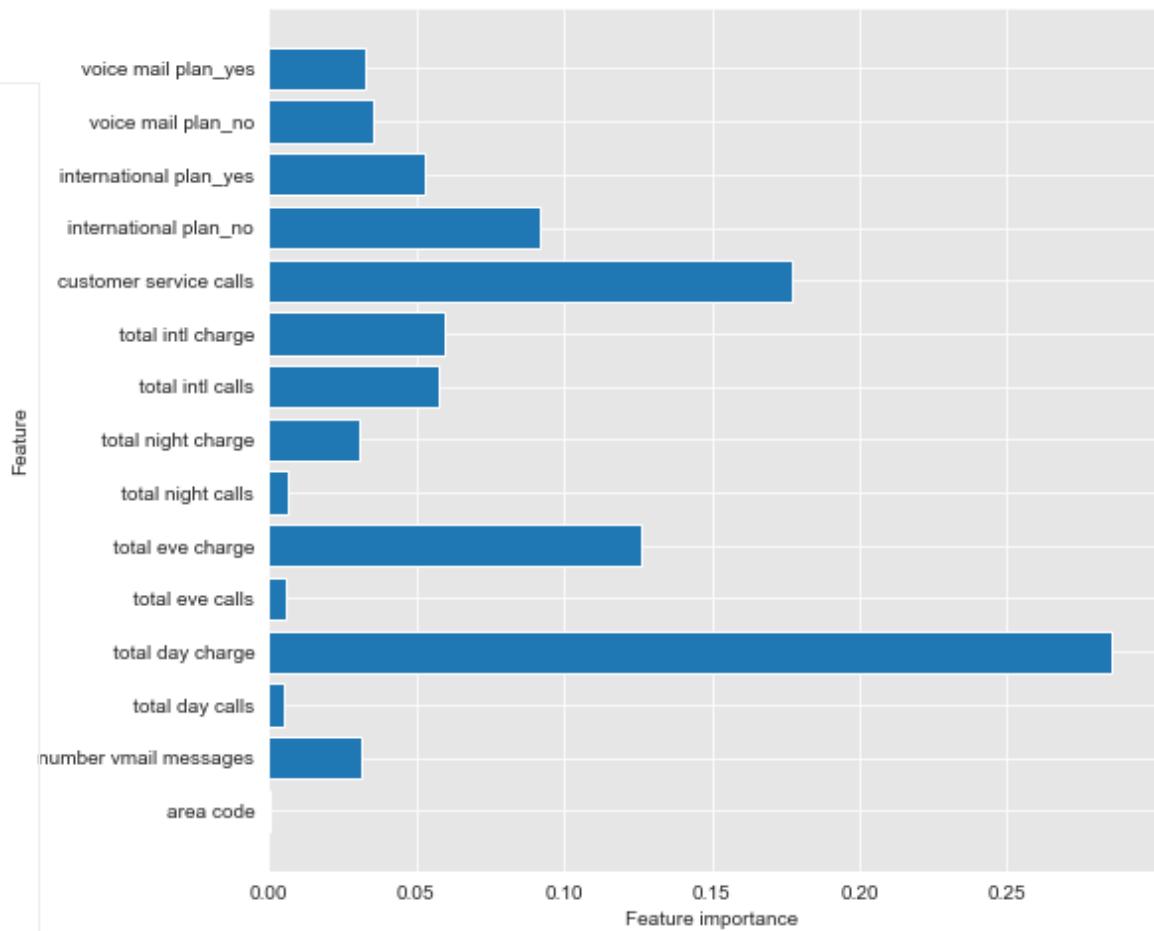
1.0 1.0

Mean Cross Validation Score:

0.7687



In [101]: 1 plot\_feature\_importances(gb\_model)



## Contents ⚙️

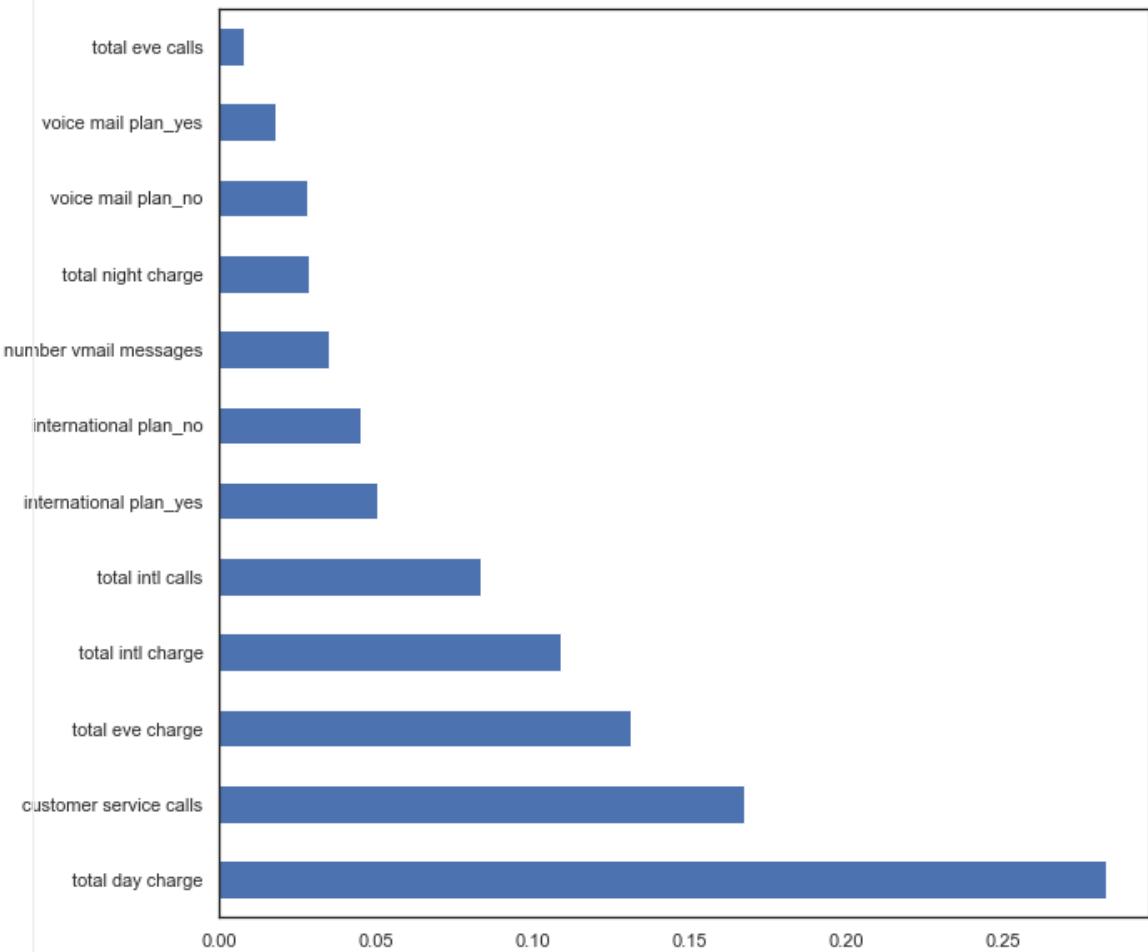
- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient BoC
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voi
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

```
In [188]: # Get Feature Importance from the classifier
1 feature_importance = gb_model.feature_importances_
2 print (gb_model.feature_importances_)
3 feat_importances = pd.Series(gb_model.feature_importances_, index=X.columns)
4 feat_importances = feat_importances.nlargest(12)
5 feat_importances.plot(kind='barh' , figsize=(10,10))
6
```

## Contents ⚙

10.2 Results  
 ▾ 11 XGBoost  
   11.1 Results  
 ▾ 12 Gaussian NB  
   12.1 Results  
 [0.00174112 0.03503597 0.00639164 0.28248003 0.00768167 0.13115795  
   0.00619365 0.02847671 0.08348926 0.108813 0.16733934 0.04515022  
   0.05019546 0.02789222 0.01796176]

Out[188]: <AxesSubplot:>



▀ 21 Conclusion

  21.1 Increases

    21.1.1 Total

    21.1.2 Customer

    21.1.3 Total

    21.1.4 Total

  21.2 Decreases

    21.2.1 Voice

    21.2.2 Internat

## 14.1 Results & comments

22 Recomend: The mean cross validation is 88.7%. We have a precision of 1 and recall of 1 for the test set. We  
 23 Future Work: don't have overfitting. It means we are able to detect the 100% of churn customers and 100% of  
                  the customers we say 'churn' was true.

The most important churn factors:

- 1) Total day charge
- 2) Customer service calls
- 3) Total evening charge

- 4) Total International charge

## Contents ⚙

- 10.2 Results

- 11 XGBoost

- 11.1 Results

- 12 Gaussian N

- 12.1 Results

- 13 AdaBoostC

- In [103]: 13.1 Results

- 14 Gradient Bo

- 14.1 Results

- 15 SVM Classi

- 15.1 Tuning

- 15.1.1 Hyp

- 15.1.1.1

- In [104]: 15.2 Results

- 15.2.1 SVI

- 15.2.2 SV

- 15.3 Results

- 16 Stacking

- 16.1 Results

- 17 SGDClassif

- In [105]: 17.1 Results

- 18 Logistic Re

- 18.1 Results

- 19 SGDClassif

- 19.1 Tuning

- 19.2 Results

- 20 Data Interp

- 20.1 Most In

- 20.1.1 Top

- 20.2 Top Po

- 20.2.1 To

- 20.2.2 Toj

- 21 Conclusion

- 21.1 Increas

- 21.1.1 Tot

- 21.1.2 Cus

- 21.1.3 Tot

- 21.1.4 Tot

- 21.2 Decreas

- 21.2.1 Voi

- 21.2.2 Inte

- 22 Recomend

- 23 Future Wor

## 15 SVM Classifier

```

1 #SVM Classifier
2 #Make a sample to run SVM faster
3
4 sample_X_train = scaled_X_train.sample(n=1000)
5 sample_y_train = y_train.sample(n=1000)

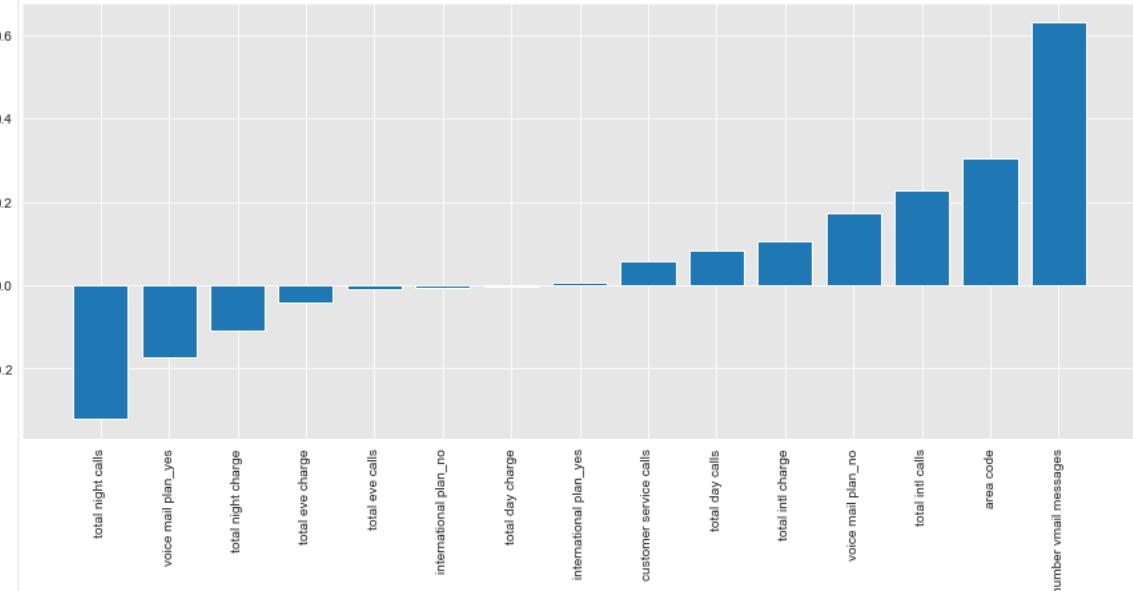
```

```

1 #SVM Classifier (Balanced Class Weight)
2 svc = SVC(kernel='linear', class_weight='balanced')
3 svc.fit(sample_X_train, sample_y_train)
4 y_pred_train = svc.predict(sample_X_train)
5 y_pred_test = svc.predict(scaled_X_test)

```

```
1 plot_coefficients(svc)
```



```
In [106]: 1 model_evaluation(sample_X_train, scaled_X_test, sample_y_train, y_test,
```

## MODEL EVALUATION METRICS:

## Confusion Matrix for train &amp; test set:

**Contents ↗ ⚙**

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	510.0	328.0	838.0
2	Churn	85.0	77.0	162.0
3	All	595.0	405.0	1000.0

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	347.0	223.0	570.0
2	Churn	65.0	32.0	97.0
3	All	412.0	255.0	667.0

## Classification Report for train &amp; test set

	Train set	precision	recall	f1-score	support
16	Stacking	False	0.86	0.61	0.71
16.1	Results	True	0.19	0.48	0.27
17	SGDClassif	accuracy		0.59	1000
17.1	Results	macro avg	0.52	0.54	0.49
18	Logistic Re	weighted avg	0.75	0.59	0.64
18.1	Results				1000

	Test set	precision	recall	f1-score	support
20.1	Most In	False	0.84	0.61	0.71
20.1.1	Top	True	0.13	0.33	0.18
20.2	Top Po	accuracy		0.57	667
20.2.1	To	macro avg	0.48	0.47	0.44
20.2.2	Toj	weighted avg	0.74	0.57	0.63
21	Conclusion				667

**21.1 Increases**

## 21.1.1 Total

## 21.1.2 Cus

## 21.1.3 Tot

## 21.1.4 Tot

**21.2 Decreases**

## 21.2.1 Voic

## 21.2.2 Inte

**22 Recomend****23 Future Wor**

## Cohen's Kappa for train and test set:

0.0523 -0.0366

## f2 score for train and test set:

0.3656 0.2488

## roc auc score for train and test set:

0.542 0.4693

[15:00:57] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[15:00:57] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[15:00:57] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[15:00:58] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

[15:00:58] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval\_metric if you'd like to restore the old behavior.

Mean Cross Validation Score:

0.8869

## Contents ⚙️

10.2 Results

▼ 11 XGBoost

11.1 Results

▼ 12 Gaussian N

12.1 Results

▼ 13 AdaBoostC

13.1 Results

▼ 14 Gradient Bo

14.1 Results

▼ 15 SVM Classi

▼ 15.1 Tuning

In [107]: Hyp

15.1.1

▼ 15.2 Results

15.2.1 SVI

15.2.2 SV

15.3 Results

▼ 16 Stacking

16.1 Results

▼ 17 SGDClassif

17.1 Results

▼ 18 Logistic Re

18.1 Results

▼ 19 SGDClassif

19.1 Tuning

19.2 Results

▼ 20 Data Interp

▼ 20.1 Most In

20.1.1 Top

▼ 20.2 Top Po

20.2.1 To

20.2.2 Toj

▼ 21 Conclusion

▼ 21.1 Increas

21.1.1 Total

21.1.2 Cus

21.1.3 Total

21.1.4 Total

▼ 21.2 Decreas

21.2.1 Voic

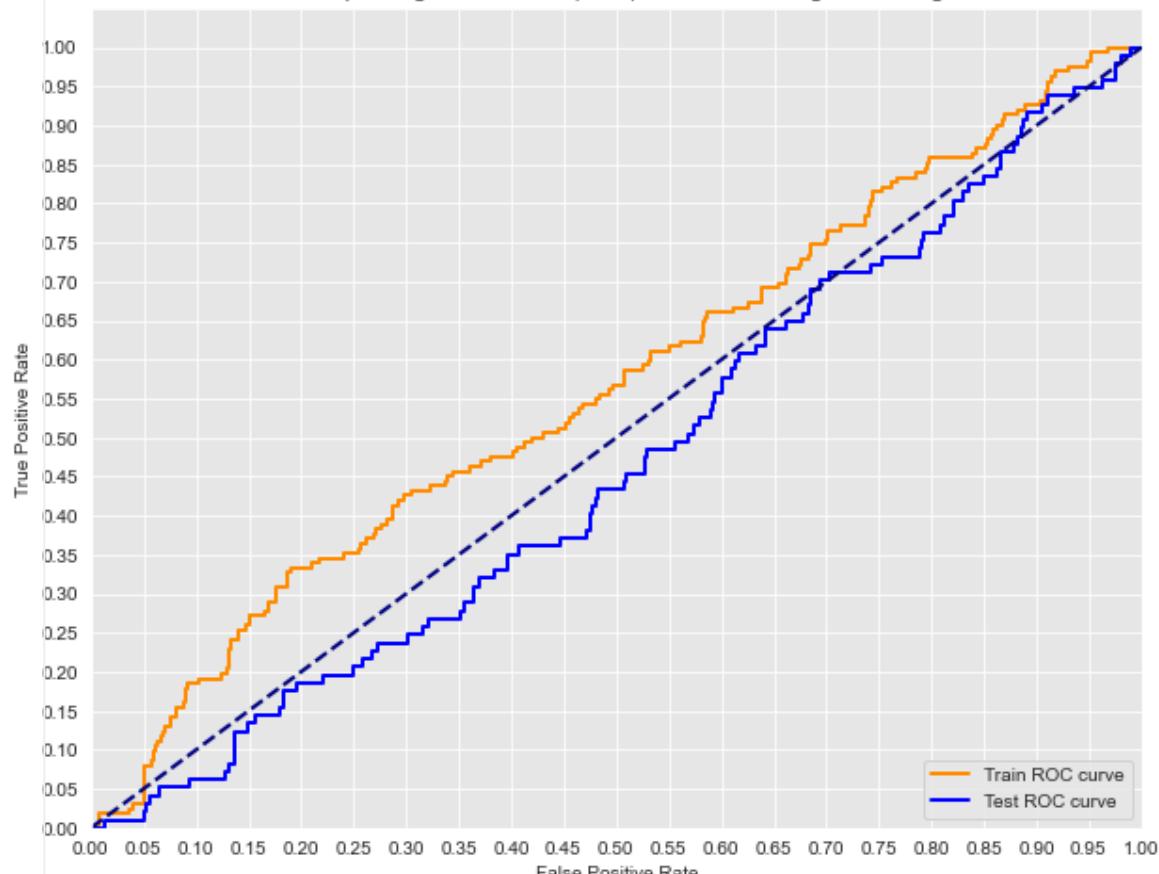
21.2.2 Inte

22 Recomend

23 Future Wor

```
1 roc_curve_and_auc(svc, sample_X_train, scaled_X_test, sample_y_train, )
```

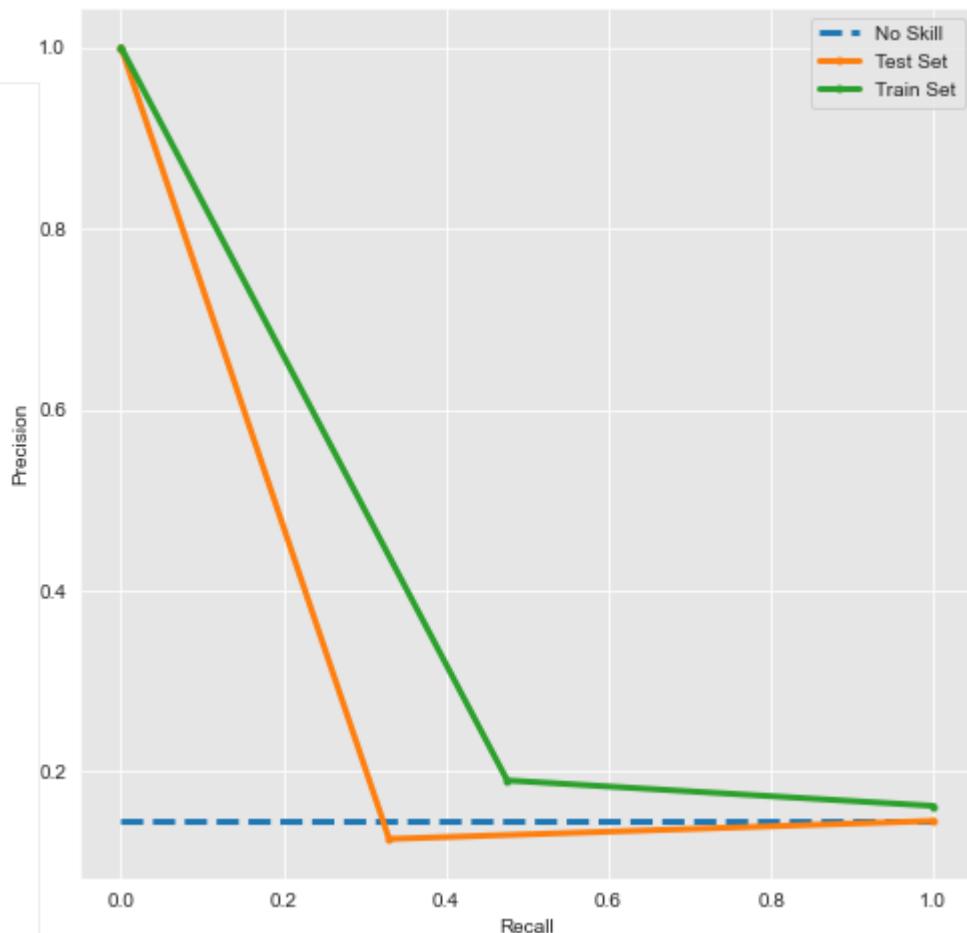
Receiver operating characteristic (ROC) Curve for Training and Testing Sets



Training AUC: 0.56507

Testing AUC: 0.46609

In [108]: 1 plot\_pr\_rc\_curve(y\_test, y\_pred\_test, sample\_y\_train, y\_pred\_train)



## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient BoC
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases
    - 21.2.1 Void
    - 21.2.2 Inter
- 22 Recomend
- 23 Future Wor

## 15.1 Tuning & Optimization

### 15.1.1 Hypertuning of SVM Classifier through GridSearchCV

In [109]: #SVM Classifier (balanced class weight and gridsearch)

```

1  param_grid = {'C': [0.1, 1, 10, 100, 1000],
2      'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
3      'kernel': ['rbf', "linear"]}
4
5
6
7 svc2 = GridSearchCV(SVC(class_weight = 'balanced'), param_grid, refit=True)
8
9
10 svc2.fit(scaled_X_train, y_train)
11 y_pred_train = svc2.predict(scaled_X_train)
12 y_pred_test = svc2.predict(X_test)
13
14 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred)
15
16 roc_curve_and_auc(svc2, scaled_X_train, scaled_X_test, y_train, y_test)
17
18 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
19
20 svc2.best_params_

```

CONFUSION MATRIX FOR train & test set:

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	2154.0	126.0	2280.0
2	Churn	37.0	349.0	386.0
3	All	2191.0	475.0	2666.0

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	570.0	0.0	570.0
2	Churn	97.0	0.0	97.0
3	All	667.0	0.0	667.0

Classification Report for train & test set

Train set

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

In [109]: svc2.best\_params\_

Out[109]:

```
{'C': 1, 'gamma': 0.1, 'kernel': 'rbf'}
```

21.1.4 Total

21.2 Decrease

21.2.1 Void 15.1.1.1 Best parameters for SVM Classifier

21.2.2 Inte

22 Recomend

23 Future Wor

```
In [115]: 1 param_grid2 = {'C': [1],
2           'gamma': [ 0.01],
3           'kernel': ['rbf']}
4
5 svc3 = GridSearchCV(SVC(class_weight = 'balanced'), param_grid2, refit=
```

7

```
8 svc3.fit(scaled_X_train, y_train)
9 y_pred_train = clf.predict(scaled_X_train)
10 y_pred_test = clf.predict(scaled_X_test)
11
12 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred)
13
14 roc_curve_and_auc(svc3, scaled_X_train, scaled_X_test, y_train, y_test)
15
16 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
```

## Contents ⚙️

10.2 Results	
▼ 11 XGBoost	11.1 Results
▼ 12 Gaussian N	12.1 Results
▼ 13 AdaBoostC	13.1 Results
▼ 14 Gradient Bc	14.1 Results
▼ 15 SVM Classi	True      0.13      0.11      0.12      382
▼ 15.1 Tuning	accuracy      0.76      2666
▼ 15.1.1 Hyp	macro avg      0.49      0.49      0.49      2666
▼ 15.2 Results	weighted avg      0.75      0.76      0.76      2666
15.2.1 SVI	
15.2.2 SV	
15.3 Results	Test set
▼ 16 Stacking	precision      recall      f1-score      support
16.1 Results	False      0.84      0.89      0.86      566
▼ 17 SGDClassif	True      0.10      0.07      0.08      101
17.1 Results	
▼ 18 Logistic Re	accuracy      0.76      667
18.1 Results	
▼ 19 SGDClassif	macro avg      0.47      0.48      0.47      667
19.1 Tuning	weighted avg      0.73      0.76      0.75      667
19.2 Results	
▼ 20 Data Interp	-----
▼ 20.1 Most In	Cohen's Kappa for train and test sets.
20.1.1 Top	
▼ 20.2 Top Po	
20.2.1 To	
20.2.2 To	
▼ 21 Conclusion	The mean cross validation is 95.02%. We have a precision of 0.10 and recall of 0.07 for the test set. This model performed poorly even after parameter optimization
▼ 21.1 Increas	21.1.1 Total
21.1.2 Cus	
21.1.3 Tot	
21.1.4 Tot	
▼ 21.2 Decrea	
21.2.1 Voic	
21.2.2 Inte	
22 Recomend	
23 Future Wor	

In [439]: 1 ! pip install delayed

```
Collecting delayed
  Downloading delayed-0.11.0b1-py2.py3-none-any.whl (19 kB)
Collecting hiredis
  Downloading hiredis-2.0.0-cp38-cp38-win_amd64.whl (18 kB)
Collecting redis
  Downloading redis-3.5.3-py2.py3-none-any.whl (72 kB)
Installing collected packages: hiredis, redis, delayed
Successfully installed delayed-0.11.0b1 hiredis-2.0.0 redis-3.5.3
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC **15.2.1 SVM with RFE**
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
    - 20.1.2 Top Po
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

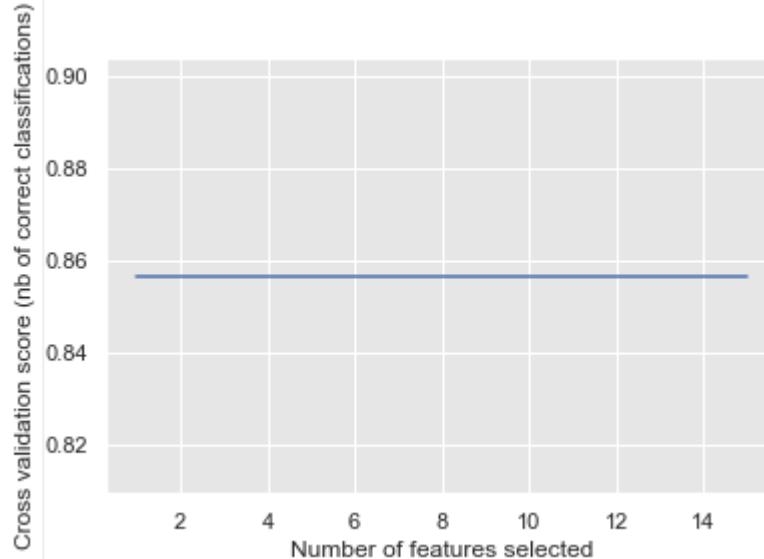
In [153]:

```

1 from sklearn.svm import SVC
2 from sklearn.model_selection import StratifiedKFold
3 from sklearn.feature_selection import RFECV
4 from sklearn.datasets import make_classification
5
6
7 # Build a classification task using 3 informative features
8
9
10 # Create the RFE object and compute a cross-validated score.
11 svc = SVC(kernel="linear")
12 # The "accuracy" scoring is proportional to the number of correct
13 # classifications
14
15 min_features_to_select = 1 # Minimum number of features to consider
16 rfecv = RFECV(estimator=svc, step=1, cv=StratifiedKFold(2),
17                 scoring='accuracy',
18                 min_features_to_select=min_features_to_select)
19 rfecv.fit(scaled_X_train, y_train)
20
21 print("Optimal number of features : %d" % rfecv.n_features_)
22
23 # Plot number of features VS. cross-validation scores
24 plt.figure()
25 plt.xlabel("Number of features selected")
26 plt.ylabel("Cross validation score (nb of correct classifications)")
27 plt.plot(range(min_features_to_select,
28                 len(rfecv.grid_scores_) + min_features_to_select),
29                 rfecv.grid_scores_)
30 plt.show()

```

Optimal number of features : 1



## Contents

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [154]: 1 rfecv.support\_

Out[154]: array([False, True, False, False])

## Contents

10.2 Results

▼ 11 XGBoost

11.1 Results

▼ 12 Gaussian N

12.1 Results

▼ 13 AdaBoostC

13.1 Results

▼ 14 Gradient Bo

14.1 Results

▼ 15 SVM Classi

15.1 Tuning

In [155]: Hyp

15.1.1

▼ 15.2 Results

15.2.1 SV

15.2.2 SV

15.3 Results

In [156]:

15.1.1

▼ 16 Stacking

16.1 Results

▼ 17 SGDClassif

17.1 Results

▼ 18 Logistic Re

18.1 Results

▼ 19 SGDClassif

19.1 Tuning

19.2 Results

▼ 20 Data Interp

20.1 Most In

20.1.1 Top

▼ 20.2 Top Po

20.2.1 To

20.2.2 To

▼ 21 Conclusion

21.1 Increases

21.1.1 Total

21.1.2 Cus

21.1.3 Total

21.1.4 Total

▼ 21.2 Decreas

21.2.1 Voic

21.2.2 Inte

22 Recomend

23 Future Wor

1 pd.DataFrame(data=rfecv.support\_.reshape(1, -1), columns=X.columns)

	area code	number vmail messages	total day calls	total day charge	total eve calls	total eve charge	total night calls	total night charge	total intl calls	ch
0	False	True	False	False	False	False	False	False	False	1

1 rfecv.ranking\_

Out[156]: array([10, 1, 8, 7, 9, 4, 13, 15, 11, 14, 5, 12, 6, 3, 2])

1 pd.DataFrame(data=rfecv.ranking\_.reshape(1, -1), columns=X.columns)

	area code	number vmail messages	total day calls	total day charge	total eve calls	total eve charge	total night calls	total night charge	total intl calls	ch
0	10	1	8	7	9	4	13	15	11	1

## 15.2.2 SVM with SelectKBest

localhost:8889/notebooks/Documents/TelecommuncationProject/Customer\_Churn\_Model.ipynb

In [158]:

```

1 from sklearn.feature_selection import f_regression, mutual_info_regression
2
3 from sklearn.svm import LinearSVC
4 from sklearn.pipeline import make_pipeline
5 from sklearn.feature_selection import SelectKBest, f_classif
6
7
8 plt.figure(1)
9 plt.clf()
10
11 X_indices = np.arange(X.shape[-1])
12
# #####
# Univariate feature selection with F-test for feature scoring
# We use the default selection function to select the four
# most significant features
13 selector = SelectKBest(f_classif, k=4)
14 selector.fit(scaled_X_train, y_train)
15 scores = -np.log10(selector.pvalues_)
16 scores /= scores.max()
17
18 plt.bar(X_indices - .45, scores, width=.2,
19         label=r'Univariate score ($-\text{Log}(p_{\text{value}})$)')
20
21
# #####
# Compare to the weights of an SVM
22 clf = make_pipeline(LinearSVC())
23 clf.fit(scaled_X_train, y_train)
24 print('Classification accuracy without selecting features: {:.3f}'
25       .format(clf.score(scaled_X_test, y_test)))
26
27 svm_weights = np.abs(clf[-1].coef_).sum(axis=0)
28 svm_weights /= svm_weights.sum()
29
30
31 plt.bar(X_indices - .25, svm_weights, width=.2, label='SVM weight')
32
33
34 clf_selected = make_pipeline(
35     SelectKBest(f_classif, k=4), LinearSVC())
36
37 clf_selected.fit(scaled_X_train, y_train)
38
39 print('Classification accuracy after univariate feature selection: {:.3f}'
40       .format(clf_selected.score(scaled_X_test, y_test)))
41
42
43 svm_weights_selected = np.abs(clf_selected[-1].coef_).sum(axis=0)
44 svm_weights_selected /= svm_weights_selected.sum()
45
46 plt.bar(X_indices[selector.get_support()] - .05, svm_weights_selected,
47         width=.2, label='SVM weights after selection')
48
49
50 plt.title("Comparing feature selection")
51 plt.xlabel('Feature number')
52 plt.yticks(())
53 plt.axis('tight')
54 plt.legend(loc='upper right')
55 plt.show()

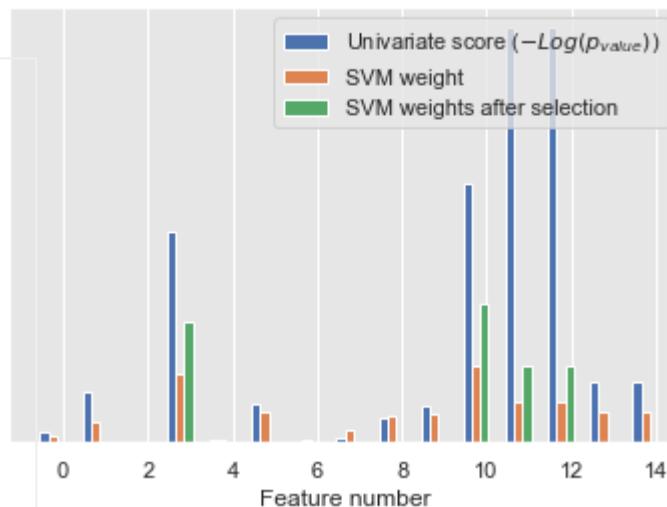
```

## Contents

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boos
  - 14.1 Results
- ▼ 15 SVM Classifi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hypothesis
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
    - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassifi
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassifi
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most Influential
  - 20.1.1 Top 10
  - ▼ 20.2 Top Positive
    - 20.2.1 Total
    - 20.2.2 Total
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases
    - 21.2.1 Void
    - 21.2.2 Intercept
- 22 Recomendations
- 23 Future Work

Classification accuracy without selecting features: 0.853  
 Classification accuracy after univariate feature selection: 0.844

Comparing feature selection



## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boosting
  - 14.1 Results
- ▼ 15 SVM Classification
  - ▼ 15.1 Tuning
    - 15.1.1 Hypothesis
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVC
    - 15.2.2 SVR
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassifiers
  - 17.1 Results
- ▼ 18 Logistic Regression
  - 18.1 Results
- ▼ 19 SGDClassifiers
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 In[150]: Out[150]
  - 20.1 Most Influential Features
    - 20.1.1 Top 10
    - 20.1.2 Top 20
  - 20.2 Top Positive Features
    - 20.2.1 Total
    - 20.2.2 Total
- ▼ 21 Conclusion
  - 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - 21.2 Decreases
    - 21.2.1 Void
    - 21.2.2 Interactions
- 22 Recomendations
- 23 Future Work

```
1 y_pred_test = clf_selected.predict(scaled_X_test)
2 print(classification_report(y_test, y_pred_test))
```

		precision	recall	f1-score	support
	False	0.86	0.98	0.91	566
	True	0.41	0.07	0.12	101
				0.84	667
	accuracy				
	macro avg	0.63	0.53	0.52	667
	weighted avg	0.79	0.84	0.79	667

## 15.3 Results & Comments

The SVC classifiers performed poorly even with RFE and SELECTK.

In [153]: 1 ! pip install sklearn.ensemble

```
ERROR: Could not find a version that satisfies the requirement sklearn.ensemble (from versions: none)
ERROR: No matching distribution found for sklearn.ensemble
```

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost 16 Stacking
- 11.1 Results
- ▼ 12 Gaussian NB
- 12.1 Results
- ▼ 13 AdaBoostC
- 13.1 Results
- ▼ 14 Gradient Boosting
- 14.1 Results
- ▼ 15 SVM Classification
- ▼ 15.1 Tuning
- ▼ 15.1.1 Hyperparameters
- 15.1.1.1 C
- 15.2 Results
- 15.2.1 SVI
- 15.2.2 SV
- 15.3 Results
- ▼ 16 Stacking
- 16.1 Results
- ▼ 17 SGDClassification
- 17.1 Results
- ▼ 18 Logistic Regression
- 18.1 Results
- ▼ 19 SGDClassif
- 19.1 Tuning
- 19.2 Results
- ▼ 20 Data Interpolation
- ▼ 20.1 Most Influential Features
- 20.1.1 Top 10
- 20.2 Top Predictors
- 20.2.1 Top 10
- 20.2.2 Top 20
- ▼ 21 Conclusion
- ▼ 21.1 Increases in Churn
- 21.1.1 Total Churn
- 21.1.2 Customer Age
- 21.1.3 Total Revenue
- 21.1.4 Total Expenses
- ▼ 21.2 Decreases in Churn
- 21.2.1 Void Recharge
- 21.2.2 Internet Service
- 22 Recomendations
- 23 Future Work

In [160]:

```

1 from sklearn.ensemble import StackingClassifier
2 from sklearn.ensemble import GradientBoostingClassifier
3 #Stacking
4
5 #Using Decisoin Tree, Random Forest, Gradient Boosting as base Learners
6
7 log_model = LogisticRegression()
8 dt_model = DecisionTreeClassifier(random_state=5,max_depth=10)
9 rf_model = RandomForestClassifier(n_estimators=100,random_state=5,max_c
10 gb_model = GradientBoostingClassifier(learning_rate=0.01,n_estimators=1
11
12
13 clf= StackingClassifier(estimators=[('tree', dt_model),('forest', rf_mc
14
15
16
17 clf.fit(scaled_X_train,y_train)
18
19 y_pred_train = clf.predict(scaled_X_train)
20 y_pred_test = clf.predict(scaled_X_test)
21
22 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred_
23
24 roc_curve_and_auc(clf, scaled_X_train, scaled_X_test, y_train, y_test)
25
26 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
27

```

## Contents

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
    - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - accuracy
  - macro avg
  - weighted avg
- ▼ 20.1 Most In
  - 20.1.1 Top
- ▼ 20.2 Top Po
  - 20.2.1 To
  - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

False	0.98	1.00	0.99	2284
True	1.00	0.86	0.92	382
accuracy			0.98	2666
macro avg	0.99	0.93	0.96	2666
weighted avg	0.98	0.98	0.98	2666

Test set

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

False	0.96	0.99	0.97	566
True	0.93	0.74	0.82	101
accuracy			0.95	667
macro avg	0.94	0.87	0.90	667
weighted avg	0.95	0.95	0.95	667

In [50]: 1 clf.get\_params

```
Out[50]: <bound method _BaseHeterogeneousEnsemble.get_params of StackingClassifier(e
stimators=[('tree',
              DecisionTreeClassifier(max_depth=10,
                                      random_state=5)),
             ('forest',
              RandomForestClassifier(max_depth=10,
                                      random_state=5)),
             ('gb',
              GradientBoostingClassifier(learning_rate=0.
                                         n_estimators=12
                                         random_state=
                                         5))],
                                         final_estimator=LogisticRegression())>
```

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - 15.1 Tuning
  - 15.1.1 Hyp
  - 15.1.1.1
  - 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - 20.1 Most In
  - 20.1.1 Top
  - 20.2 Top Po
  - 20.2.1 To
  - 20.2.2 Toj
- ▼ 21 Conclusion
  - 21.1 Increas
  - 21.1.1 Total
  - 21.1.2 Cus
  - 21.1.3 Total
  - 21.1.4 Total
  - 21.2 Decrea
  - 21.2.1 Voic
  - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

## 16.1 Results & Comments

The mean cross validation is 95%. We have a precision of 0.93 and recall of 0.74 for the test set. We don't have overfitting. It means we are able to detect the 74% of churn customers and 93% of the customers we say 'churn' was true.

## 17 SGDClassifier (Lasso & Ridge regularizations)

```
In [161]: 1 #SCDClassifier (balanced class weight)
2
3 scd = SGDClassifier(penalty='elasticnet', class_weight='balanced')
4 scd.fit(scaled_X_train, y_train)
5 y_pred_train = scd.predict(scaled_X_train)
6 y_pred_test = scd.predict(scaled_X_test)
7
8
9 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred)
10 roc_curve_and_auc(scd, scaled_X_train, scaled_X_test, y_train, y_test)
11
12 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
13
14
```

## Contents ☰⚙️

- 10.2 Results
  - ▼ 11 XGBoost
    - 11.1 Results
  - ▼ 12 Gaussian N
    - 12.1 Results
  - ▼ 13 AdaBoostC
    - 13.1 Results
  - ▼ 14 Gradient BoC
    - 14.1 Results
  - ▼ 15 SVM Classif
    - 15.1 Tuning
      - 15.1.1 Hyp
      - 15.1.1.1
    - 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
  - ▼ 16 Stacking
    - 16.1 Results
  - ▼ 17 SGDClassif
    - 17.1 Results
  - ▼ 18 Logistic Re
    - 18.1 Results
  - ▼ 19 SGDClassif
    - 19.1 Tuning
    - 19.2 Results
  - ▼ 20 Data Interp
    - ▼ 20.1 Most In
      - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 Toj
  - ▼ 21 Conclusion
    - ▼ 21.1 Increas
      - 21.1.1 Total
      - 21.1.2 Cus
      - 21.1.3 Total
      - 21.1.4 Total
    - ▼ 21.2 Decreas
      - 21.2.1 Voic
      - 21.2.2 Inte
  - 22 Recomend
  - 23 Future Wor

## MODEL EVALUATION METRICS:

### Confusion Matrix for train & test set:

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	1356.0	928.0	2284.0
2	Churn	77.0	305.0	382.0
3	All	1433.0	1233.0	2666.0

		Predicted	No Churn	Churn	All
▼ 16	Stacking	0	Actual	NaN	NaN
16.1	Results	1	No Churn	354.0	212.0
▼ 17	SGDClassif	2	Churn	19.0	82.0
17.1	Results	3	All	373.0	294.0
▼ 18	Logistic Re				667.0

## Classification Report for train & test set

## Train set

	precision	recall	f1-score	support
False	0.95	0.59	0.73	2284
True	0.25	0.80	0.38	382
accuracy			0.62	2666
macro avg	0.60	0.70	0.55	2666
weighted avg	0.65	0.63	0.58	2666

## Test set

	precision	recall	f1-score	support
False	0.95	0.63	0.75	566
True	0.28	0.81	0.42	101
accuracy			0.65	667
macro avg	0.61	0.72	0.58	667
weighted avg	0.85	0.65	0.70	667

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
      - 15.3 Results
  - ▼ 16 Stacking
    - 16.1 Results
  - ▼ 17 SGDClassif
    - 17.1 Results
  - ▼ 18 Logistic Re
    - 18.1 Results
  - ▼ 19 SGDClassif
    - 19.1 Tuning
    - 19.2 Results
  - ▼ 20 Data Interp
    - ▼ 20.1 Most In
      - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 Toj
  - ▼ 21 Conclusion
    - ▼ 21.1 Increas
      - 21.1.1 Total
      - 21.1.2 Cus
      - 21.1.3 Total
      - 21.1.4 Total
    - ▼ 21.2 Decreas
      - 21.2.1 Voic
      - 21.2.2 Inte
  - 22 Recomend
  - 23 Future Wor

Cohen's Kappa for train and test set:

0.2034 0.245

f2 score for train and test set:

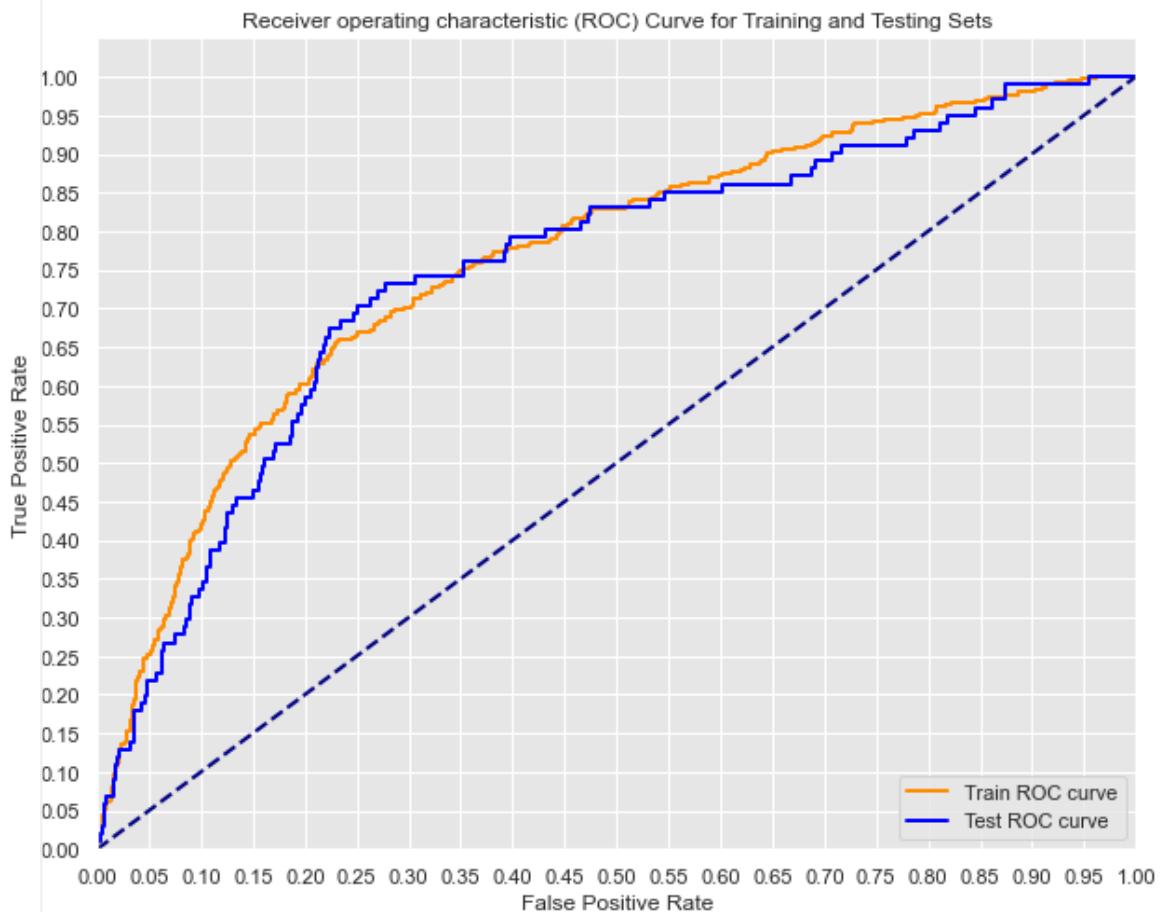
0.5523 0.5874

roc auc score for train and test set:

0.6961 0.7187

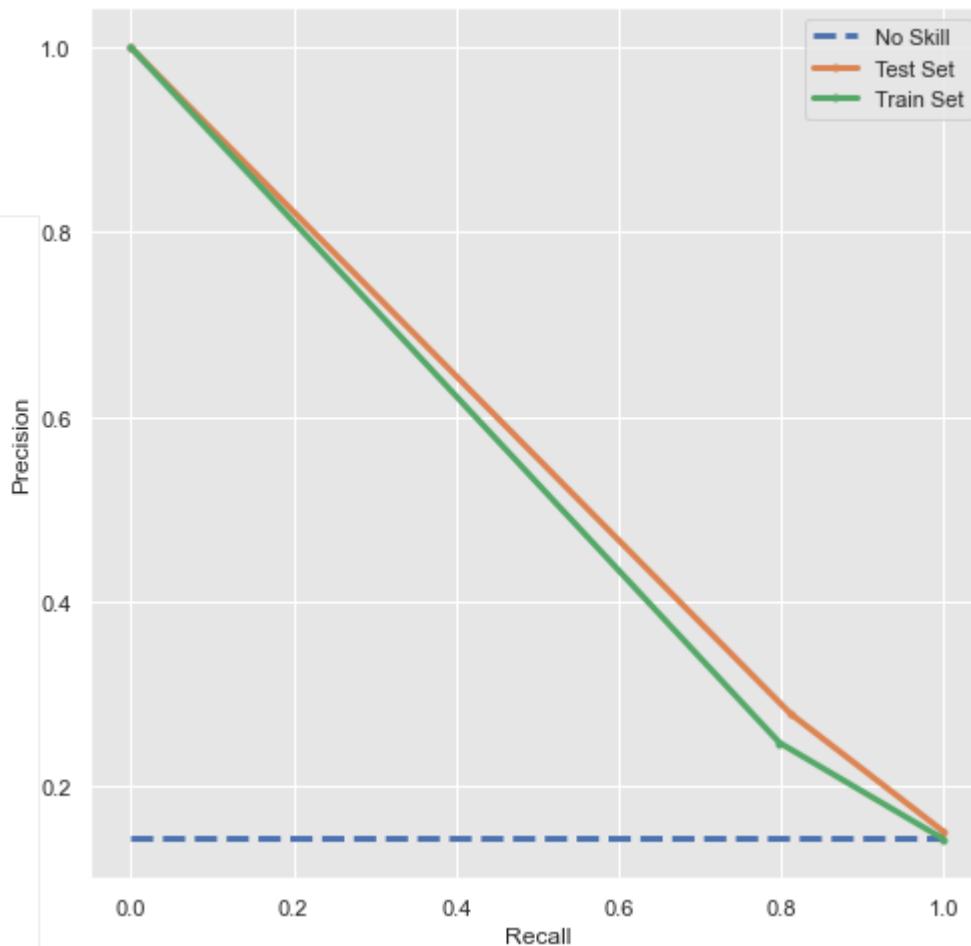
Mean Cross Validation Score:

0.9502



Training AUC: 0.76274

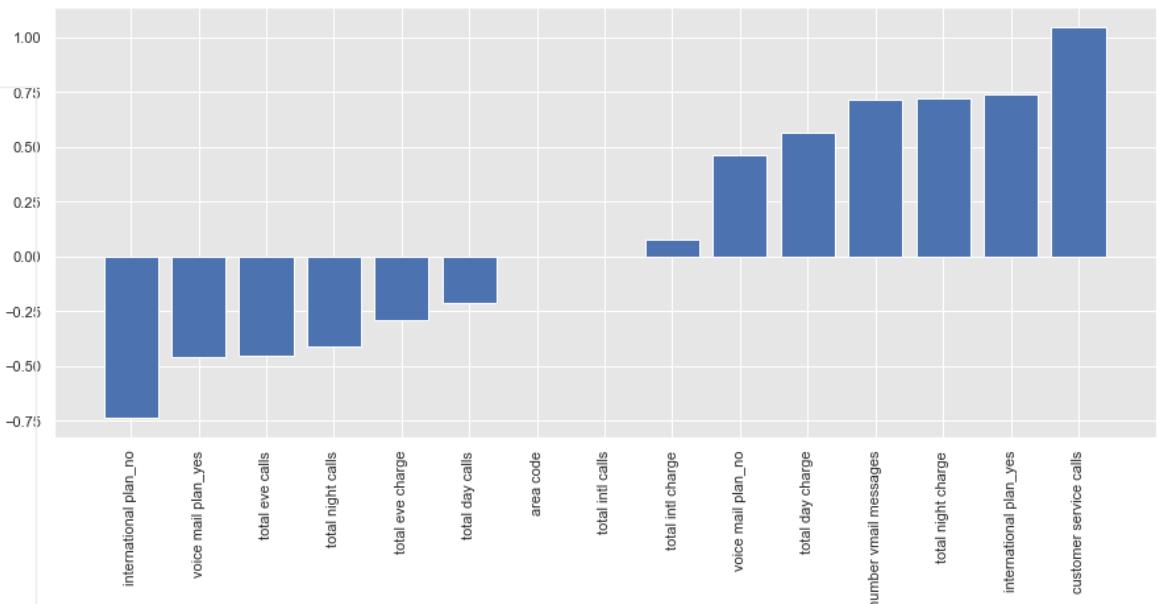
Testing AUC: 0.74873



## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient BoC
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
  - ▼ 15.1.1 Hyp
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
    - 20.1.2 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [162]: 1 plot\_coefficients(scd)



## Contents ⚙️

10.2 Results
▼ 11 XGBoost
11.1 Results
▼ 12 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 Gradient Bc
14.1 Results
▼ 15 SVM Classi
▼ 15.1 Tuning
15.1.1 Hyp
15.1.1.1
▼ 15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results
▼ 18 Logistic Re
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
▼ 20.1 Most In
20.1.1 Top
20.1.2 Top
▼ 20.2 Top Po
20.2.1 To
20.2.2 To
▼ 21 Conclusion
▼ 21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
▼ 21.2 Decrea
21.2.1 Vol
21.2.2 Inte
22 Recomend
23 Future Wor

## 17.1 Results & Comments

The mean cross validation is 95%. We have a precision of 0.28 and recall of 0.81 for the test set.  
We don't have overfitting. It means we are able to detect the 81% of churn customers and 28% of the customers we say 'churn' was true.

## 18 Logistic Regression - RandomizedSearchCV

```

In [164]: 1 from sklearn.model_selection import RandomizedSearchCV
          2 from scipy.stats import uniform
          3
          4 logistic = LogisticRegression(solver='saga', tol=1e-2, max_iter=200,
                                         random_state=0)
          5
          6 distributions = dict(C=uniform(loc=0, scale=4),
                                 penalty=['l2', 'l1'])
          7
          8 clf = RandomizedSearchCV(logistic, distributions, random_state=0)
          9 search = clf.fit(scaled_X_train, y_train)
         10
         11 search.best_params_

```

Out[164]: {'C': 1.75034884505077, 'penalty': 'l2'}

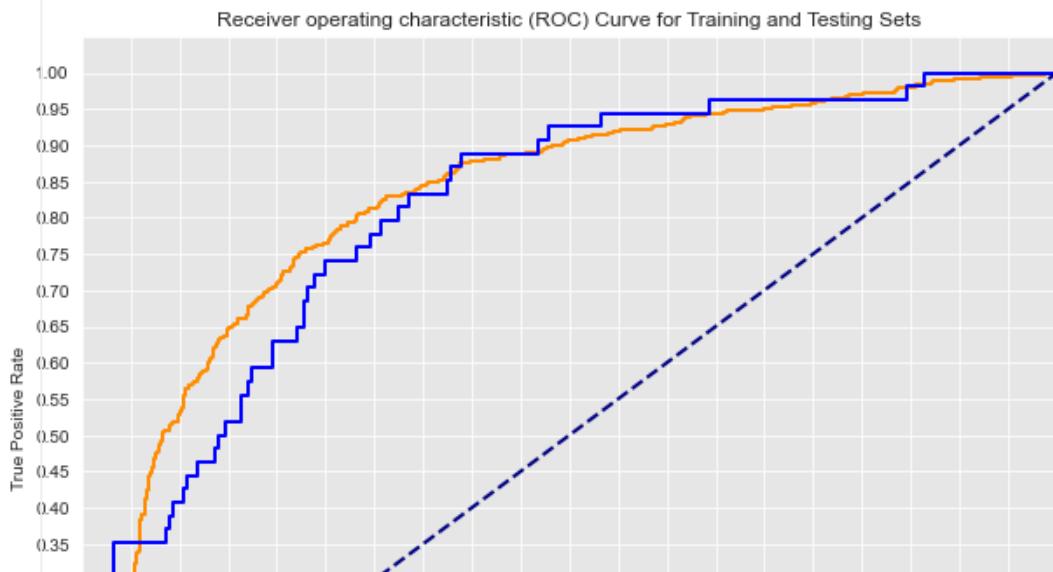
```
In [165]: 1 logistic = LogisticRegression(solver='saga', tol=1e-2, max_iter=200,
2                                         random_state=0)
3 distributions = dict(C=[1.5], penalty=['l2'])
4
5 clf = RandomizedSearchCV(logistic, distributions, random_state=0)
6 search = clf.fit(scaled_X_train, y_train)
```

## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost:
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
    - 20.1.2 To
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas This model had a mean cross validation of 85.66%. Only 20% of the time the churn was true with
    - 21.1.1 Tot 61 % precision.
    - 21.1.2 Cus
    - 21.1.3 Tot
    - 21.1.4 Tot
  - ▼ 21.2 Decrea
    - 21.2.1 Voi
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

Mean Cross Validation Score:

0.8566



## 18.1 Results & Comments

- ▼ 21 Conclusion
  - ▼ 21.1 Increas This model had a mean cross validation of 85.66%. Only 20% of the time the churn was true with
    - 21.1.1 Tot 61 % precision.
    - 21.1.2 Cus
    - 21.1.3 Tot
    - 21.1.4 Tot
  - ▼ 21.2 Decrea
    - 21.2.1 Voi
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

## 19 SGDClassifier - RandomizedSearchCV

In [167]:

```

1 from time import time
2 import scipy.stats as stats
3 from sklearn.utils.fixes import loguniform
4
5 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
6 from sklearn.datasets import load_digits
7 from sklearn.linear_model import SGDClassifier
8
9 # build a classifier
10 clf = SGDClassifier(loss='hinge', penalty='elasticnet',
11                      fit_intercept=True)
12
13
14 # Utility function to report best scores
15 def report(results, n_top=3):
16     for i in range(1, n_top + 1):
17         candidates = np.flatnonzero(results['rank_test_score'] == i)
18         for candidate in candidates:
19             print("Model with rank: {0}".format(i))
20             print("Mean validation score: {0:.3f} (std: {1:.3f})"
21                  .format(results['mean_test_score'][candidate],
22                         results['std_test_score'][candidate]))
23             print("Parameters: {0}".format(results['params'][candidate]))
24             print("")
25
26
27 # specify parameters and distributions to sample from
28 param_dist = {'average': [True, False],
29                 'l1_ratio': stats.uniform(0, 1),
30                 'alpha': loguniform(1e-4, 1e0)}
31
32 # run randomized search
33 n_iter_search = 20
34 random_search = RandomizedSearchCV(clf, param_distributions=param_dist,
35                                     n_iter=n_iter_search)
36
37 start = time()
38 random_search.fit(scaled_X_train, y_train)
39 print("RandomizedSearchCV took %.2f seconds for %d candidates"
40       " parameter settings." % ((time() - start), n_iter_search))
41 report(random_search.cv_results_)
42

```

RandomizedSearchCV took 0.88 seconds for 20 candidates parameter settings.

Model with rank: 1

Mean validation score: 0.863 (std: 0.006)

Parameters: {'alpha': 0.0006626483822186211, 'average': False, 'l1\_ratio': 0.503491101662517}

Model with rank: 2

Mean validation score: 0.860 (std: 0.008)

Parameters: {'alpha': 0.0005634122292363634, 'average': False, 'l1\_ratio': 0.7403534721013764}

Model with rank: 3

## Contents

10.2 Results
▼ 11 XGBoost
11.1 Results
▼ 12 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 Gradient Bo
14.1 Results
▼ 15 SVM Classi
▼ 15.1 Tuning
▼ 15.1.1 Hyp
15.1.1.1
▼ 15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results
▼ 18 Logistic Re
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
▼ 20.1 Most In
20.1.1 Top
▼ 20.2 Top Po
20.2.1 To
20.2.2 Toj
▼ 21 Conclusion
▼ 21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
▼ 21.2 Decreas
21.2.1 Voi
21.2.2 Inte
22 Recomend
23 Future Wor

```
Mean validation score: 0.857 (std: 0.001)
Parameters: {'alpha': 0.04172285339373269, 'average': False, 'l1_ratio': 0.6586706021204276}

Model with rank: 3
Mean validation score: 0.857 (std: 0.001)
Parameters: {'alpha': 0.062422867830419115, 'average': False, 'l1_ratio': 0.9302673127586086}
```

## Contents ⚙️⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boosting
  - 14.1 Results
- ▼ 15 SVM Classification
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hypothesis
    - 15.1.1.1 Results
  - 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassification
  - 17.1 Results
- ▼ 18 Logistic Regression
  - 18.1 Results
- ▼ 19 SGDClassifiers
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interpolation
  - ▼ 20.1 Most Influential Features
    - 20.1.1 Top 10
    - In [168]: 20.2 Top 10 Po
    - 20.2.1 Total
    - 20.2.2 Total
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases
    - 21.2.1 Void
    - 21.2.2 Intercept
- 22 Recomendations
- 23 Future Work

## 19.1 Tuning & Optimization

In [168]:

```
1 random_search.best_params_
```

Out[168]:

```
{'alpha': 0.0006626483822186211,
 'average': False,
 'l1_ratio': 0.503491101662517}
```

```
In [169]: # build a classifier
  1 clf = SGDClassifier(loss='hinge', penalty='elasticnet',
  2                         fit_intercept=True)
  3 param_dist = {'alpha': [0.0006626483822186211], 'average': [False], 'l1_
  4 ratio': [0.503491101662517]}
  5
  6 # run randomized search
  7 n_iter_search = 20
  8 random_search = RandomizedSearchCV(clf, param_distributions=param_dist,
  9                                     n_iter=n_iter_search)
 10
 11 start = time()
 12 rs=random_search.fit(scaled_X_train, y_train)
 13 print("RandomizedSearchCV took %.2f seconds for %d candidates"
 14         " parameter settings." % ((time() - start), n_iter_search))
 15 report(random_search.cv_results_)
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results

```
In [170]: # model evaluation
  1 y_pred_train = rs.predict(scaled_X_train)
  2 y_pred_test = rs.predict(scaled_X_test)
  3
  4
  5 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred_
  6 test)
  7 roc_curve_and_auc(rs, scaled_X_train, scaled_X_test, y_train, y_test)
  8 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
```

Test set					
	precision	recall	f1-score	support	
20.1 Most In	False	0.85	1.00	0.92	566
20.1.1 Top	True	0.00	0.00	0.00	101
20.2 Top Po	accuracy			0.85	667
20.2.1 To	macro avg	0.42	0.50	0.46	667
20.2.2 Toj	weighted avg	0.72	0.85	0.78	667

-----

Cohen's Kappa for train and test set:  
0.0045 0.0

f2 score for train and test set:  
0.0033 0.0

roc auc score for train and test set:  
0.5013 0.5

Mean Cross Validation Score:  
0.6972

## 19.2 Results & Comments

This model performed poorly.

# 20 Data Interpretation

## Contents ↻ ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

```
In [189]: N
1 dt_param_grid2 = {
2     'criterion': ['gini'],
3     'max_depth': [6],
4     'min_samples_split': [4],
5     'min_samples_leaf': [10]
6 }
7
8 dt_grid_search2 = GridSearchCV(tree_clf, dt_param_grid2, cv=10, return_
9
10 # Fit to the data
11
12 clf = dt_grid_search2 .fit(scaled_X_train, y_train)
13
14
15 #Examine best parameters
16
17 # Mean training score
18 dt_gs_training_score2 = np.mean(clf.cv_results_['mean_train_score'])
19
20 # Mean test score
21 dt_gs_testing_score2 = clf.score(scaled_X_test, y_test)
22
23
24 print(f"Mean Training Score: {dt_gs_training_score2 :.2%}")
25 print(f"Mean Test Score: {dt_gs_testing_score2 :.2%}")
26
27 # Model Performance
28 # Test set predictions
29
30 y_pred_test = clf.predict(scaled_X_test)
31 y_pred_train = clf.predict(scaled_X_train)
32
33 print(classification_report(y_test, y_pred_test))
34
35 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred_
36
37
38
39 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
40
41 plot_feature_importances(tree_clf)
```

	true	0.0%	0.0%	0.0%	101
accuracy				0.94	667
macro avg	0.92	0.83	0.87		667
weighted avg	0.94	0.94	0.93		667

#### MODEL EVALUATION METRICS:

-----

#### Confusion Matrix for train & test set:

	Predicted	No Churn	Churn	All
0	Actual	Nan	Nan	Nan
1	No Churn	2262.0	22.0	2284.0
2	Churn	104.0	278.0	382.0

	All	2366.0	300.0	2666.0
--	-----	--------	-------	--------

	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	558.0	8.0	566.0
2	Churn	22.0	52.0	74.0

## Contents ↴

### 20.1.1 Top 5 Important Features Affecting Churn Rate

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

In [196]:

```

1 feature_importance = tree_clf.feature_importances_
2 print (tree_clf.feature_importances_)
3 feat_importances = pd.Series(tree_clf.feature_importances_, index=X.columns)
4 feat_importances = feat_importances.nlargest(5)
5 feat_importances.plot(kind='barh', figsize=(10,10), color="b")
6
7 plt.xlabel('Feature importance', fontsize=15)
8 plt.ylabel('Positive Churn Factors', fontsize=15)
9 plt.title('Top 5 important features that affect positively on the churn rate')

```

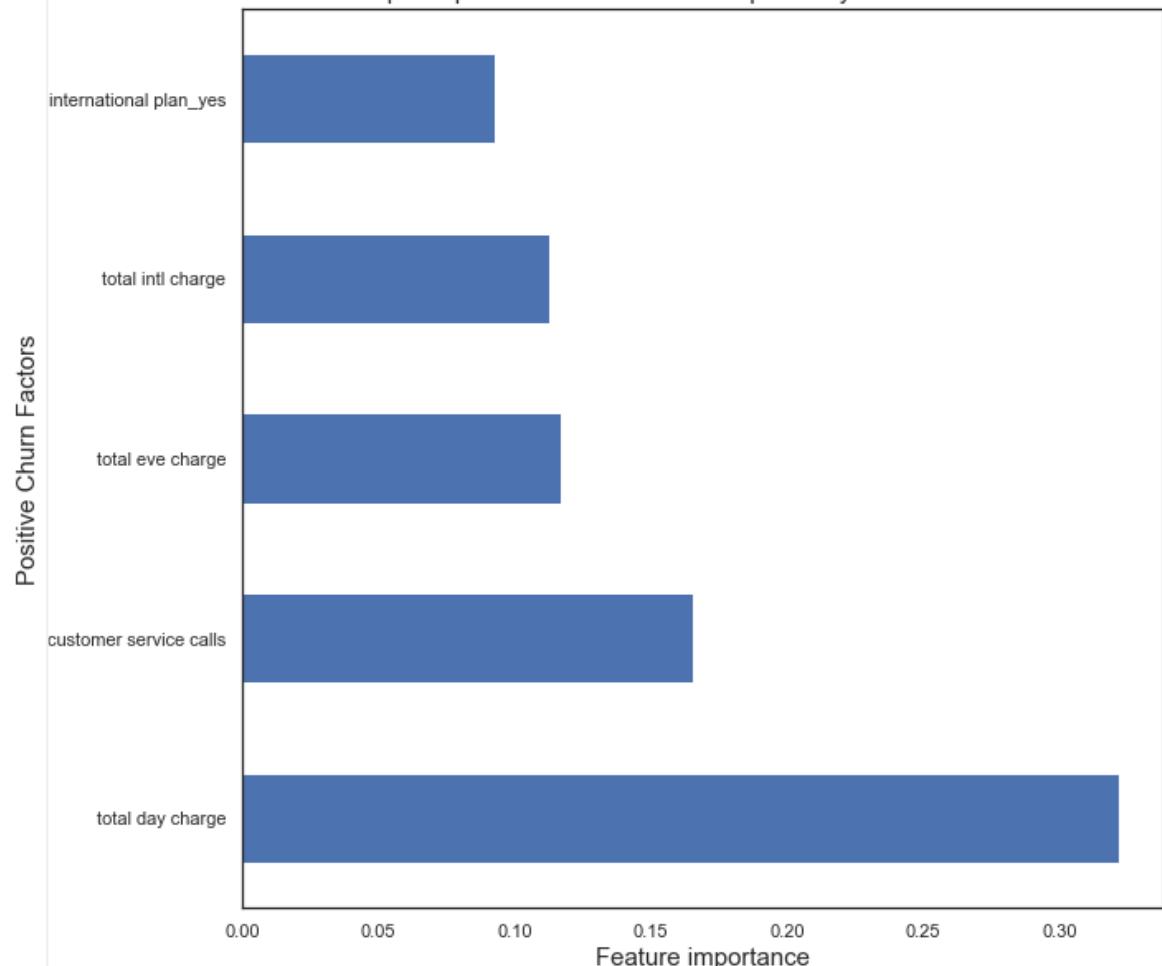
**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Boosting
  - 14.1 Results
- ▼ 15 SVM Classification
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassification
  - 17.1 Results
- ▼ 18 Logistic Regression
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interpolation
  - ▼ 20.1 Most Influential Features
    - 20.1.1 Top 5
    - 20.1.2 Top 10
  - ▼ 20.2 Top Positive Correlations
    - 20.2.1 Total
    - 20.2.2 Total
- ▼ 21 Conclusion
  - ▼ 21.1 Increases in Churn Rate
    - 21.1.1 Total
    - 21.1.2 Customer Service Calls
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases in Churn Rate
    - 21.2.1 Voice Minutes
    - 21.2.2 Internet Usage
- 22 Recomendations
- 23 Future Work

```
[0.          0.          0.          0.32193796  0.0134261  0.11687289
 0.00422496  0.02131784  0.08281552  0.11242512  0.1652558   0.
 0.0928464   0.          0.0688774 ]
```

**Out[196]:** Text(0.5, 1.0, 'Top 5 important features that affect positively on the churn rate')

Top 5 important features that affect positively on the churn rate



## 20.2 Top Positive and Negative Features According to Logistic Regression

```
In [178]: N
1 lr = LogisticRegression(fit_intercept=False, C=100,
2                           solver='liblinear', class_weight='balanced')
3 lr.fit(scaled_X_train, y_train)
4 y_pred_train = lr.predict(scaled_X_train)
5 y_pred_test = lr.predict(scaled_X_test)
6
7
8 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred)
9
10 roc_curve_and_auc(lr, scaled_X_train, scaled_X_test, y_train, y_test)
11
12 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)
```

## Contents ⚙️

10.2 Results	
▼ 11 XGBoost	
11.1 Results	
▼ 12 Gaussian N	
12.1 Results	
▼ 13 AdaBoostC	
13.1 Results	
▼ 14 Gradient Bo	
14.1 Results	
▼ 15 SVM Classi	
▼ 15.1 Tuning	
▼ 15.1.1 Hyp	
15.1.1.1	
▼ 15.2 Results	
15.2.1 SVI	
15.2.2 SV	
15.3 Results	
▼ 16 Stacking	
16.1 Results	
▼ 17 SGDClassif	
17.1 Results	
▼ 18 Logistic Re	
18.1 Results	
▼ 19 SGDClassif	
19.1 Tuning	
19.2 Results	
▼ 20 Data Interp	
▼ 20.1 Most In	
20.1.1 Top	
▼ 20.2 Top Po	
20.2.1 To	
20.2.2 To	
▼ 21 Conclusion	
▼ 21.1 Increas	
21.1.1 Total	
21.1.2 Cus	
21.1.3 Total	
21.1.4 Total	
▼ 21.2 Decreas	
21.2.1 Voic	
21.2.2 Inte	
22 Recomendat	
23 Future Wor	

	Actual	No Churn	Churn	All
1	No Churn	1426.0	858.0	2284.0
2	Churn	48.0	334.0	382.0
3	All	1474.0	1192.0	2666.0

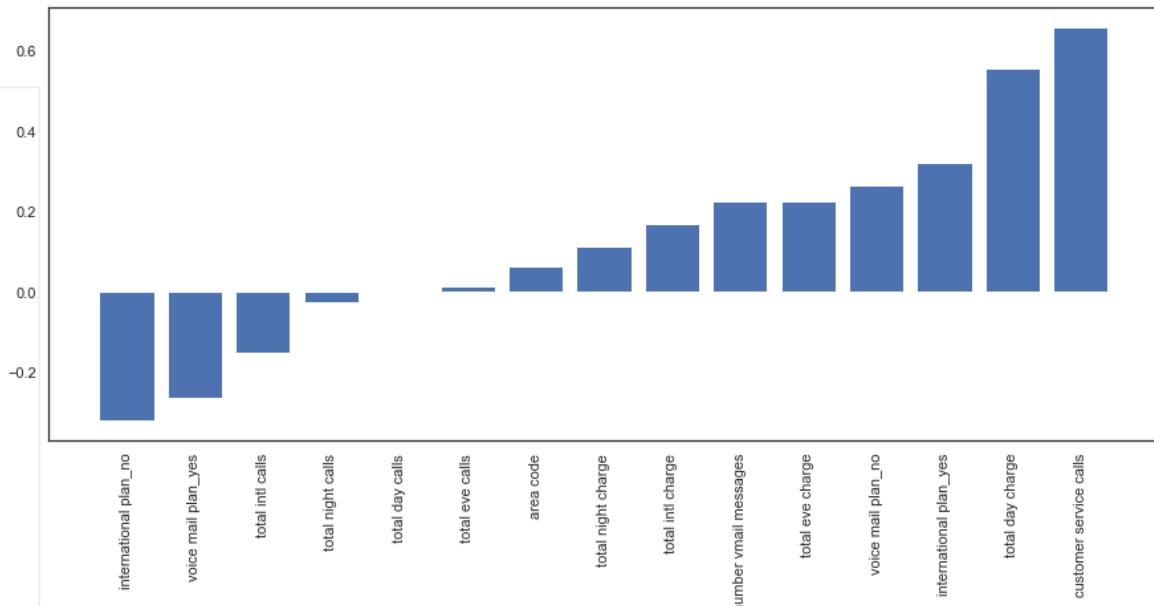
	Predicted	No Churn	Churn	All
0	Actual	NaN	NaN	NaN
1	No Churn	357.0	209.0	566.0
2	Churn	10.0	91.0	101.0
3	All	367.0	300.0	667.0

Classification Report for train & test set

Train set

	precision	recall	f1-score	support
False	0.97	0.62	0.76	2284
True	0.28	0.87	0.42	382

In [179]: 1 plot\_coefficients(lr)



## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
- In [180]: 1 lr.coef\_.round(2)
- Out[180]: array([[ 0.06, 0.23, 0. , 0.56, 0.01, 0.23, -0.03, 0.11, -0.15, 0.17, 0.66, -0.32, 0.32, 0.26, -0.26]])
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - In [181]: 1
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - In [182]: 1
    - 21.1.2 Cus
    - 21.1.3 Tot
    - 21.1.4 Tot
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor

```
1 from collections import OrderedDict
2 coefficients = pd.Series(lr.coef_[0], index=X.columns.values).to_dict()
3
4 abs_coefficients = {k: abs(v) for k, v in coefficients.items()}
5 coefs_ranked = OrderedDict(sorted(abs_coefficients.items(), key=lambda kv: kv[1]))
6 top_5 = list(coefs_ranked)[:5]
7 print("The 5 most important coefficients are:")
8 print(" / ".join(top_5))
```

The 5 most important coefficients are:

international plan\_no / voice mail plan\_yes / total intl calls / total night calls / total day calls

```
1 coefs_pos = {k: abs(v) for k, v in coefs_ranked.items() if v > 0}
2 coefs_neg = {k: abs(v) for k, v in coefs_ranked.items() if v < 0}
3 coefs_pos = OrderedDict(sorted(coefs_pos.items(), key=lambda kv: kv[1]),
```

## 20.2.1 Top Negative Factors affecting Churn Rate

In [183]:

```

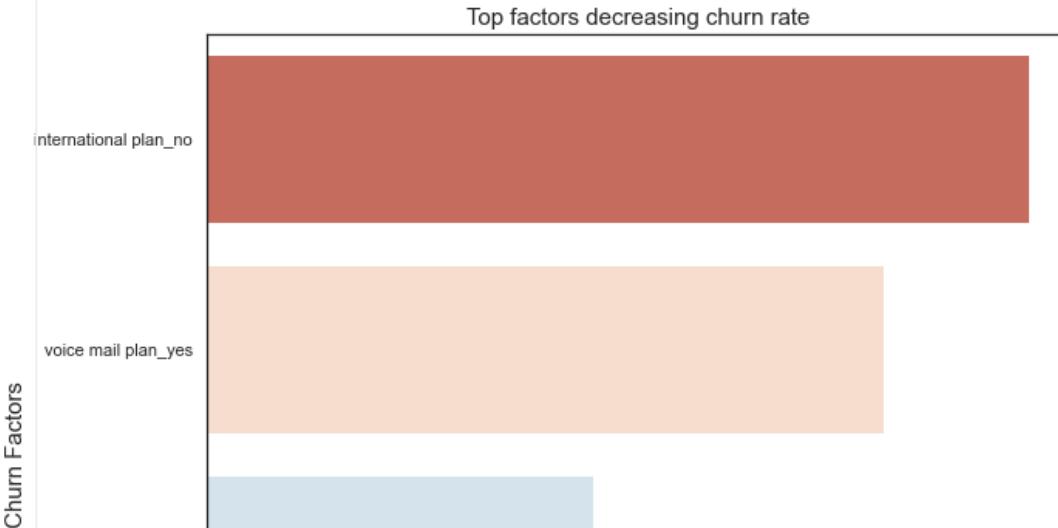
1 sns.set(rc={"figure.figsize":(10,10)})
2 sns.set_style("white")
3 sns.barplot(x=list(coefs_neg.values()), y=list(coefs_neg.keys()), palette="magma")
4 plt.title('Top factors decreasing churn rate', fontsize=15)
5 plt.xlabel("Importance Index", fontsize=15)
6 plt.ylabel('Churn Factors', fontsize=15)
7 print("\n\nBiggest Negative Factors on Churn")

```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
  - ▼ 16 Stacking
    - 16.1 Results
  - ▼ 17 SGDClassif
    - 17.1 Results
  - ▼ 18 Logistic Re
    - 18.1 Results
  - ▼ 19 SGDClassif
    - 19.1 Tuning
    - 19.2 Results
  - ▼ 20 Data Interp
    - ▼ 20.1 Most In
      - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 To
  - ▼ 21 Conclusion
    - ▼ 21.1 Increas
      - 21.1.1 Total
      - 21.1.2 Cus
      - 21.1.3 Total
      - 21.1.4 Total
    - ▼ 21.2 Decreas
      - 21.2.1 Voic
      - 21.2.2 Inte
  - 22 Recomend
  - 23 Future Wor

Biggest Negative Factors on Churn

**20.2.2 Top Positive Fators Affectig Churn Rate**

In [184]:

```

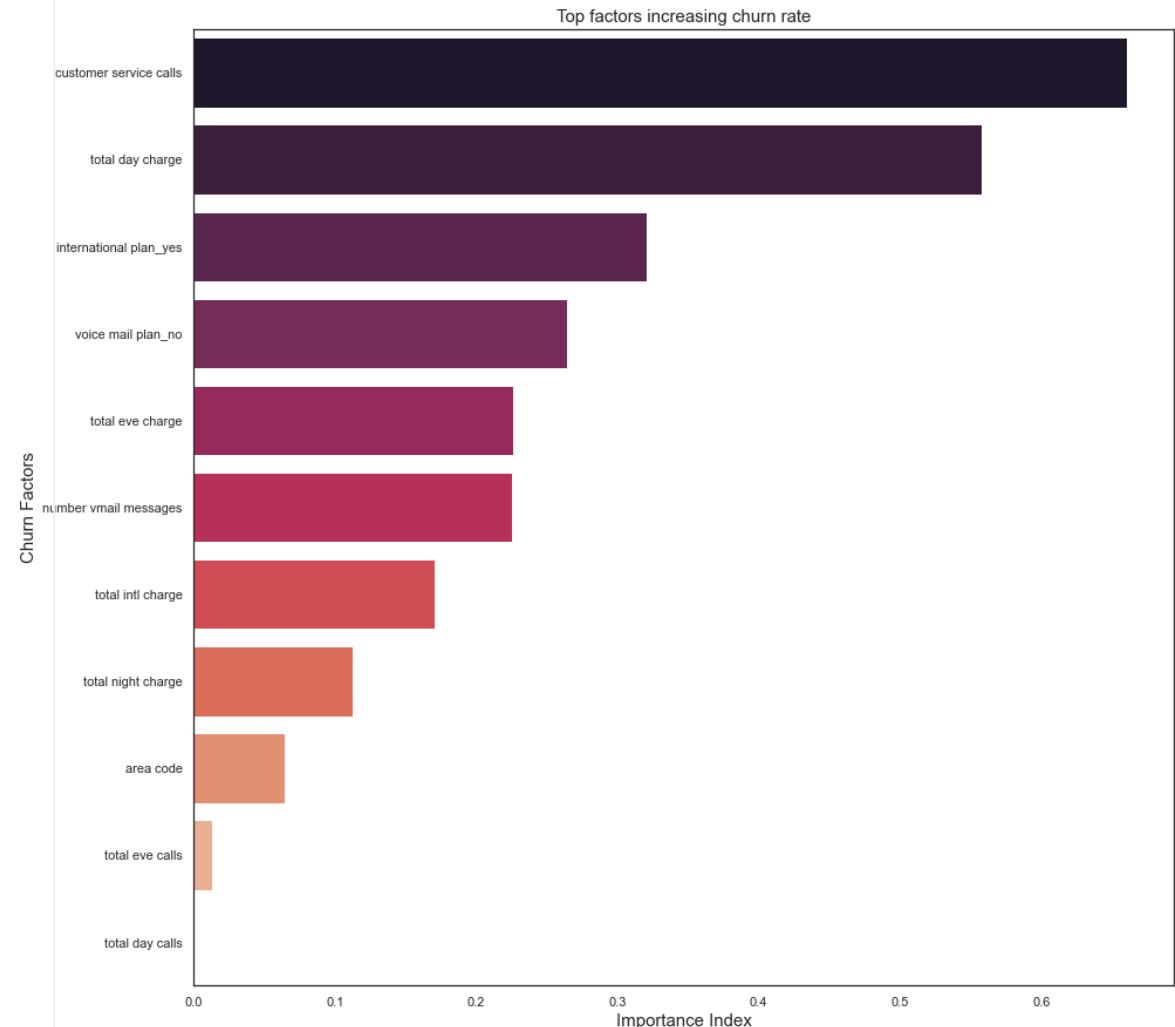
1 sns.set(rc={"figure.figsize":(15,15)})
2 sns.set_style("white")
3 sns.barplot(x=list(coefs_pos.values()), y=list(coefs_pos.keys()), palette="magma")
4 plt.title('Top factors increasing churn rate', fontsize=15)
5 plt.xlabel("Importance Index", fontsize=15)
6 plt.ylabel('Churn Factors', fontsize=15)
7 print("\n\nBiggest Positive Factors on Churn rate")

```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
    - 15.3 Results
  - ▼ 16 Stacking
    - 16.1 Results
  - ▼ 17 SGDClassif
    - 17.1 Results
  - ▼ 18 Logistic Re
    - 18.1 Results
  - ▼ 19 SGDClassif
    - 19.1 Tuning
    - 19.2 Results
  - ▼ 20 Data Interp
    - ▼ 20.1 Most In
      - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 To
  - ▼ 21 Conclusion
    - ▼ 21.1 Increas
      - 21.1.1 Total
      - 21.1.2 Cus
      - 21.1.3 Total
      - 21.1.4 Total
    - ▼ 21.2 Decreas
      - 21.2.1 Voic
      - 21.2.2 Inte
  - 22 Recomend
  - 23 Future Wor

Biggest Positive Factors on Churn rate



## Contents ⚙

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian NB
  - 12.1 Results
- ▼ 13 AdaBoostC The best performing machine learning model to predict the churn rate is for Stacking Classifier with recall of 74% and precision of 93%. However, we can not interpret this model and extract the important features.
  - 13.1 Results
- ▼ 14 Gradient Boosting
  - 14.1 Results
- ▼ 15 SVM Classifier NB Gaussian, Adaboost and Gradient Boost Classifier had the optimal recall and precision of 100%
  - ▼ 15.1 Tuning
  - ▼ 15.1.1 Hyp We decided to use DecisionTree classifier which performed well with 65% recall and 87% precision to interpret the data and extract the features.
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results Adaboost and GradientBoost, DecisionTree and RandoForest interpreted the data similarly.
- ▼ 16 Stacking Factors that contribute to churn of customers are:
  - 16.1 Results
- ▼ 17 SGDClassifier
  - 17.1 Results
- ▼ 18 Logistic Regression
  - 18.1 Results
- ▼ 19 SGDClassifier
  - 19.1 Tuning
  - 19.2 Results Factors that contribute to the retention of customers and churn prevention are:
- ▼ 20 Data Interpolation
  - 20.1 Most International calls
  - 20.1.1 Top 1) International plan = no
  - 20.2 Top Percentage 2) Voice mail plan = yes
  - 20.2.1 To 3) Total International calls
  - 20.2.2 Total
- ▼ 21 Conclusion
  - ▼ 21.1 Increasing Churn factors
  - 21.1.1 Total day charge
  - 21.1.2 Customer service calls
  - 21.1.3 Total evening charge
  - 21.1.4 Total international calls
- ▼ 21.2 Decreasing Churn factors
  - 21.2.1 Voice mail plan
  - 21.2.2 International calls
- 22 Recomendations
- 23 Future Work

In [10]:

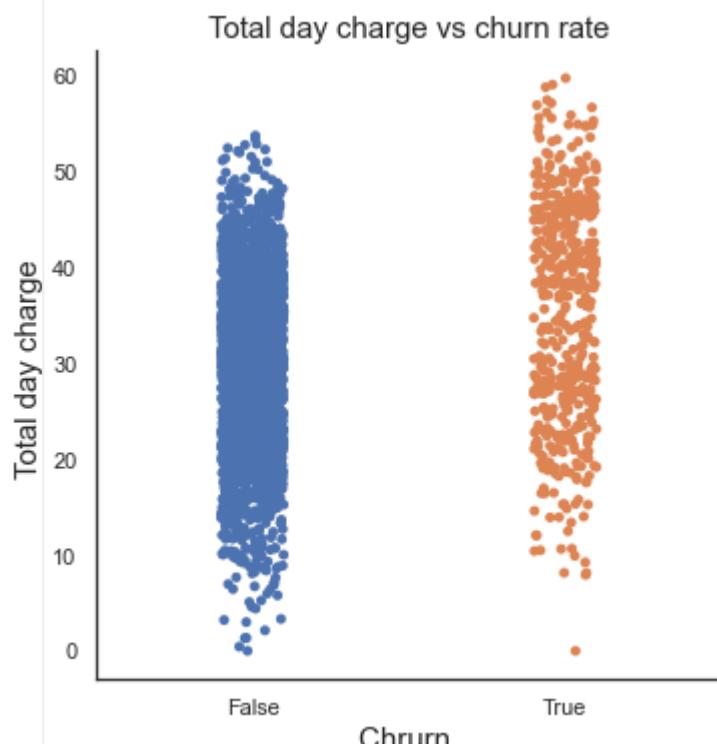
```

1 sns.set_style("white")
2 g=sns.catplot(data=df, x="churn", y="total day charge")
3 plt.title('Total day charge vs churn rate', fontsize=15)
4 plt.xlabel("Churn ", fontsize=15)
5 plt.ylabel('Total day charge ', fontsize=15)
6
7 spots = zip(g.ax.patches)
8 for spot in spots:
9     height = spot[0].get_height()
10    g.ax.text(spot[0].get_x(), height+3, '{:1.2f}'.format(churn.value_c

```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
      - 15.3 Results
  - ▼ 16 Stacking
    - 16.1 Results
  - ▼ 17 SGDClassif
    - 17.1 Results
  - ▼ 18 Logistic Re
    - 18.1 Results
  - ▼ 19 SGDClassif
    - 19.1 Tuning
    - 19.2 Results
  - ▼ 20 Data Interp
    - ▼ 20.1 Most In
      - 20.1.1 Top
    - ▼ 20.2 Top Po
      - 20.2.1 To
      - 20.2.2 To
  - ▼ 21 Conclusion
    - ▼ 21.1 Increas
      - 21.1.1 Total
      - 21.1.2 Cus
      - 21.1.3 Total
      - 21.1.4 Total
    - ▼ 21.2 Decreas
      - 21.2.1 Voic
      - 21.2.2 Inte
  - 22 Recomend
  - 23 Future Wor



In [266]:

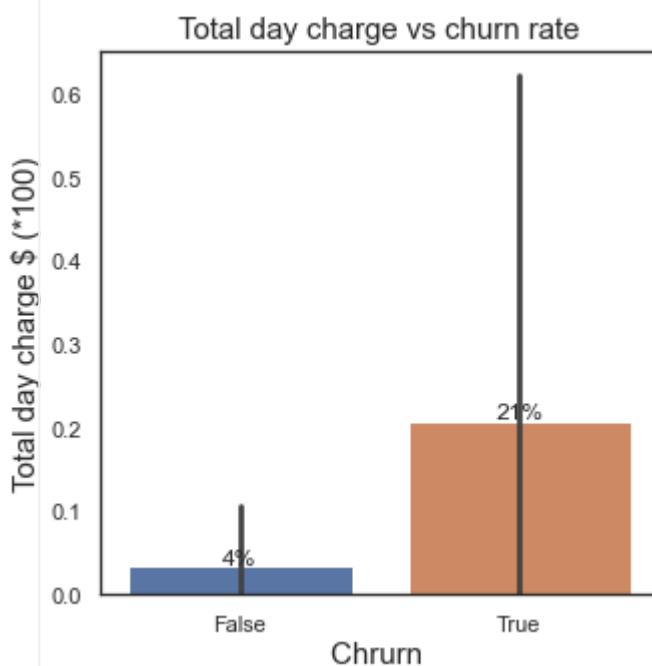
```

1 sns.set(rc={"figure.figsize":(5,5)})
2 sns.set_style("white")
3 ax=sns.barplot( x="churn", y="total day charge", data=df, estimator=lan
4 sns.set(rc={"figure.figsize":(5,5)})
5 plt.title('Total day charge vs churn rate', fontsize=15)
6 plt.xlabel("Chrurn ", fontsize=15)
7 plt.ylabel('Total day charge $ (*100)', fontsize=15)
8 for p in ax.patches:
9     width = p.get_width()
10    height = p.get_height()
11    x, y = p.get_xy()
12    ax.annotate(f'{height:.0%}', (x + width/2, y + height*1.02), ha='center')

```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



In [12]:

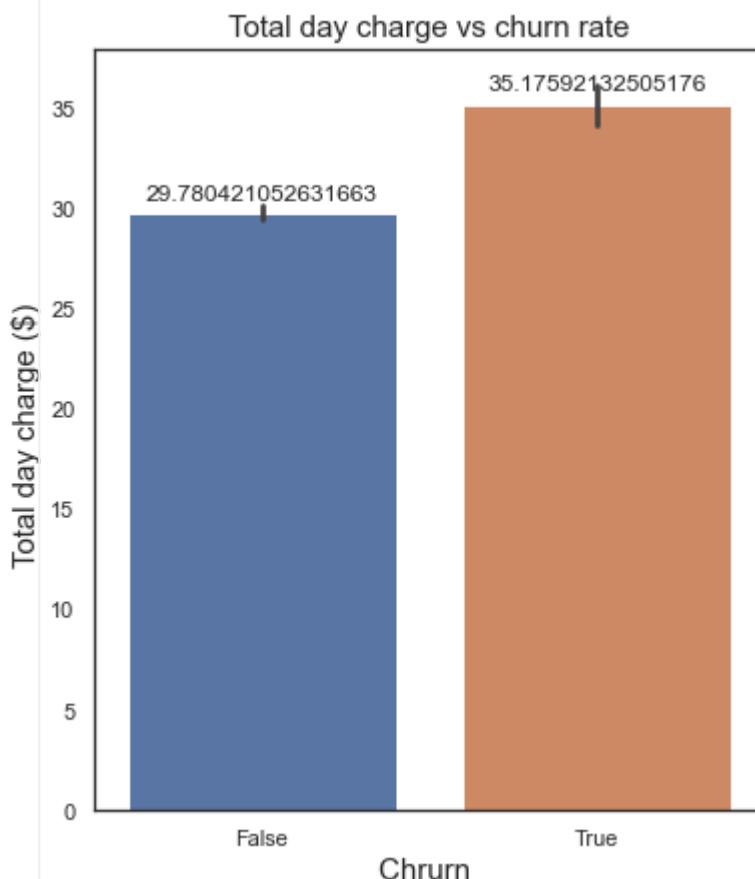
```

1 sns.set(rc={"figure.figsize":(6,7)})
2 sns.set_style("white")
3 ax=sns.barplot(data=df, x="churn", y="total day charge")
4 plt.title('Total day charge vs churn rate', fontsize=15)
5 plt.xlabel("Churn", fontsize=15)
6 plt.ylabel('Total day charge ($)', fontsize=15)
7 for p in ax.patches:
8     width = p.get_width()
9     height = p.get_height()
10    x, y = p.get_xy()
11    ax.annotate(f'{height}', (x + width/2, y + height*1.02), ha='center')
12

```

**Contents**

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



In [42]:

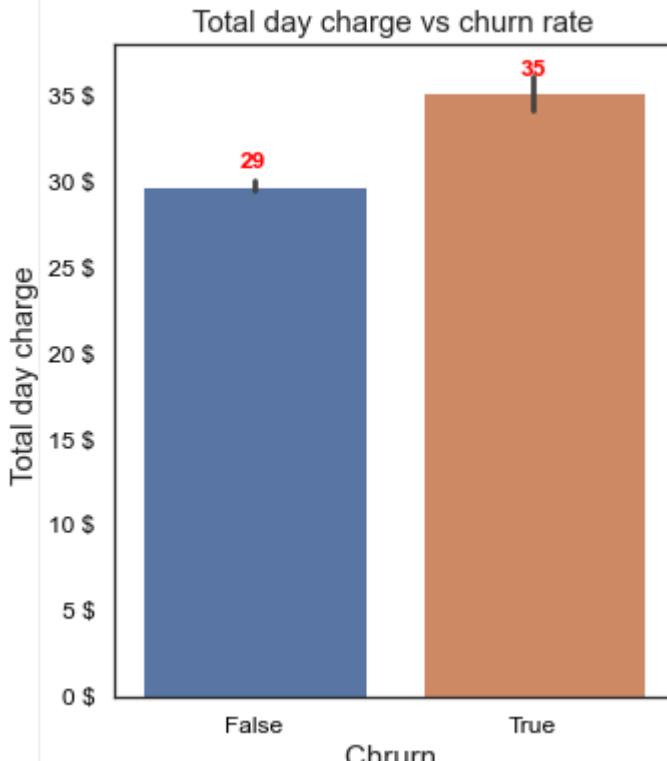
```

1 sns.set(rc={"figure.figsize":(5,6)})
2 sns.set_style("white")
3 ax=sns.barplot(data=df, x="churn", y="total day charge")
4 plt.title('Total day charge vs churn rate', fontsize=15)
5 plt.xlabel("Churn ", fontsize=15)
6 plt.ylabel('Total day charge ', fontsize=15)
7
8 yticks = np.arange(0, len(df['total day charge']), 1)
9
10
11
12
13 ylabel=["0 $", "5 $", "10 $", "15 $", "20 $", "25 $", "30 $", "35 $"]
14 xlabel=["False", "True"]
15 ax.set_yticklabels(ylabel, c='black', fontsize=12)
16 ax.set_xticklabels(xlabel, c='black', fontsize=12)
17
18 for i in ax.patches:
19     ax.text(i.get_x() + .4, i.get_height() + 1, int(i.get_height()), color='red')
20
21

```

## Contents

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
  - ▼ 15.1.1 Hyp
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



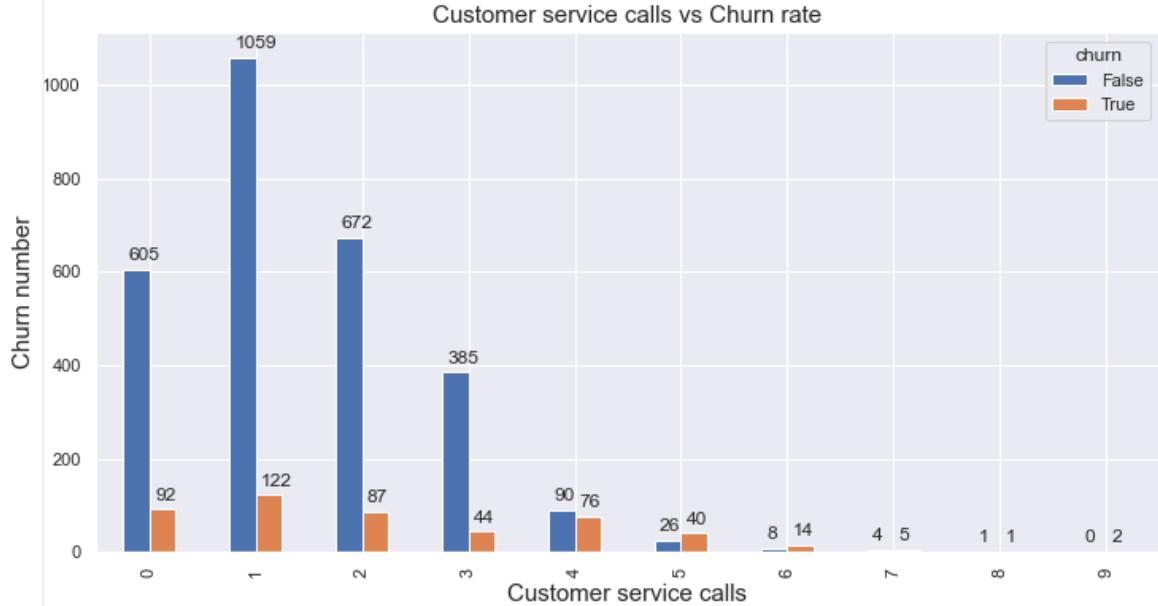
### 21.1.2 Customer service calls

```
In [216]: 1 cs=df.groupby(["customer service calls", "churn"]).size().unstack().plot(kind='bar', stacked=True)
2 for i in cs.patches:
3     cs.text(i.get_x()+0.05, i.get_height()+20,int(i.get_height()))
4
5 plt.xlabel('Customer service calls', fontsize=15)
6 plt.ylabel('Churn number', fontsize=15)
7 plt.title('Customer service calls vs Churn rate', fontsize=15)
```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 GaussianNB
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bo
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
    - ▼ 15.2 Results
      - 15.2.1 SVI
      - 15.2.2 SV
      - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - 21.1.1 Total
    - 21.1.2 Customer
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreases
    - 21.2.1 Voicemail
    - 21.2.2 Internet
- 22 Recomendations
- 23 Future Work

Out[216]: Text(0.5, 1.0, 'Customer service calls vs Churn rate')



In [300]:

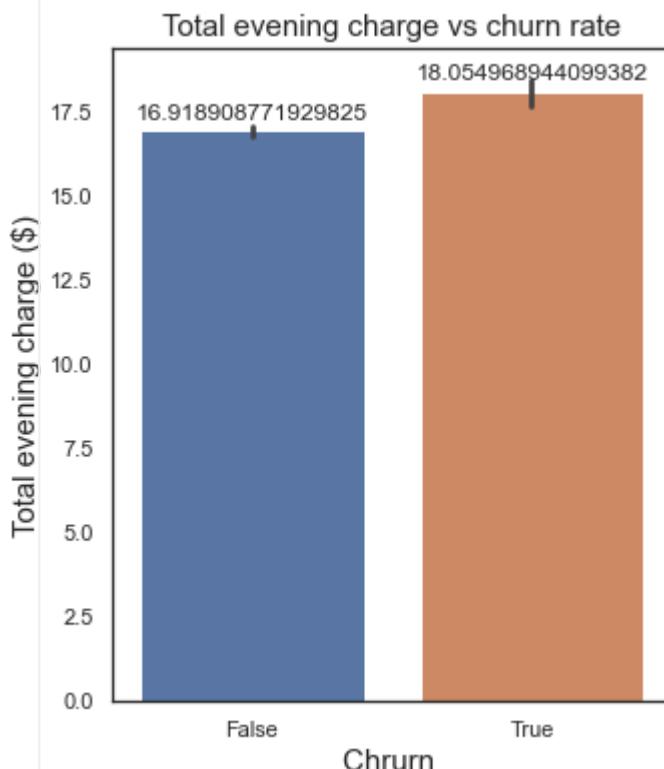
```

1 sns.set(rc={"figure.figsize":(5,6)})
2 sns.set_style("white")
3 ax=sns.barplot(data=df, x="churn", y="total eve charge")
4 #ax2=ax.twinx()
5 plt.title('Total evening charge vs churn rate', fontsize=15)
6 plt.xlabel("Churn", fontsize=15)
7 plt.ylabel('Total evening charge ($)', fontsize=15)
8
9
10 #ax2.set_yticklabels(( '$', '$', '$', '$', '$'))
11 for p in ax.patches:
12     width = p.get_width()
13     height = p.get_height()
14     x, y = p.get_xy()
15     ax.annotate(f'{height}', (x + width/2, y + height*1.02), ha='center')
16
17

```

## Contents ⚙️

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - ▼ 15.1.1 Hyp
      - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



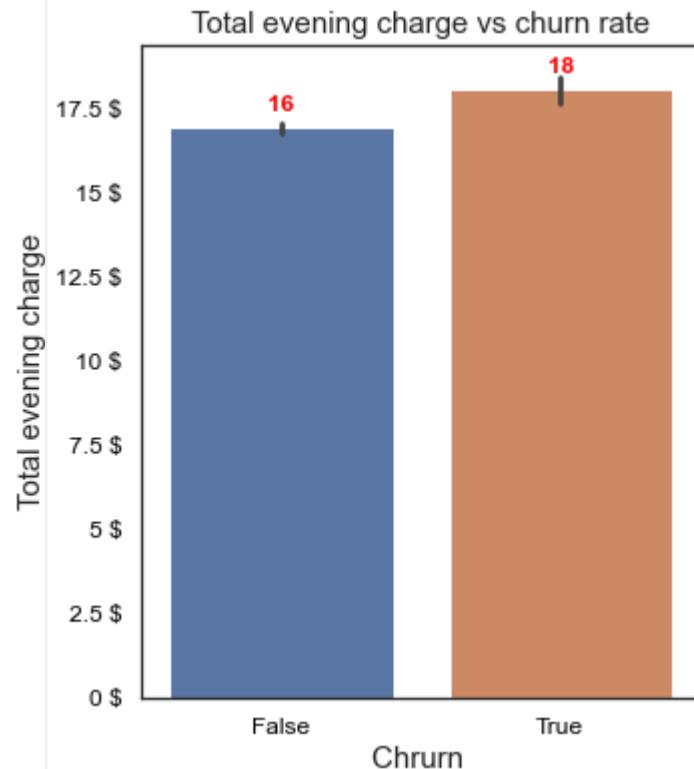
## Contents ⚙

- In [45]: 10.2 Results
- ▼ 11 XGBoost 11.1 Results
- ▼ 12 Gaussian N 12.1 Results
- ▼ 13 AdaBoostC 13.1 Results
- ▼ 14 Gradient Bo 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
  - ▼ 15.1.1 Hyp 15.1.1.1
  - ▼ 15.2 Results
  - 15.2.1 SVI
  - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking 16.1 Results
- ▼ 17 SGDClassif 17.1 Results
- ▼ 18 Logistic Re 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In 20.1.1 Top
  - ▼ 20.2 Top Po 20.2.1 To 20.2.2 To
- ▼ 21 Conclusion
  - ▼ 21.1 Increas 21.1.1 Total
  - 21.1.2 Cus
  - 21.1.3 Total
  - 21.1.4 Total
  - ▼ 21.2 Decre 21.2.1 Voic
  - 21.2.2 Inter
- 22 Recomend
- 23 Future Wor

```

1sns.set(rc={"figure.figsize":(5,6)})
2sns.set_style("white")
3ax=sns.barplot(data=df, x="churn", y="total eve charge")
4
5ticks = np.arange(0, len(df['total eve charge']), 1)
6
7
8plt.title('Total evening charge vs churn rate', fontsize=15)
9plt.xlabel("Churn", fontsize=15)
10plt.ylabel('Total evening charge', fontsize=15)
11
12label=["0 $", "2.5 $", "5 $", "7.5 $", "10 $", "12.5 $", "15 $", "17.5 $"]
13xlabel=["False", "True"]
14ax.set_yticklabels(label, c='black', fontsize=12)
15ax.set_xticklabels(xlabel, c='black', fontsize=12)
16
17for i in ax.patches:
18    ax.text(i.get_x()+.4,i.get_height()+.5,int(i.get_height()), color='red')
19
20
21

```



## 21.1.4 Total internatioanl charge

In [223]:

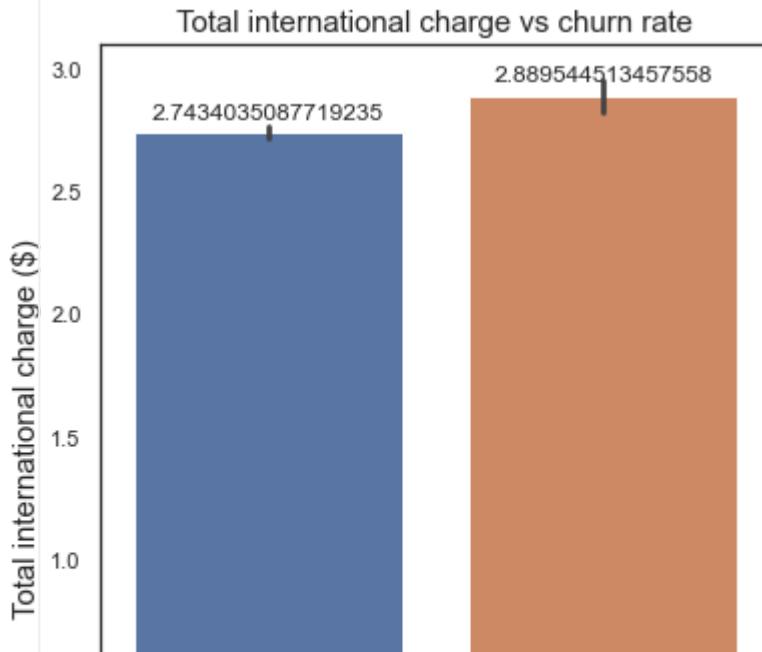
```

1 sns.set(rc={"figure.figsize":(6,7)})
2 sns.set_style("white")
3 ax=sns.barplot(data=df, x="churn", y="total intl charge")
4 plt.title('Total international charge vs churn rate', fontsize=15)
5 plt.xlabel("Churn ", fontsize=15)
6 plt.ylabel('Total international charge ($)' , fontsize=15)
7
8 for p in ax.patches:
9     width = p.get_width()
10    height = p.get_height()
11    x, y = p.get_xy()
12    ax.annotate(f'{height}', (x + width/2, y + height*1.02), ha='center')

```

### Contents

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decrea
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



In [32]:

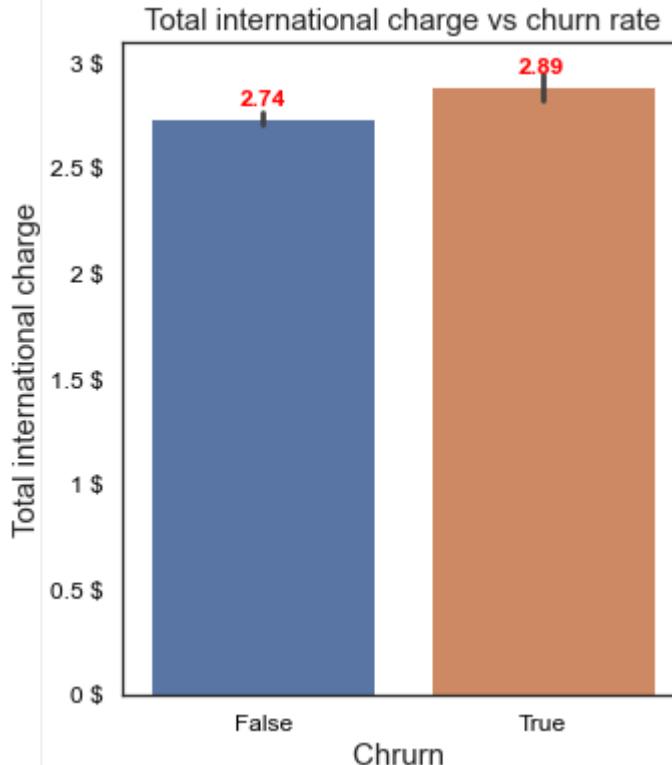
```

1 sns.set(rc={"figure.figsize":(5,6)})
2 sns.set_style("white")
3 ax=sns.barplot(data=df, x="churn", y="total intl charge")
4 plt.title('Total international charge vs churn rate', fontsize=15)
5 plt.xlabel("Churn ", fontsize=15)
6 plt.ylabel('Total international charge', fontsize=15)
7
8 yticks = np.arange(0, len(df['total intl charge']), 1)
9
10
11
12
13 ylabel=["0 $", "0.5 $", "1 $", "1.5 $", "2 $", "2.5 $", "3 $"]
14 xlabel=["False", "True"]
15 ax.set_yticklabels(ylabel, c='black', fontsize=12)
16 ax.set_xticklabels(xlabel, c='black', fontsize=12)
17
18 for p in ax.patches:
19     width = p.get_width()
20     height = p.get_height()
21     x, y = p.get_xy()
22     ax.annotate(f'{round(height,2)}', (x + width/2, y + height*1.02), color='red')
23

```

## Contents

- 10.2 Results
- ▼ 11 XGBoost
  - 11.1 Results
- ▼ 12 Gaussian N
  - 12.1 Results
- ▼ 13 AdaBoostC
  - 13.1 Results
- ▼ 14 Gradient Bc
  - 14.1 Results
- ▼ 15 SVM Classi
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
    - 15.1.1.1
  - ▼ 15.2 Results
    - 15.2.1 SVI
    - 15.2.2 SV
  - 15.3 Results
- ▼ 16 Stacking
  - 16.1 Results
- ▼ 17 SGDClassif
  - 17.1 Results
- ▼ 18 Logistic Re
  - 18.1 Results
- ▼ 19 SGDClassif
  - 19.1 Tuning
  - 19.2 Results
- ▼ 20 Data Interp
  - ▼ 20.1 Most In
    - 20.1.1 Top
  - ▼ 20.2 Top Po
    - 20.2.1 To
    - 20.2.2 Toj
- ▼ 21 Conclusion
  - ▼ 21.1 Increas
    - 21.1.1 Total
    - 21.1.2 Cus
    - 21.1.3 Total
    - 21.1.4 Total
  - ▼ 21.2 Decreas
    - 21.2.1 Voic
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



## Contents ⚙

10.2 Results
▼ 11 XGBoost
11.1 Results
▼ 12 Gaussian N
12.1 Results
▼ 13 AdaBoostC
13.1 Results
▼ 14 Gradient Bo
14.1 Results
▼ 15 SVM Classi
▼ 15.1 Tuning
▼ 15.1.1 Hyp
15.1.1.1
▼ 15.2 Results
15.2.1 SVI
15.2.2 SV
15.3 Results
▼ 16 Stacking
16.1 Results
▼ 17 SGDClassif
17.1 Results
▼ 18 Logistic Re
18.1 Results
▼ 19 SGDClassif
19.1 Tuning
19.2 Results
▼ 20 Data Interp
▼ 20.1 Most In
20.1.1 Top
▼ 20.2 Top Po
20.2.1 To
20.2.2 Toj
▼ 21 Conclusion
▼ 21.1 Increas
21.1.1 Total
21.1.2 Cus
21.1.3 Total
21.1.4 Total
▼ 21.2 Decreas
21.2.1 Voic
21.2.2 Inte
22 Recomend
23 Future Wor

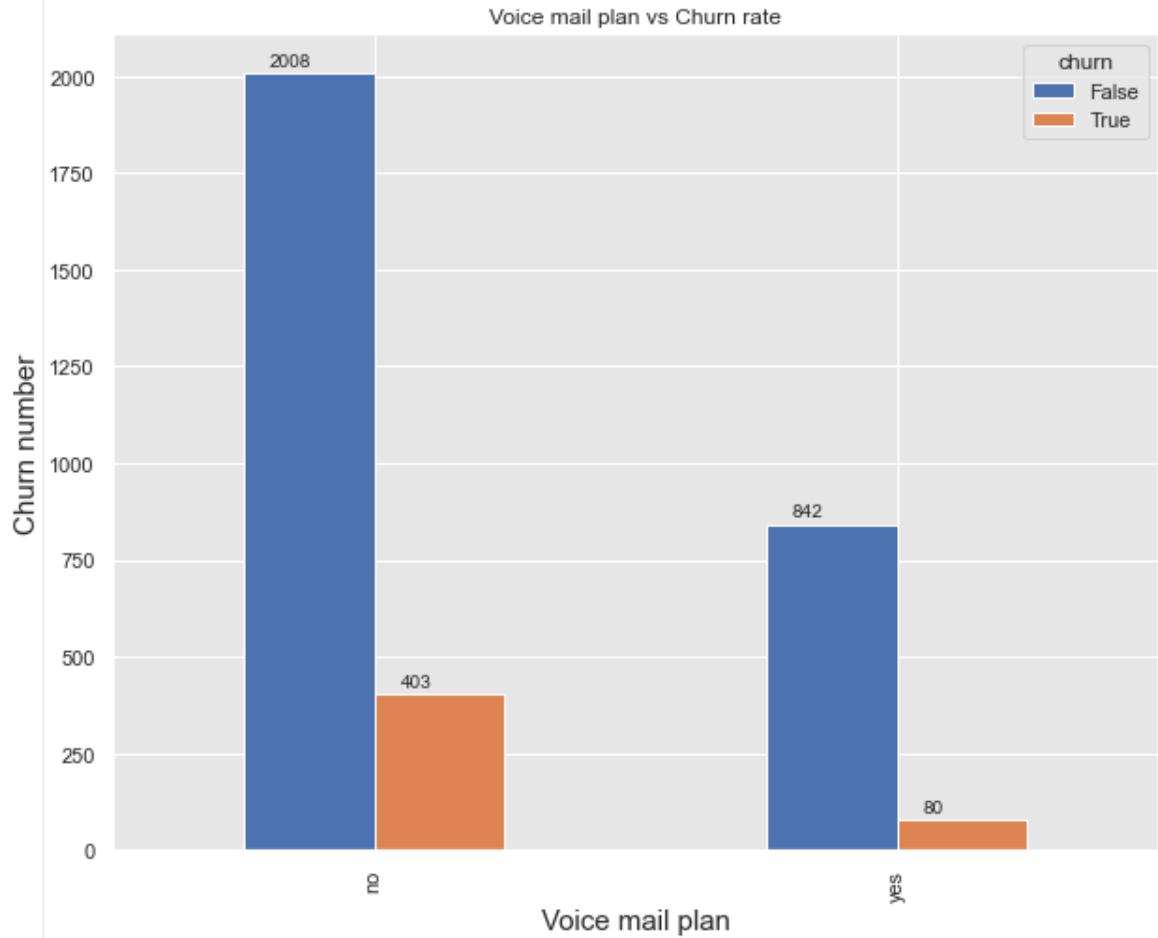
## 21.2 Decreasing churn factors

### 21.2.1 Voice mail plan

```
In [256]: 1 ax=df.groupby(["voice mail plan","churn"]).size().unstack().plot(kind='bar')
2 for i in ax.patches:
3     ax.text(i.get_x()+0.05, i.get_height()+20,str(i.get_height()))
4
5 plt.xlabel('Voice mail plan', fontsize=15)
6 plt.ylabel('Churn number', fontsize=15)
7 plt.title('Voice mail plan vs Churn rate')
```

## Contents ⚙

- 10.2 Results
- 11 XGBoost
- 11.1 Results
- 12 Gaussian N
- 12.1 Results
- 13 AdaBoostC
- 13.1 Results
- 14 Gradient Bo
- 14.1 Results
- 15 SVM Classi
- 15.1 Tuning
- 15.1.1 Hyp
- 15.1.1.1
- 15.2 Results
- 15.2.1 SVI
- 15.2.2 SV
- 15.3 Results
- 16 Stacking
- 16.1 Results
- 17 SGDClassif
- 17.1 Results
- 18 Logistic Re
- 18.1 Results
- 19 SGDClassif
- 19.1 Tuning
- 19.2 Results
- 20 Data Interp
- 20.1 Most In
- 20.1.1 Top
- 20.2 Top Po
- In [228]: 1 df.columns
- 21 Conclusion
- 21.1 Increase
- 21.1.1 Total
- 21.1.2 Cus
- 21.1.3 Total
- 21.1.4 Total
- 21.2 Decrease
- 21.2.1 Voic
- 21.2.2 Inte
- 22 Recomend
- 23 Future Wor



```
In [228]: 1 df.columns
```

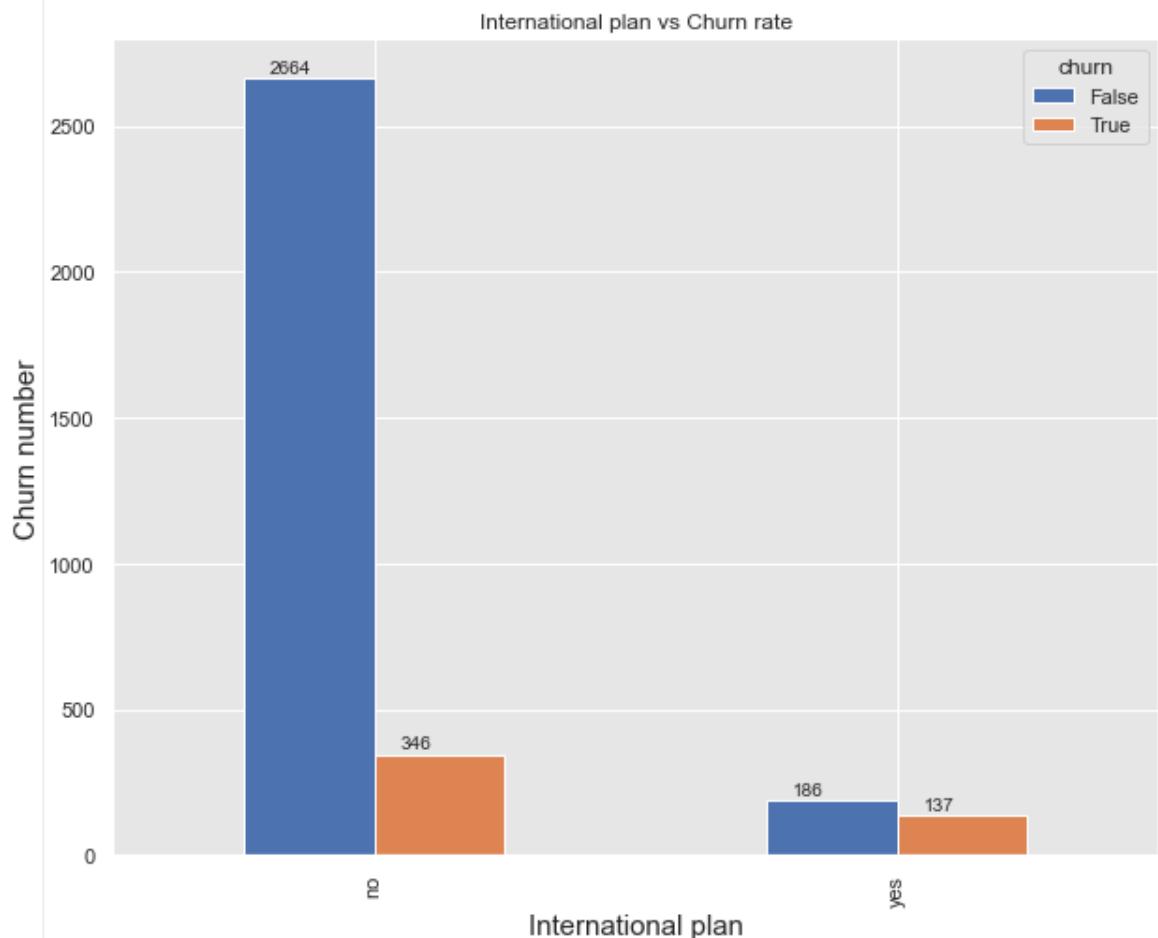
```
Out[228]: Index(['area code', 'number vmail messages', 'total day calls',
       'total day charge', 'total eve calls', 'total eve charge',
       'total night calls', 'total night charge', 'total intl calls',
       'total intl charge', 'customer service calls', 'churn',
       'international plan_no', 'international plan_yes', 'voice mail plan_no',
       'voice mail plan_yes'],
      dtype='object')
```

## 21.2.2 International plan

```
In [113]: 1 ac=df.groupby(['international plan','churn']).size().unstack().plot(kind='bar')
          2 for i in ac.patches:
          3     ac.text(i.get_x()+0.05, i.get_height()+20,str(i.get_height()))
          4 plt.xlabel('International plan', fontsize=15)
          5 plt.ylabel('Churn number', fontsize=15)
          6 plt.title('International plan vs Churn rate')
```

## Contents ⚙

- 10.2 Results
- 11 XGBoost
- 11.1 Results
- 12 Gaussian N
- 12.1 Results
- 13 AdaBoostC
- 13.1 Results
- 14 Gradient Bo
- 14.1 Results
- 15 SVM Classi
- 15.1 Tuning
- 15.1.1 Hyp
- 15.1.1.1
- 15.2 Results
- 15.2.1 SVI
- 15.2.2 SV
- 15.3 Results
- 16 Stacking
- 16.1 Results
- 17 SGDClassif
- 17.1 Results
- 18 Logistic Re
- 18.1 Results
- 19 SGDClassif
- 19.1 Tuning
- 19.2 Results
- 20 Data Interp
- 20.1 Most In
- 20.1.1 Top
- 20.2 Top Po
- 20.2.1 To
- 20.2.2 To
- 21 Conclusion



## 22 Recomendations

- 21.1 Increases
  - 21.1.1 Total
  - 21.1.2 Customer 1) Decrease call charge for day and special charges for evening calls.
  - 21.1.3 Total
  - 21.1.4 Total 2) Have a better international plan.
- 21.2 Decreases
  - 21.2.1 Voice 3) Customer service calls is a measure usually for bad experience. Look what are the main causes
  - 21.2.2 Internet of complaints and improve the quality of phone service.
- 22 Recomendations
- 23 Future Work 4) Keep up the good work with the voice mail plan which is the best feature that keeps customers in the company.

- 5) Compete with other companies on this strongly recommended feature "voice mail plan" and make better offers for the customers. This can attract customers from competing rivals

## 23 Future Work

### Contents

- More important data is needed to predict the churn of customers.
- 10.2 Results 1) Competitor Information
- ▼ 11 XGBoost
  - 11.1 Results Customer churn usually occurs when a rival company makes a good offer to our customers
- ▼ 12 Gaussian N
  - 12.1 Results Comparison of other companies phone charges (day, night, international) and internet service will be good factor to be used in the models to predict the churn of customers.
- ▼ 13 AdaBoostC
  - 13.1 Results 2) Internet package information
- ▼ 14 Gradient Bo
  - 14.1 Results Internet speed information, charge of internet per usage, the internet technology (fiber optic), internet package offer details including speed, data usage, TV streaming ...etc. All this information is important to include in our customer churn models.
  - ▼ 15.1 Tuning
    - 15.1.1 Hyp
      - 15.1.1.1 3) Contract Information
  - ▼ 15.2 Results
    - 15.2.1 SVI Month to month contract, 6 month contract, a year contact or 2 years contract. This information is important to know and definitely affects the churn or retention of customers.
    - 15.2.2 SV
    - 15.3 Results
- ▼ 16 Stacking 5) Payment method
  - 16.1 Results
- ▼ 17 SGDClassif
  - How the customers pay for the service: mailed check, automatic credit card payments, electronic check, automatic bank transfers. The less complex the payment method is, the more likely
- ▼ 18 Logistic Re
  - 18.1 Results customers will not churn???
- ▼ 19 SGDClassif
  - 6) Complaint data
    - 19.1 Tuning
    - 19.2 Results Text data from e-mails and customer service phone calls can be used to find the possible
- ▼ 20 Data Interp
  - complaints and solve them before they cause customer churn.
  - ▼ 20.1 Most In
    - 20.1.1 Top 7) Tenure data
  - ▼ 20.2 Top Po
    - 20.2.1 To Tenure data of customers is important throughout the years which can show us loyalty and churn
    - 20.2.2 To of customers.
- ▼ 21 Conclusion
  - ▼ 21.1 Increases
    - In [1]: Total ► 1
    - 21.1.1 Cus
    - 21.1.2 Tot
    - 21.1.3 Tot
    - 21.1.4 Tot
  - ▼ 21.2 Decreases
    - 21.2.1 Voi
    - 21.2.2 Inte
- 22 Recomend
- 23 Future Wor