

Table of Contents

1 Abstract
2 Business Problem For a Telecommunication Company
Contents ↗ ⚙
3 Import Libaries
1 Abstract
2 Business Problem
3 Import Libaries
4 Load Data
5 Data Exploration
6 Data Engineering
6.1 Group state column into south and north
6.1.1 Comments
6.2 Drop unnecessary columns
6.3 Drop highly correlated features
6.4 One Hot Encoder Categorical Columns
6.5 Train test split
6.6 Standardize Features
6.7 Building analysis functions
7 Logistic Regression
7.1 Optimization
7.2 Best performing parameters
7.3 Results & Comments
8 Decision Tree Classifier
8.1 Hyperparameter Tuning
8.1.1 Select Best Parameters for Decision tree classifier
8.2 Results & Comments
9 Random Forest Classifier
9.1 Hyperparameters
9.2 Best Parameters
9.3 Results & Comments
10 KNN Classifier
10.1 Tuning & Optimization
10.1.1 Select Best K for KNN Machine Learning Model
10.1.1.1 K=3
10.1.2 Iterate over n values
10.1.2.1 K=1
10.1.2.2 Iterate over 1 < n < 26
10.1.2.2.1 K=1
10.2 Results & Comments
11 XGBoost
11.1 Results & comments
12 Gaussian NB
12.1 Results & Comments
13 AdaBoostClassifier
13.1 Results & Comments
14 Gradient Boosting Classifier
14.1 Results & Comments
15 SVM Classifier
15.1 Tuning & Optimization
15.1.1 Hypertuning of SVM Classifier through GridSearchCV
15.1.1.1 Best parameters 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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comments
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comments
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comments
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal K
 - ▼ 10.1.1 Select Best Model
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterations
 - 10.1.2.1 K=1
 - 10.2 Results & Comments
- ▼ 11 XGBoost Classifier
 - 11.1 Results & Comments
- ▼ 12 Gaussian NB
- 12.1 Results & Comments
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comments
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comments
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal C
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model
 - 15.2 Results & Comments

1 Abstract

Customer churn in the telecommunication industry. In this project, machine learning models are used to predict possible customer churn in a telecommunication industry. Our models are powerful enough to determine churn factors. Therefore, we can give a company the advantage of taking preventive measures before losing customers.

2 Business Problem For a Telecommunication Company

Determine important factors affecting churn.

Determine factors increasing churn rate.

Determine factors decreasing churn rate.

3 Import Libaries

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state of origin
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Descion Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Parameters
 - ▼ 10.1.1 Select Euclidean Distance
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate over different values of K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Parameters
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

In [229]:

```

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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Common values
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standarize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Type
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

```

57 import seaborn as sns
58 sns.set_style('darkgrid', {'axes.facecolor': '0.9'})
59
60 from sklearn.preprocessing import LabelEncoder, OneHotEncoder, StandardScaler
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, accuracy_score
64 from sklearn.metrics import f1_score, fbeta_score, r2_score, roc_auc_score
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_precision_score
74
75 from sklearn.metrics import plot_confusion_matrix
76

```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building baseline models
- 7 Logistic Regression
 - 7.1 Optimization
 - In [280]:
 - 7.2 Best performance
 - 7.3 Results & Conclusion
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results & Conclusion
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Conclusion
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over error function
 - 10.1.2.1 K=1
 - 10.2 Results & Conclusion
- 11 XGBoost
 - 11.1 Results & Conclusion
- 12 Gaussian NB
 - 12.1 Results & Conclusion
- 13 AdaBoostClassifier
 - 13.1 Results & Conclusion
- 14 Gradient Boosting
 - 14.1 Results & Conclusion
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

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]:

Contents ↗

	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	AL	118	510	391-8027		yes	no	0 223.4
6	MA	121	510	355-9993		no	yes	24 218.2
7	MO	147	415	329-9001		yes	no	0 157.0
8	LA	117	408	335-4719		no	no	0 184.5
9	WV	141	415	330-8173		yes	yes	37 258.6

0 rows × 21 columns

In [232]: 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
turn bool
type: object

Data Exploration & Visualization

In [233]: 1 #Data Visualization

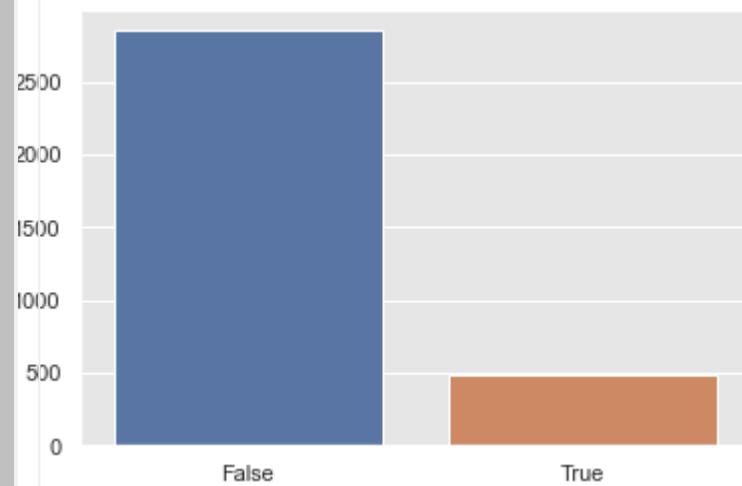
```
2
3 y = df["churn"].value_counts()
4 #print(y)
5 sns.barplot(y.index, y.values)
```

```
6
7 #sns.countplot(x='churn', data=churn_data)
8 #plt.show()
```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state column
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 8 Decision Tree Classifier
 - 8.1 Hyperparameter Tuning
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameter Tuning
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison

Out[233]: AxesSubplot:>



1 y_True = df["churn"][df["churn"] == True]

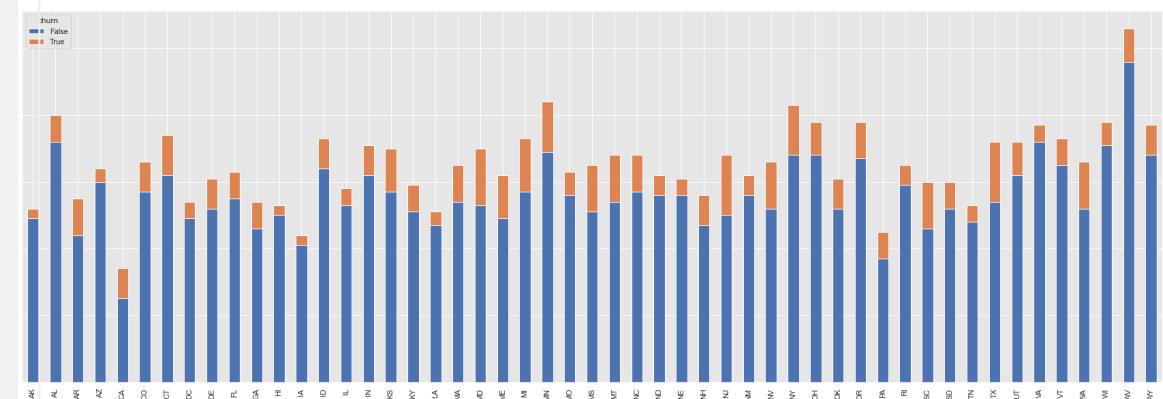
2 print ("Churn Percentage = "+str((y_True.shape[0] / df["churn"].shape[0]) * 100))

Churn Percentage = 14.491449144914492

1 #Churn by state

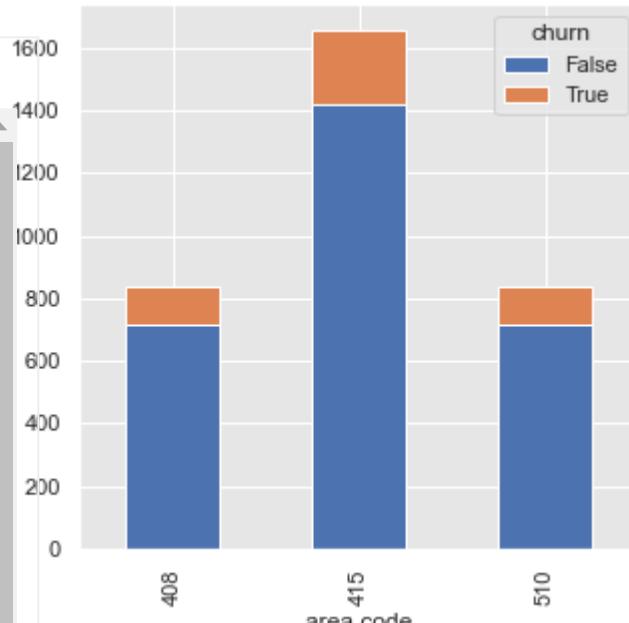
2 df.groupby(["state", "churn"]).size().unstack().plot(kind='bar', stacked=True)

Out[235]: AxesSubplot:xlabel='state'>



In [236]: 1 #Churn by area code
2 df.groupby(["area code", "churn"]).size().unstack().plot(kind='bar', st

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



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state country
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.1.2 Iterate different K values
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

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  
 5   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 ⚙

1 Abstract	1	state	3333 non-null	object
2 Business Problem	2	account length	3333 non-null	int64
3 Import Libraries	3	area code	3333 non-null	int64
4 Load Data	4	phone number	3333 non-null	object
5 Data Exploration	5	international plan	3333 non-null	object
6 Data Engineering	5	voice mail plan	3333 non-null	object
6.1 Group state by region	7	number vmail messages	3333 non-null	int64
6.1.1 Comment on results	8	total day minutes	3333 non-null	float64
6.2 Drop unnecessary columns	9	total day calls	3333 non-null	int64
6.3 Drop highly correlated columns	10	total day charge	3333 non-null	float64
6.4 One Hot Encoding	11	total eve minutes	3333 non-null	float64
6.5 Train test split	12	total eve calls	3333 non-null	int64
6.6 Standardize Features	13	total eve charge	3333 non-null	float64
6.7 Building analysis	14	total night minutes	3333 non-null	float64
7 Logistic Regression	15	total night calls	3333 non-null	int64
7.1 Optimization	16	total night charge	3333 non-null	float64
7.2 Best performance	17	total intl minutes	3333 non-null	float64
7.3 Results Comparison	18	total intl calls	3333 non-null	int64
8 Descion Tree Classification	19	total intl charge	3333 non-null	float64
8.1 Hyperparameter Tuning	20	customer service calls	3333 non-null	int64
8.1.1 Select Best Model		churn	3333 non-null	bool
8.2 Results Comparison		dtypes: bool(1), float64(8), int64(8), object(4)		
8.2.1 Memory usage:		memory usage:	524.2+ KB	

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

dtypes: object

```
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
turn                 0
type: int64
```

```
In [240]: 1 # check for duplicated rows
2 df.duplicated().sum()
3
```

Contents ↗⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state by country
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.1.2 Iterate over K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

In [15]:

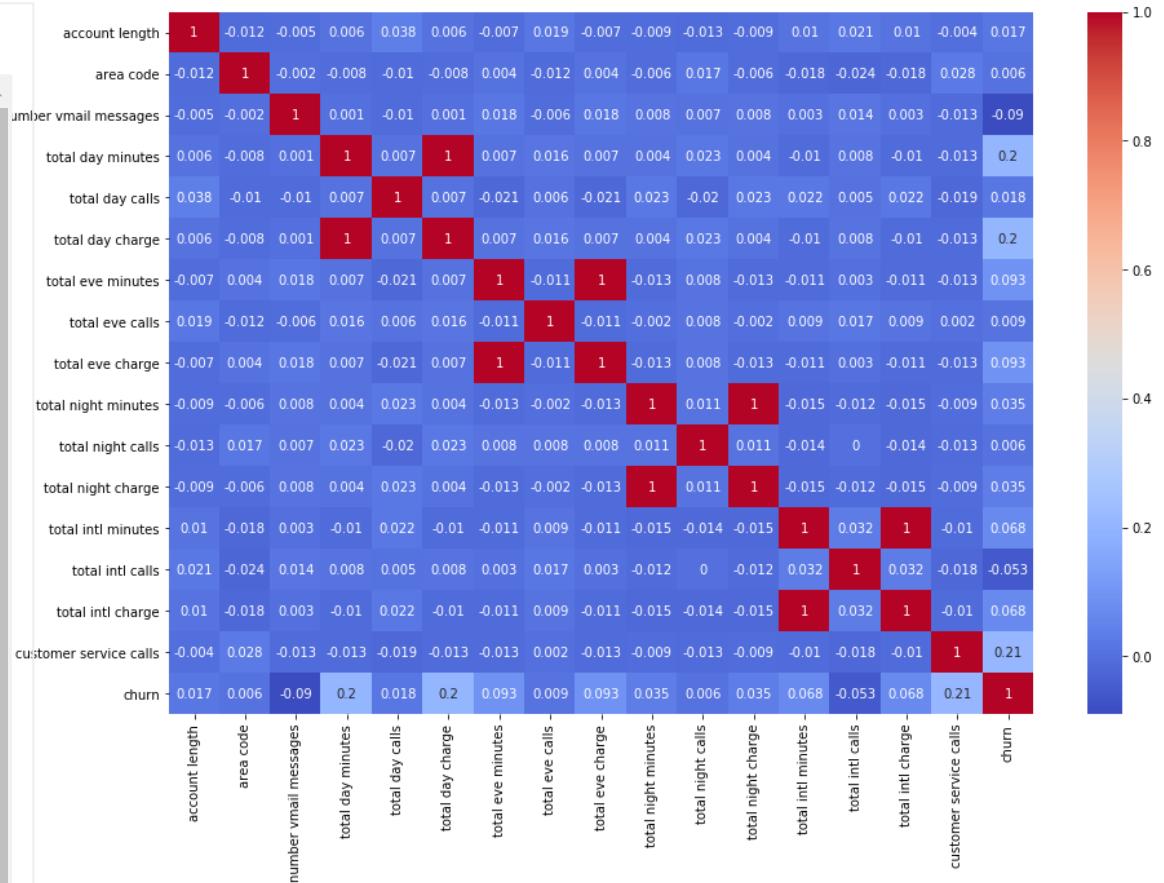
```

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

```

Contents ⚙

1	Abstract
2	Business Problem
3	Import Libraries
4	Load Data
5	Data Exploration
6	Data Engineering
6.1	Group state columns
6.1.1	Comments
6.2	Drop unnecessary columns
6.3	Drop highly correlated columns
6.4	One Hot Encoding
6.5	Train test split
6.6	Standardize Features
6.7	Building analysis
7	Logistic Regression
7.1	Optimization
7.2	Best performance
7.3	Results Comparison
8	Decision Tree Classifier
8.1	Hyperparameters
8.1.1	Select Best Parameters
8.2	Results Comparison
9	Random Forest Classifier
9.1	Hyperparameters
9.2	Best Parameters
9.3	Results & Comparison
10	KNN Classifier
10.1	Tuning & Optimization
10.1.1	Select Best Parameters
10.1.1.1	K=3
10.1.2	Iterate on K
10.1.2.1	K=1
10.2	Results & Comparison
11	XGBoost
11.1	Results & Comparison
12	Gaussian NB
12.1	Results & Comparison
13	AdaBoostClassifier
13.1	Results & Comparison
14	Gradient Boosting
14.1	Results & Comparison
15	SVM Classifier
15.1	Tuning & Optimization
15.1.1	Hypertuning
15.1.1.1	Best Parameters



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
type: object

```

Contents ⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state by country
 - 6.1.1 Comment on results
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Number of Neighbors
 - ▼ 10.1.1 Select Error Type
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate over different K values
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Parameters
 - ▼ 15.1.1 Hyperparameter Tuning
 - 15.1.1.1 Best Parameters

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

6 Data Engineering

In [268]: 1 df.groupby(['state', 'churn']).sum().head()

Contents C ↗

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state column
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 7.1 Optimization
- 7.2 Best performance
- 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Metrics
 - 10.1.1.1 K=3
 - 10.1.2 Iterations
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

		account length	area code	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes
state	churn							
AK	False	4640	21567	471	8796.3	4684	1495.39	9088.6
	True	414	1245	0	479.7	270	81.55	494.1
AL	False	7065	31063	545	13235.4	7156	2250.02	13770.8
	True	777	3387	29	1645.4	772	279.70	1866.2
AR	False	4173	19506	367	7721.0	4484	1312.59	8541.4

6. Group state column into south and north

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()

```
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()

Out[183]:

	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

Contents ↗

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
- 6.1 7.1 Optimization
- 7.2 Best performance
- 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
- ▼ 8.1 Hyperparameters

rows × 23 columns

Comments

to include state_north and state_south in the models but both features are insignificant. I decided to remove state feature.

In [83]: 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
10.1.2.1 K=1	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
turn	bool
state_south	bool
state_north	bool
type	object

6.2 Drop unnecessary columns

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

Contents

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
- 7 Logistic Regression
- 8 Descision Tree Classifier
- 9 Random Forest
- 10 KNN Classifier
- 11 XGBoost
- 12 Gaussian NB
- 13 AdaBoostClassifier
- 14 Gradient Boosting
- 15 SVM Classifier

```
1 df=df.drop(['state'], axis=1)
```

```
1 df=df.drop(['account length'], axis=1)
```

Drop highly correlated features

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

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

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

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

```
1 df.dtypes
```

state	object
area code	int64
international plan	object
voice mail plan	object
number vmail messages	int64
total day calls	int64
total day charge	float64
total eve calls	int64
total eve charge	float64
total night calls	int64
total night charge	float64
total intl calls	int64
total intl charge	float64
customer service calls	int64
turn	bool
type	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 ⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state data
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Model
 - ▼ 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Model
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

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

```
0    3010
1    323
Name: international plan, dtype: int64
```

Contents ⚙

1 Abstract	4
2 Business Problem	5
3 Import Libraries	5
4 Load Data	3
5 Data Exploration	7
6 Data Engineering	9
6.1 Group state of account	5
6.1.1 Comment	2
6.2 Drop unnecessary columns	3
6.3 Drop highly correlated columns	5
6.4 One Hot Encoding	3
6.5 Train test split	4
6.6 Standardize Features	5
6.7 Building analysis	5
7 Logistic Regression	9
7.1 Optimization	7
7.2 Best performance	5
7.3 Results Comparison	1
8 Descion Tree Classifier	7
8.1 Hyperparameters	3
8.1.1 Select Best Model	66
8.2 Results Comparison	9
9 Random Forest Classifier	9
9.1 Hyperparameters	5
9.2 Best Parameters	63
9.3 Results & Comparison	4
10 KNN Classifier	2
10.1 Tuning & Optimal K	3
10.1.1 Select Error Metric	1
10.1.1.1 K=3	5
10.1.2 Iterate over K	61
10.1.2.1 K=1	9
10.2 Results & Comparison	3
11 XGBoost	1
11.1 Results & Comparison	7
12 Gaussian NB	4
12.1 Results & Comparison	3
13 AdaBoostClassifier	5
13.1 Results & Comparison	4
14 Gradient Boosting	3
14.1 Results & Comparison	2
15 SVM Classifier	1
15.1 Tuning & Optimal C	3
15.1.1 Hyperparameter Tuning	3
15.1.1.1 Best Model	2

```
Name: state, dtype: int64
0    2411
1     922
Name: voice mail plan, dtype: int64
```

Contents ↴ 6.4 One Hot Encoder Categorical Columns

1 Abstract
 2 Business Problem
 3 Import Libraries
 4 Load Data
 5 Data Exploration
 ▶ 6 Data Engineering
 ▶ 6.1 Group state column
 6.1.1 Comment
 6.2 Drop unnecessary columns
 6.3 Drop highly correlated columns
 6.4 One Hot Encoding
 6.5 Train test split
 6.6 Standardize Features
 6.7 Building analysis
 ▶ 7 Logistic Regression
 7.1 Optimization
 7.2 Best performance
 7.3 Results Comparison
 ▶ 8 Decision Tree Classifier
 ▶ 8.1 Hyperparameters
 8.1.1 Select Best Model
 8.2 Results Comparison
 ▶ 9 Random Forest Classifier
 9.1 Hyperparameters
 9.2 Best Parameters
 9.3 Results & Comparison
 ▶ 10 KNN Classifier
 ▶ 10.1 Tuning & Optimization
 ▶ 10.1.1 Select Best Model
 In [138]:
 ▶ 10.1.1 K=3
 ▶ 10.1.2 Iterate over K values
 Out[140]:
 10.2 Results & Comparison
 ▶ 11 XGBoost
 11.1 Results & Comparison
 ▶ 12 Gaussian NB
 12.1 Results & Comparison
 ▶ 13 AdaBoostClassifier
 13.1 Results & Comparison
 ▶ 14 Gradient Boosting
 14.1 Results & Comparison
 ▶ 15 SVM Classifier
 ▶ 15.1 Tuning & Optimization
 ▶ 15.1.1 Hyperparameter Tuning
 ▶ 15.1.1 Best Model

```
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
```

In [140]:
 ▶ 10.1.1 K=3
 ▶ 10.1.2 Iterate over K values
 Out[140]:
 10.2 Results & Comparison
 ▶ 11 XGBoost
 11.1 Results & Comparison
 ▶ 12 Gaussian NB
 12.1 Results & Comparison
 ▶ 13 AdaBoostClassifier
 13.1 Results & Comparison
 ▶ 14 Gradient Boosting
 14.1 Results & Comparison
 ▶ 15 SVM Classifier
 ▶ 15.1 Tuning & Optimization
 ▶ 15.1.1 Hyperparameter Tuning
 ▶ 15.1.1 Best Model

	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	0
1	415	26	123	27.47	103	16.62	103	11.45	3	0
2	415	0	114	41.38	110	10.30	104	7.32	5	0
3	408	0	71	50.90	88	5.26	89	8.86	7	0
4	415	0	113	28.34	122	12.61	121	8.41	3	0

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', 'total intl charge', 'customer service calls', 'churn'], dtype='object')

Contents ⚙️⚙️

1 Abstract																																																												
2 Business Problem																																																												
In [141]:	1 df = pd.concat([pd.DataFrame(df_cont), df_cat], axis=1)																																																											
3 Import Libraries																																																												
4 Load Data																																																												
In [106]:	1 df.head()																																																											
6 Data Engineering																																																												
Out[196]:	<table border="1"> <thead> <tr> <th>account length</th><th>area code</th><th>number vmail messages</th><th>total day minutes</th><th>total day calls</th><th>total day charge</th><th>total eve minutes</th><th>total eve calls</th><th>total eve charge</th></tr> </thead> <tbody> <tr> <td>0</td><td>128</td><td>415</td><td>25</td><td>265.1</td><td>110</td><td>45.07</td><td>197.4</td><td>99</td><td>16</td></tr> <tr> <td>1</td><td>107</td><td>415</td><td>26</td><td>161.6</td><td>123</td><td>27.47</td><td>195.5</td><td>103</td><td>16</td></tr> <tr> <td>2</td><td>137</td><td>415</td><td>0</td><td>243.4</td><td>114</td><td>41.38</td><td>121.2</td><td>110</td><td>10</td></tr> <tr> <td>3</td><td>84</td><td>408</td><td>0</td><td>299.4</td><td>71</td><td>50.90</td><td>61.9</td><td>88</td><td>5</td></tr> <tr> <td>4</td><td>75</td><td>415</td><td>0</td><td>166.7</td><td>113</td><td>28.34</td><td>148.3</td><td>122</td><td>12</td></tr> </tbody> </table>	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
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																																																			
7 Logistic Regression	rows × 72 columns																																																											
8 Descion Tree Classifier																																																												
8.1 Hyperparameters																																																												
8.1.1 Select Best Model																																																												
In [82]:	1 df.columns																																																											
9 Random Forest Classifier																																																												
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')																																																											
10 KNN Classifier																																																												
10.1 Tuning & Optimize																																																												
10.1.1 Select Ensemble																																																												
10.1.1.1 K=3																																																												
10.1.2 Iterate cross-validation																																																												
10.1.2.1 K=1																																																												
10.2 Results & Conclusion																																																												
In [198]:	1 df.dtypes																																																											
11 XGBoost																																																												
11.1 Results & Conclusion																																																												
12 Gaussian NB																																																												
12.1 Results & Conclusion																																																												
13 AdaBoostClassifier																																																												
13.1 Results & Conclusion																																																												
14 Gradient Boosting																																																												
14.1 Results & Conclusion																																																												
15 SVM Classifier																																																												
15.1 Tuning & Optimization																																																												
15.1.1 Hyperparameter Tuning																																																												
15.1.1.1 Best Model																																																												
	length: 72, dtype: object																																																											

In [50]: 1 df.shape
2

Out[50]: (3333, 16)

In [51]: 1 df.shape

Contents Out[51]: (3333, 16)

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration

6.1 Train test split

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

```
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

6.2 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)
```

6.3 Building analysis functions

- 1 KNN Classification
- 2 ▼ 10.1 Tuning & Optimization
 - 3 ▼ 10.1.1 Select Estimator
 - 4 10.1.1.1 K=3
 - 5 ▼ 10.1.2 Iterate over different K values
 - 6 10.1.2.1 K=1
 - 7 10.2 Results & Conclusion
 - 8 ▼ 11 XGBoost
 - 9 11.1 Results & Conclusion
 - 10 ▼ 12 Gaussian NB
 - 11 12.1 Results & Conclusion
 - 12 ▼ 13 AdaBoostClassifier
 - 13 13.1 Results & Conclusion
 - 14 ▼ 14 Gradient Boosting
 - 15 14.1 Results & Conclusion
 - 16 ▼ 15 SVM Classifier
 - 17 15.1 Tuning & Optimization
 - 18 ▼ 15.1.1 Hyperparameter Tuning
 - 19 15.1.1.1 Best Model Selection

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)))

```

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='Train 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 ↗

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Euclidean Distance Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

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

```

```
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 ↗ ↘

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Best Model
 - 10.1.1.1 K=3
 - 10.1.2 Iterate & Select Best Model
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

Logistic Regression

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    model_log = logreg.fit(scaled_X_train, y_train)
12    print(model_log)
13
14    # Predict
15    y_hat_test = logreg.predict(scaled_X_test)
16
17    y_score = logreg.fit(scaled_X_train, y_train).decision_function(sca
18
19    fpr, tpr, thresholds = roc_curve(y_test, y_score)
20
21    print('AUC for {}: {}'.format(names[n], auc(fpr, tpr)))
22    print('-----')
23    lw = 2
24    plt.plot(fpr, tpr, color=colors[n],
25              lw=lw, label='ROC curve {}'.format(names[n]))
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)
JC for None: 0.81423403870501

logisticRegression(C=1e+20, class_weight='balanced', fit_intercept=False)
JC for Balanced: 0.814794718755652

logisticRegression(C=1e+20, class_weight={0: 1, 1: 2}, fit_intercept=False)
JC for 2 to 1: 0.8150479291011034

logisticRegression(C=1e+20, class_weight={0: 1, 1: 10}, fit_intercept=False)
JC for 10 to 1: 0.8137457044673538

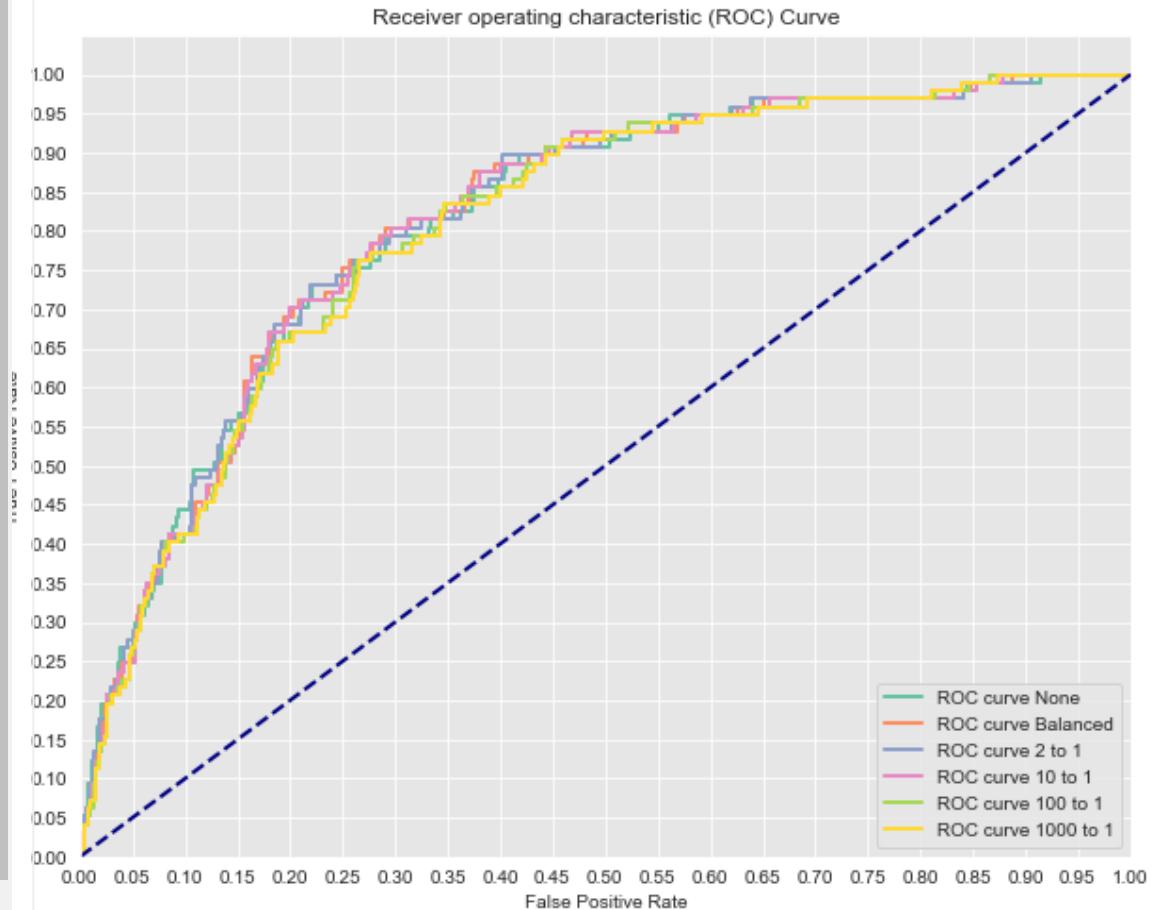
15.1.1 Best

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

```
LogisticRegression(C=1e+20, class_weight={0: 1, 1: 1000}, fit_intercept=False)
AUC for 1000 to 1: 0.8059504431181045
```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Euclidean Distance
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



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)
```

Contents ⚙️⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Common features
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

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(scaled_X_test)
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'C={}'.format(c)))
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')
JC for 0.001: 0.8137999638270935

logisticRegression(C=0.01, fit_intercept=False, solver='liblinear')
JC for 0.01: 0.8139989148128052

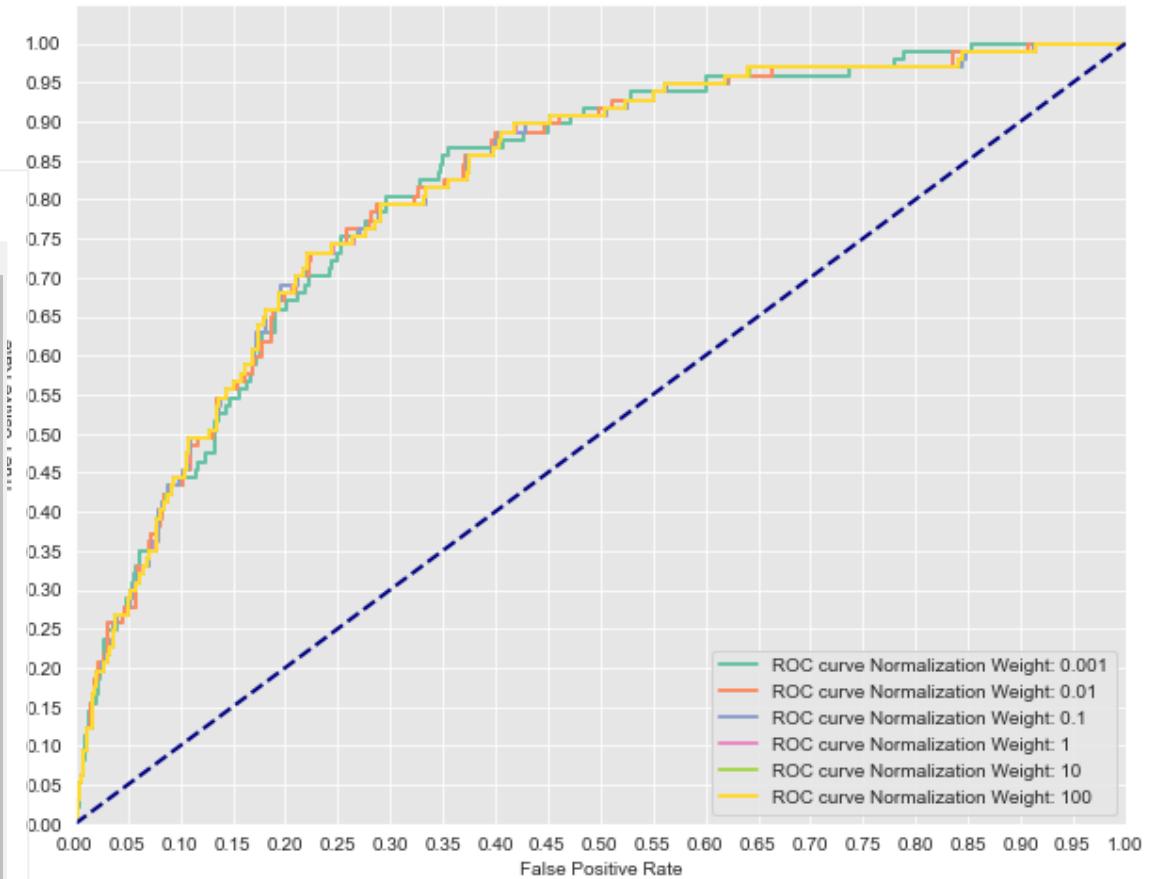
logisticRegression(C=0.1, fit_intercept=False, solver='liblinear')
JC for 0.1: 0.8140893470790378

logisticRegression(C=1, fit_intercept=False, solver='liblinear')
JC for 1: 0.8142159522517635

logisticRegression(C=10, fit_intercept=False, solver='liblinear')
JC for 10: 0.8142340387050099

logisticRegression(C=100, fit_intercept=False, solver='liblinear')
JC for 100: 0.81423403870501

Receiver operating characteristic (ROC) Curve



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performing parameters
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Best Model
 - 10.1.1.1 K=3
 - 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

Best performing parameters for Logistic Regression Model

```
In [64]: N
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 ⚙️

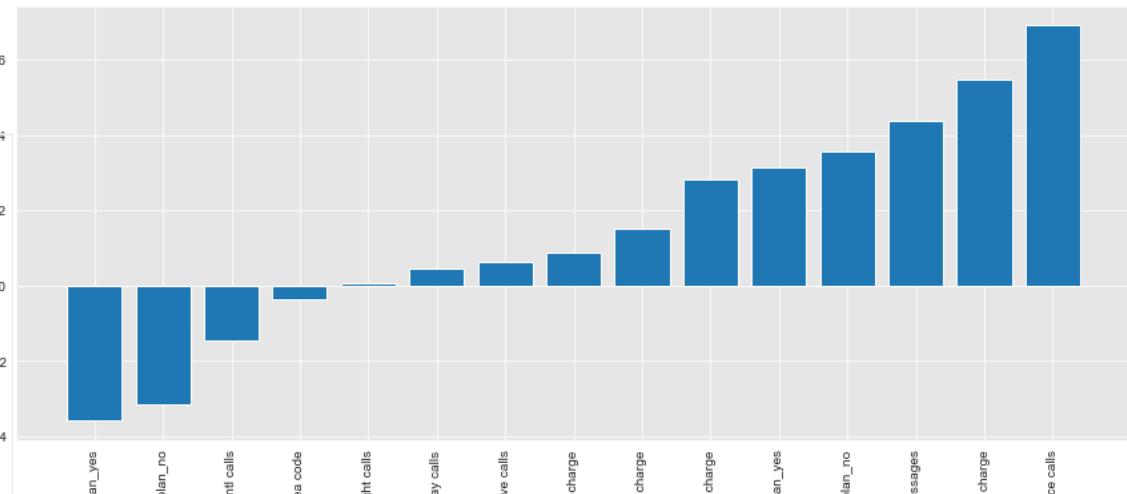
- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state & gender
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classification
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classification
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

Classification Report for train & test set

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

```
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]])
```

In [66]: 1 plot_coefficients(clf)



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building a Model
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance set.
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 Minimizing error function
 - 10.1.2 Iterate over different k values
 - 10.1.2.1 Reducing error function
 - 10.2 Results & Comparison
- 11 XGBoost Classifier
 - 11.1 Results & Comparison
- 12 Gaussian NB Classifier
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameter Tuning
 - 15.1.1.1 Best Parameters

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. It means 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.

Important positive churn factors:

Customer service calls

I day charge

Number voicemail messages

e mail plan _no

national plan_yes

I evening charge

10.1.2.1 Reducing error function

Important negative churn features:

e mail plan _yes

national plan_no

I International call

area code

8 Decision Tree Classifier

In [67]:

```

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.

Hyperparameter Tuning for decision tree using GridSearch

```

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} different permutations.")

```

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_error=False)
4
5 # Fit to the data
6 dt_grid_search.fit(scaled_X_train, y_train)
7
8 #Examine best parameters
9
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}
```

Select Best Parameters for Descion tree classifier

Contents

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - Out[70]:**
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over error function
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

```
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)

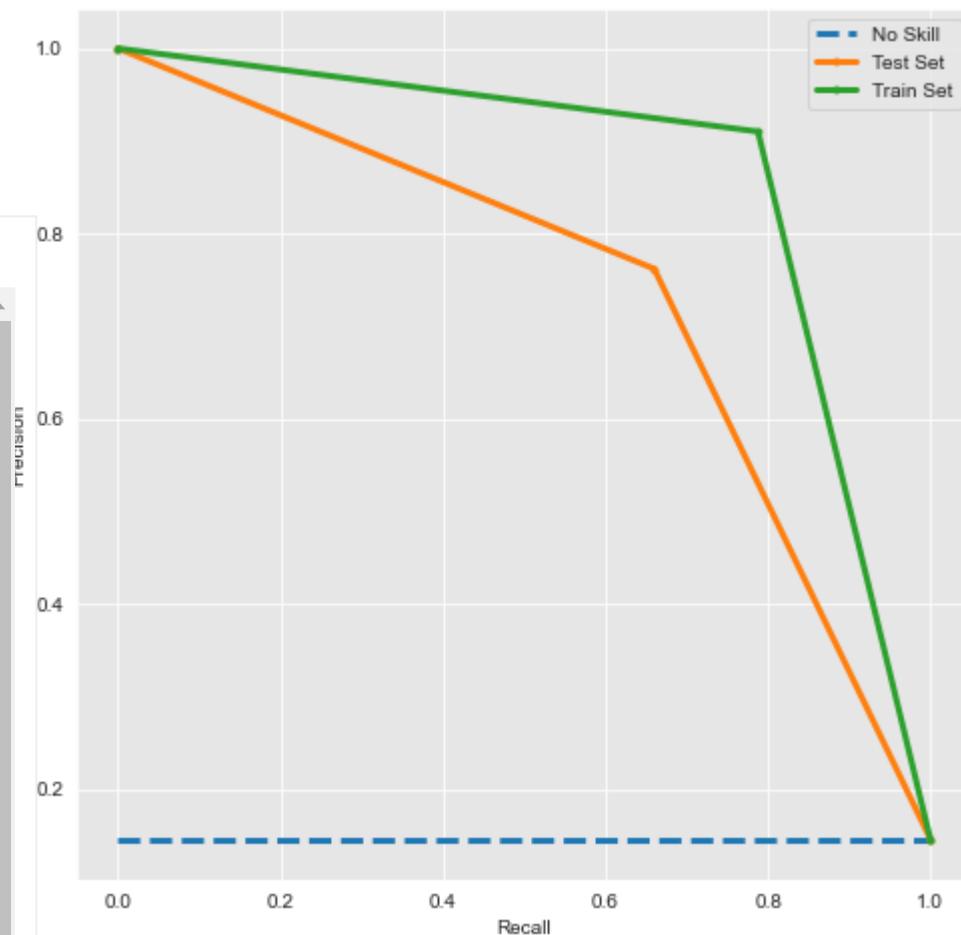
mean Training Score: 95.79%
mean Test Score: 92.05%
precision    recall   f1-score   support
False        0.94      0.96      0.95      570
True         0.76      0.66      0.71       97
accuracy          0.92      0.83      0.92      667
macro avg      0.85      0.81      0.83      667
weighted avg    0.92      0.92      0.92      667

MODEL EVALUATION METRICS:
15.1.1 Best
-----confusion Matrix for train & 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

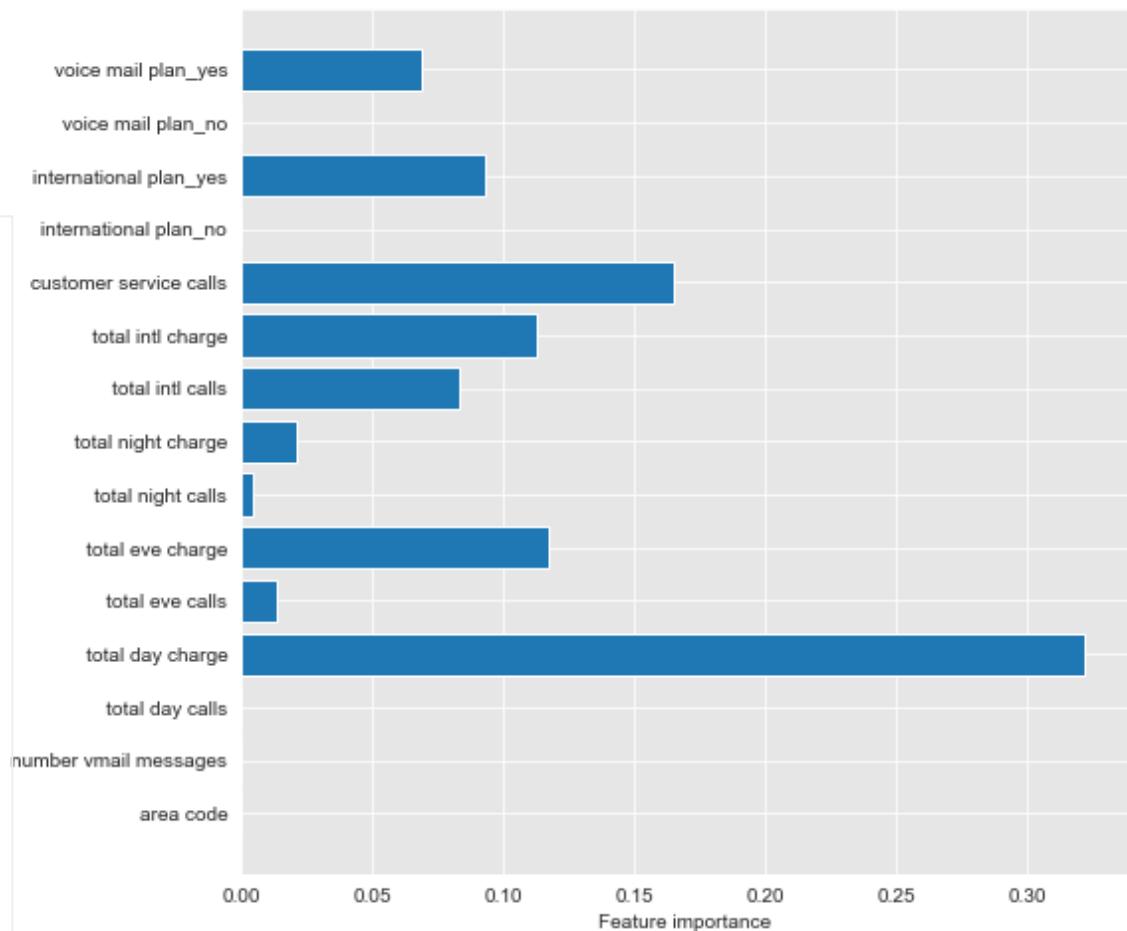
Contents ⚙️

1 Abstract	Predicted	No Churn	Churn	All	
2 Business Problem	Actual	NaN	NaN	NaN	
3 Import Libraries	No Churn	550.0	20.0	570.0	
4 Load Data	Churn	33.0	64.0	97.0	
5 Data Exploration	All	583.0	84.0	667.0	
6 Data Engineering					
6.1 Group state of account				Classification Report for train & test set	
6.1.1 Comments					
6.2 Drop unnecessary columns					
6.3 Drop highly correlated					
6.4 One Hot Encoding					
6.5 Train test split	precision	recall	f1-score	support	
6.6 Standardize Features	False	0.96	0.99	0.98	2280
6.7 Building analysis	True	0.91	0.79	0.84	386
7 Logistic Regression	accuracy			0.96	2666
7.1 Optimization	macro avg	0.94	0.89	0.91	2666
7.2 Best performance	weighted avg	0.96	0.96	0.96	2666
8 Decision Tree Classifier					
8.1 Hyperparameters				Classification Report for test set	
8.1.1 Select Best Model	precision	recall	f1-score	support	
8.2 Results Comparison	False	0.94	0.96	0.95	570
9 Random Forest Classifier	True	0.76	0.66	0.71	97
9.1 Hyperparameters	accuracy			0.92	667
9.2 Best Parameters	macro avg	0.85	0.81	0.83	667
9.3 Results & Comparison	weighted avg	0.92	0.92	0.92	667
10 KNN Classifier					
10.1 Tuning & Optimal Model				Kappa's Kappa for train and test set:	
10.1.1 Select Euclidean Distance	0.8203	0.6615			
10.1.2 Iterate cross validation	0.8203	0.6615			
10.1.2.1 K=1	0.8203	0.6615			
10.2 Results & Comparison	0.8094	0.678			
11 XGBoost	roc_auc score for train and test set:				
11.1 Results & Comparison	0.8872	0.8124			
12 Gaussian NB	mean Cross Validation Score:				
12.1 Results & Comparison	0.9373				
13 AdaBoostClassifier					
13.1 Results & Comparison					
14 Gradient Boosting					
14.1 Results & Comparison					
15 SVM Classifier					
15.1 Tuning & Optimal Model				Hypertuning for train and test set:	
15.1.1 Best Model	0.8872	0.8124			



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Parameters
 - ▼ 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Parameters
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimizing
 - ▼ 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimizing
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

In [72]:

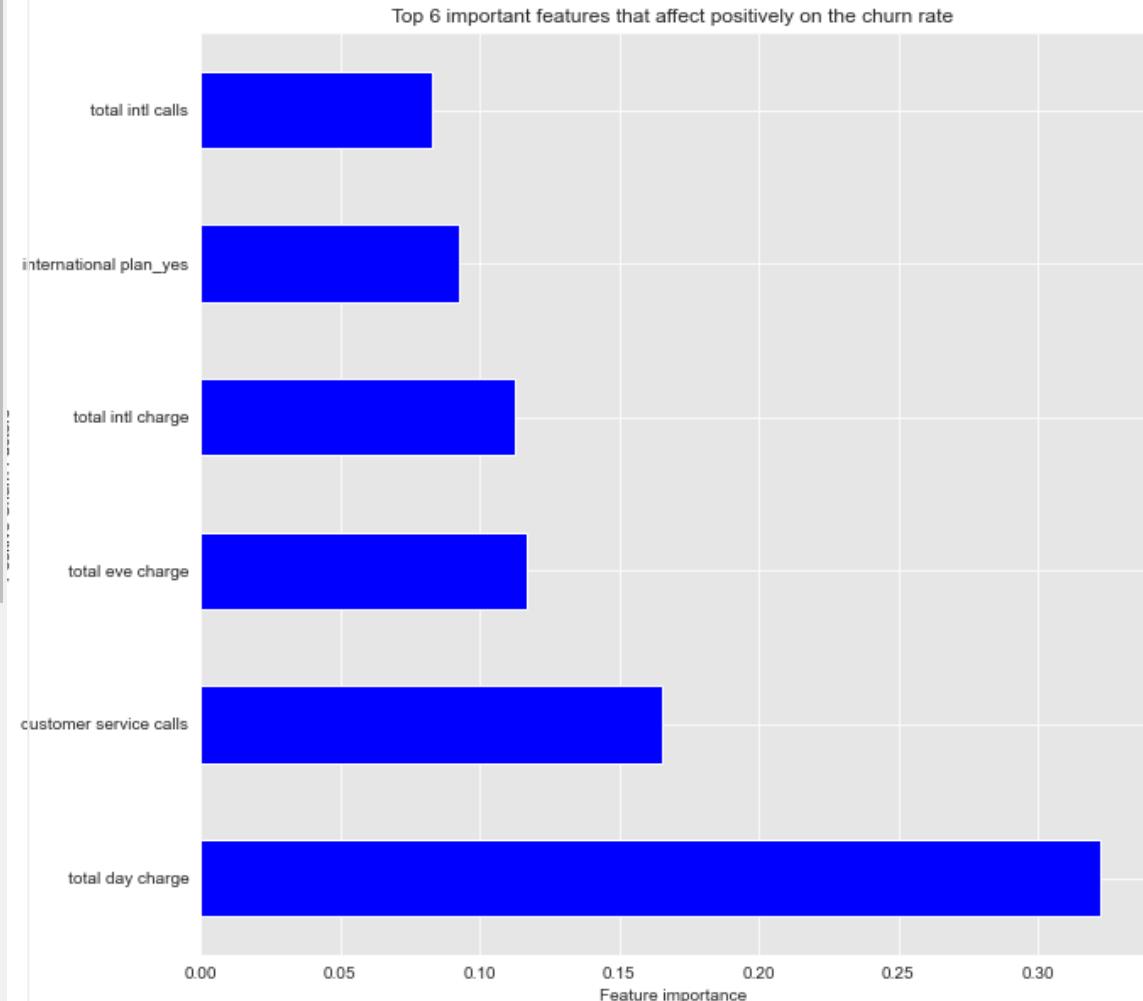
```

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 ⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state by city
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
- 7 Logistic Regression
- 8 Descion Tree Classifer
- 9 Random Forest Classifer
- 10 KNN Classifier
- 11 XGBoost
- 12 Gaussian NB
- 13 AdaBoostClassifer
- 14 Gradient Boosting
- 15 SVM Classifer

8.1 Results Comments

We have a precision of 0.87 and recall of 0.65 for the test set. It means we were able to detect the churn customers 65% but 87% of the customers we say 'churn' was true.

Top 5 factors affecting te churn rate:

1) Total day charge 2) Customer service calls 3) Total evening charge 4) Total international charge
5) International plan_yes 4) Total international calls

9 Random Forest Classifer

```

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%}")

```

Mean Cross Validation Score: 93.96%

```
In [77]: 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

Contents ⚙️

1 Abstract	
2 Business Problem	
3 Import Libraries	
4 Load Data	
5 Data Exploration	
6 Data Engineering	
6.1 Group state	
6.1.1 Comment	
6.2 Drop unnece	
6.3 Drop higly co	
6.4 One Hot Enc	
6.5 Train test sp	
6.6 Standarize F	
6.7 Building anal	
7 Logistic Regress	
7.1 Optimization	
7.2 Best perform	
7.3 Results Com	
8 Descion Tree Cl	
8.1 Hyperparamet	
8.1.1 Select Be	
8.2 Results Com	
9 Random Forest C	
9.1 Hyperparamet	
9.2 Best Paramet	
9.3 Results & Co	
10 KNN Classifier	
10.1 Tuning & Opt	
10.1.1 Select E	
10.1.1.1 K=3	
10.1.2 Iterate c	
10.1.2.1 K=1	
10.2 Results & C	
11 XGBoost	
11.1 Results & co	
12 Gaussian NB	
12.1 Results & C	
13 AdaBoostClassi	
13.1 Results & C	
14 Gradient Boosti	
14.1 Results & co	
15 SVM Classifier	
15.1 Tuning & Opt	
15.1.1 Hypertun	
15.1.1.1 Best	

Hyperparameter Tuning for Random Forest using GridSearch

```
In [78]: 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_traini
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_train_s
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_
ean Training Score: 91.03%
ean Test Score: 93.85%
est Parameter Combination Found During Grid Search:
'criterion': 'entropy',
'max_depth': None,
'min_samples_leaf': 1,
'min_samples_split': 5}
```

9.2 Best Parameters for Random Forest

```
In [79]: dt_param_grid2 = {
    'criterion': ['entropy'],
    'max_depth': [None],
    'min_samples_split': [5],
    'min_samples_leaf': [1]
}
dt_grid_search2 = GridSearchCV(rf_clf, dt_param_grid2, cv=10, return_train_time=True)
# Fit to the data
clf=dt_grid_search2.fit(scaled_X_train, y_train)

#Examine best parameters
# Mean training score
dt_gs_training_score2 = np.mean(clf.cv_results_['mean_train_score'])

# Mean test score
dt_gs_testing_score2 = clf.score(scaled_X_test, y_test)

print(f"Mean Training Score: {dt_gs_training_score2 :.2%}")
print(f"Mean Test Score: {dt_gs_testing_score2 :.2%}")

# Model Performance
# Test set predictions
y_pred_test =clf.predict(scaled_X_test)
y_pred_train =clf.predict(scaled_X_train)
```

Mean Training Score: 98.82%
 Mean Test Score: 93.25%

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Error Metrics
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best

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

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Euclidean Distance
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters
 - accuracy

	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 & test set:**

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	2280.0	0.0	2280.0
Churn	31.0	355.0	386.0
All	2311.0	355.0	2666.0

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	563.0	7.0	570.0
Churn	38.0	59.0	97.0
All	601.0	66.0	667.0

Classification Report for train & test set

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

test set	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

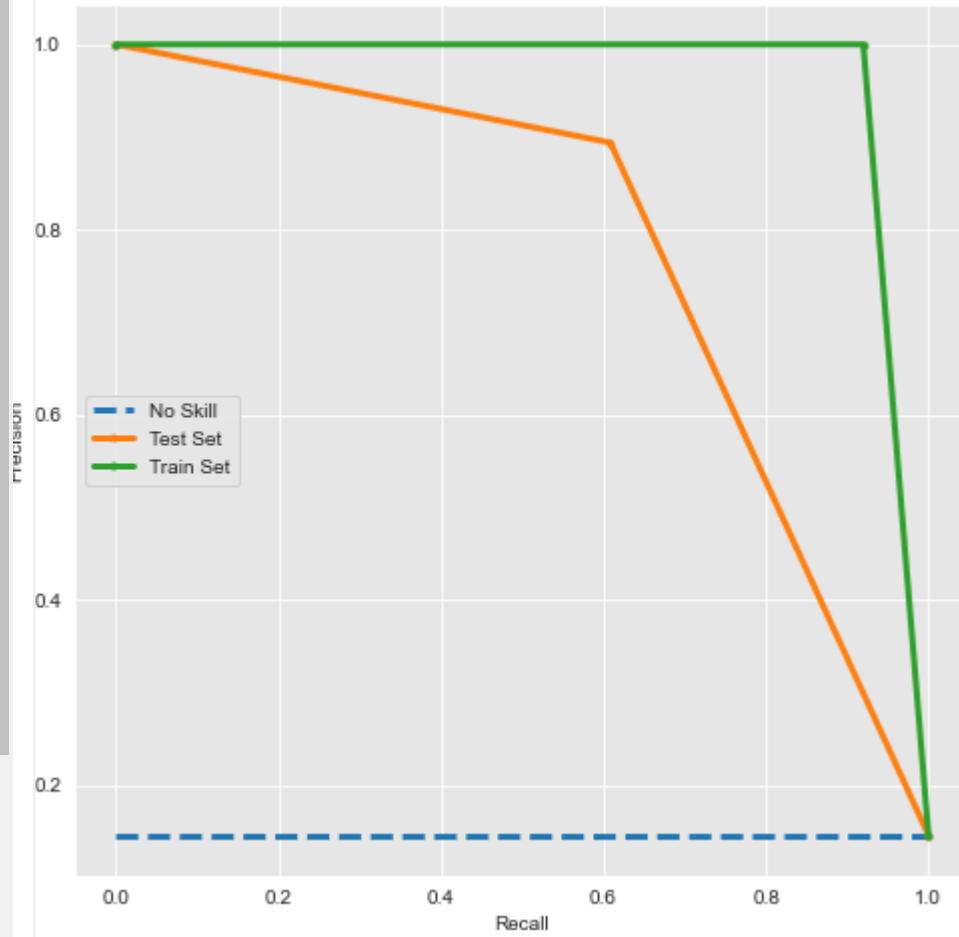
macro avg	0.92	0.80	0.84	667
weighted avg	0.93	0.93	0.93	667

Cohen's Kappa for train and test set:

0.9514 0.6871

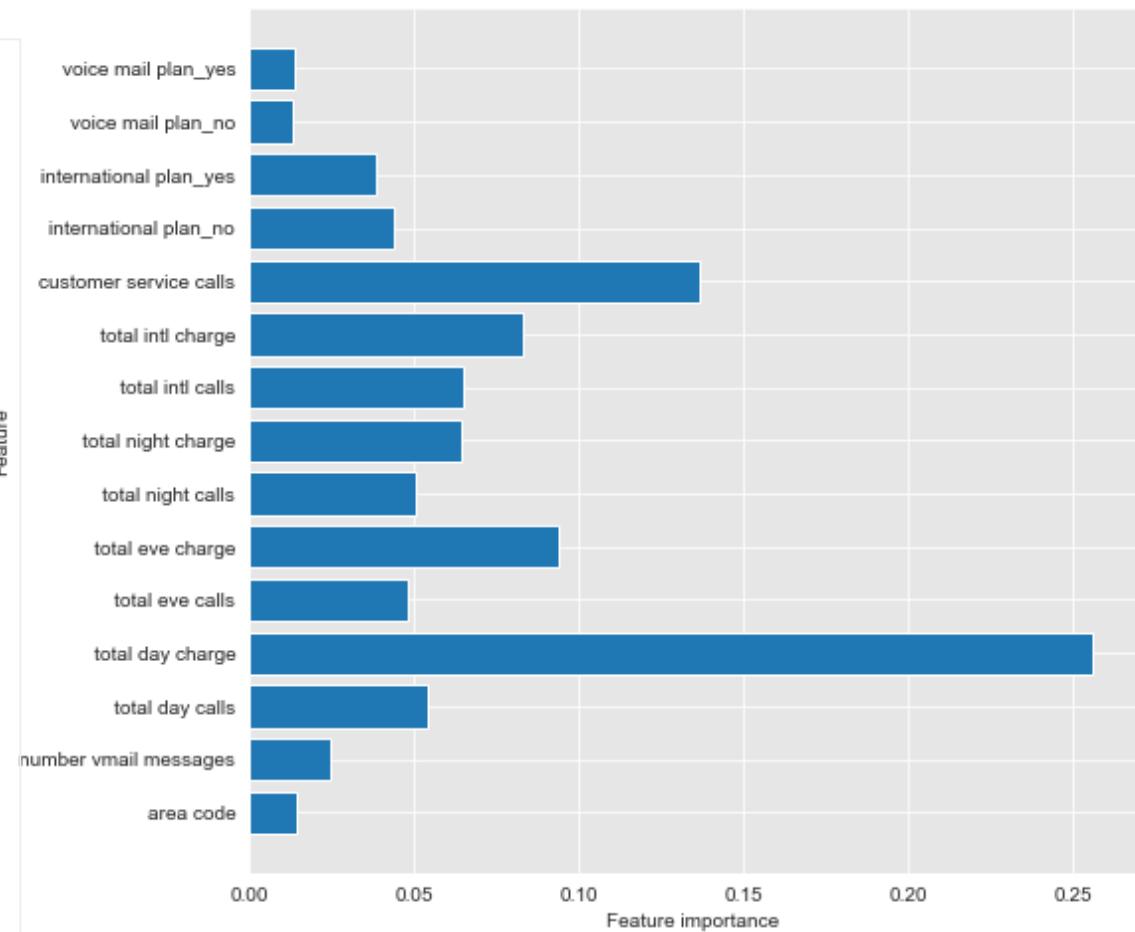
Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Parameters
 - ▼ 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Parameters
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Parameters
 - ▼ 10.1.1 Select Error Metrics
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
 - ▼ 11 XGBoost
 - 11.1 Results & comparison
 - ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
 - ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
 - ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
 - ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Parameters
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



In [81]:

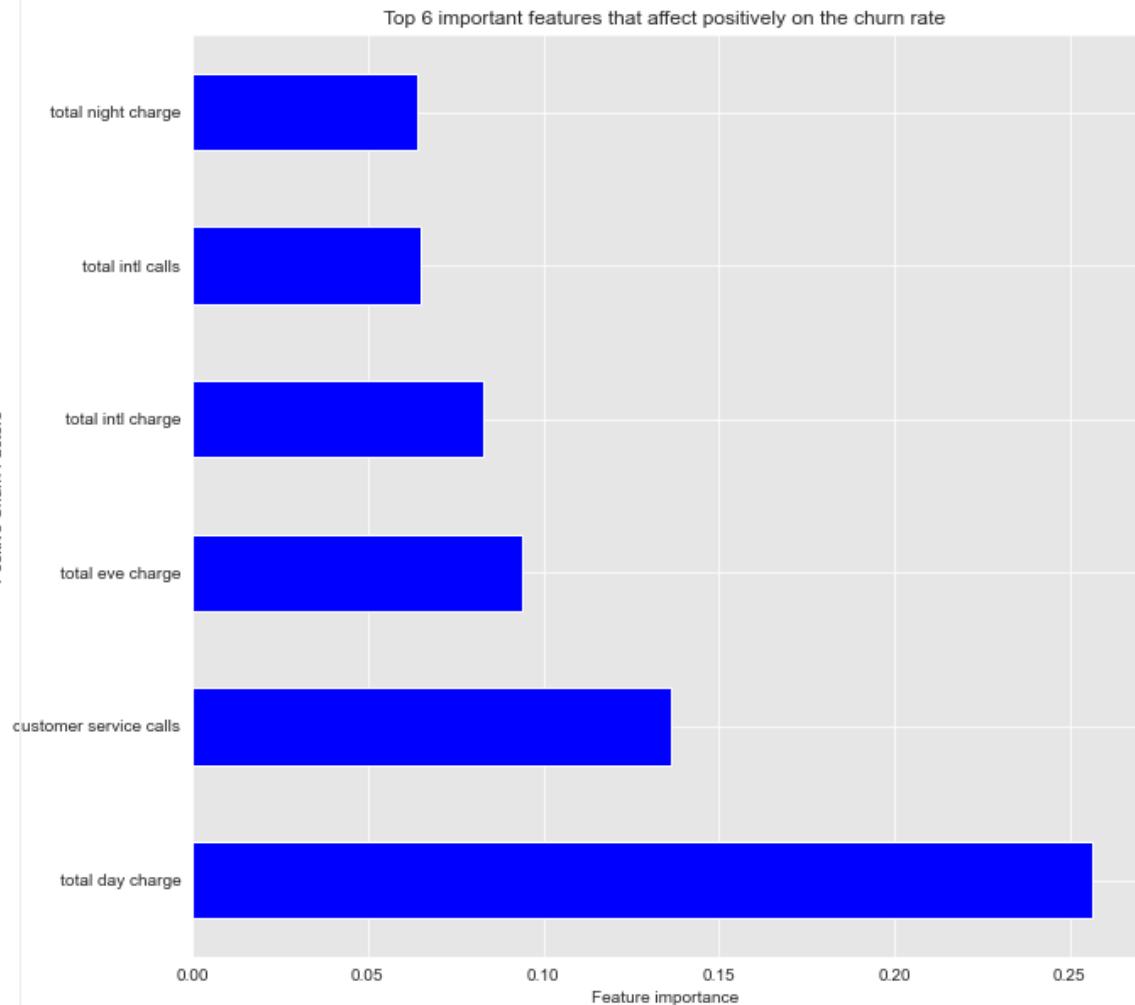
```

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

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - Out[81]:**
 - 6.1 Group state column
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 7.1 Optimization
- 7.2 Best performance
- 7.3 Results Comparison
- 8 Decision Tree Classifier
 - Out[81]:**
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



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.

The top 5 factors affecting te churn rate:

Contents

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
- 7 Logistic Regress
- 8 Descion Tree Clas
- 9 Random Forest C
- 10 KNN Classifier
- 11 XGBoost
- 12 Gaussian NB
- 13 AdaBoostClassi
- 14 Gradient Boosti
- 15 SVM Classifier

1) Total day charge 2) Customer service calls 3) Total evening charge 4) Total international charge
5) Total international calls 6) Total night charge

10 KNN Classifier

- 6.1 Group state c
- 6.1.1 Commen
- 6.2 Drop unnece
- 6.3 Drop higly co
- 6.4 One Hot Enc
- 6.5 Train test spl
- 6.6 Standarize F
- 6.7 Building anal
- 7.1 Optimization
- 7.2 Best perform
- 7.3 Results Com
- 8.1 Hyperparamet
- 8.1.1 Select Be
- 8.2 Results Com
- 9.1 Hyperparamet
- 9.2 Best Paramet
- 9.3 Results & Co
- 10.1 Tuning & Opt
- 10.1.1 Select E
- 10.1.1.1 K=3
- 10.1.2 Iterate c
- 10.1.2.1 K=1
- 10.2 Results & C
- 11.1 Results & co
- 12.1 Results & C
- 13.1 Results & C
- 14.1 Results & co
- 15.1 Tuning & Opt
- 15.1.1 Hypertun
- 15.1.1.1 Best

In [82]: N ▾

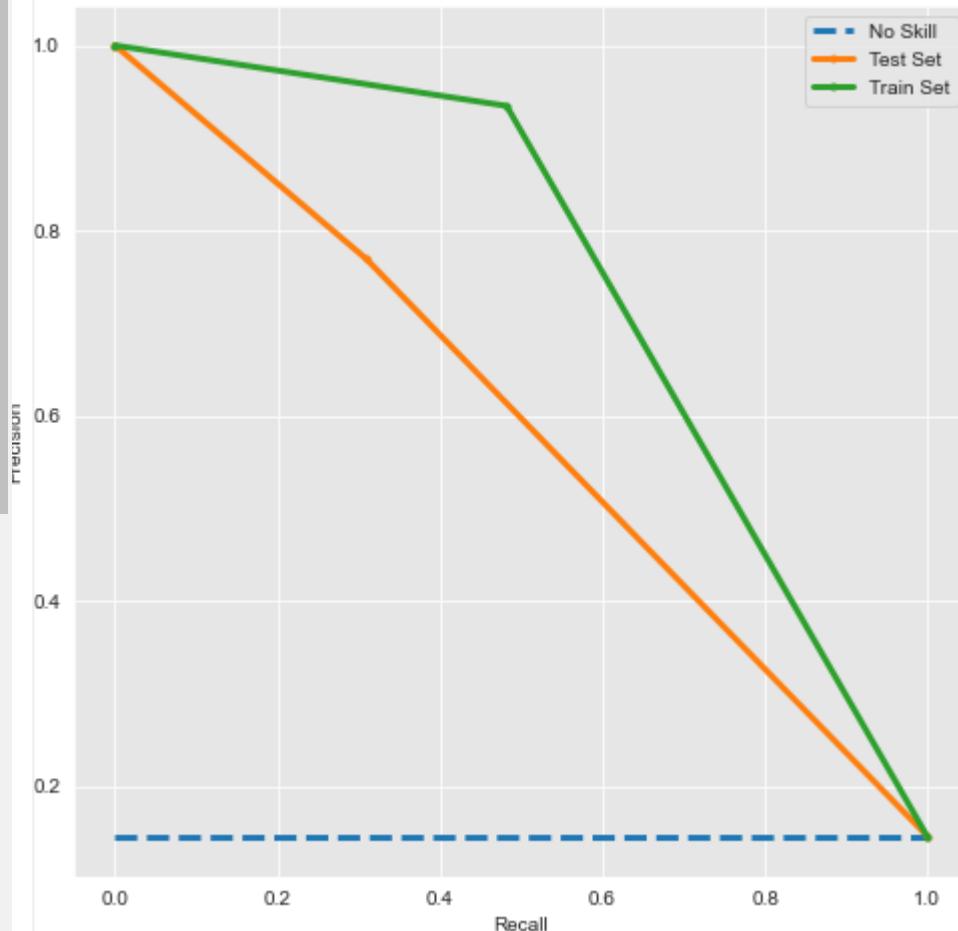
```

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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model



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:

Predicted	No Churn	Churn	All
Actual	Nan	Nan	Nan
No Churn	2267.0	13.0	2280.0
Churn	200.0	186.0	386.0
All	2467.0	199.0	2666.0

Predicted	No Churn	Churn	All
Actual	Nan	Nan	Nan
No Churn	561.0	9.0	570.0
Churn	67.0	30.0	97.0
All	628.0	39.0	667.0

Classification Report for train & test set

train set

	precision	recall	f1-score	support	
7.1 Optimization					
7.2 Best perform	False	0.92	0.99	0.96	2280
7.3 Results Com	True	0.93	0.48	0.64	386

accuracy

accuracy				
macro avg	0.93	0.74	0.80	2666

weighted avg

weighted avg	0.92	0.92	0.91	2666
--------------	------	------	------	------

Random Forest Classifier

test set

	precision	recall	f1-score	support	
9.1 Hyperparam					
9.2 Best Param	False	0.89	0.98	0.94	570
9.3 Results & Co	True	0.77	0.31	0.44	97

accuracy

accuracy				
macro avg	0.83	0.65	0.69	667

weighted avg

weighted avg	0.88	0.89	0.86	667
--------------	------	------	------	-----

Results & Conclusions

XGBoost

Results & Conclusions

Gaussian NB

Results & Conclusions

AdaBoostClassifier

Results & Conclusions

Gradient Boosting

Results & Conclusions

SVM Classifier

Tuning & Optimize

Hyperparameter

Best Model

Kappa's Kappa for train and test set:

0.5961 0.3903

AUC score for train and test set:

0.5336 0.3513

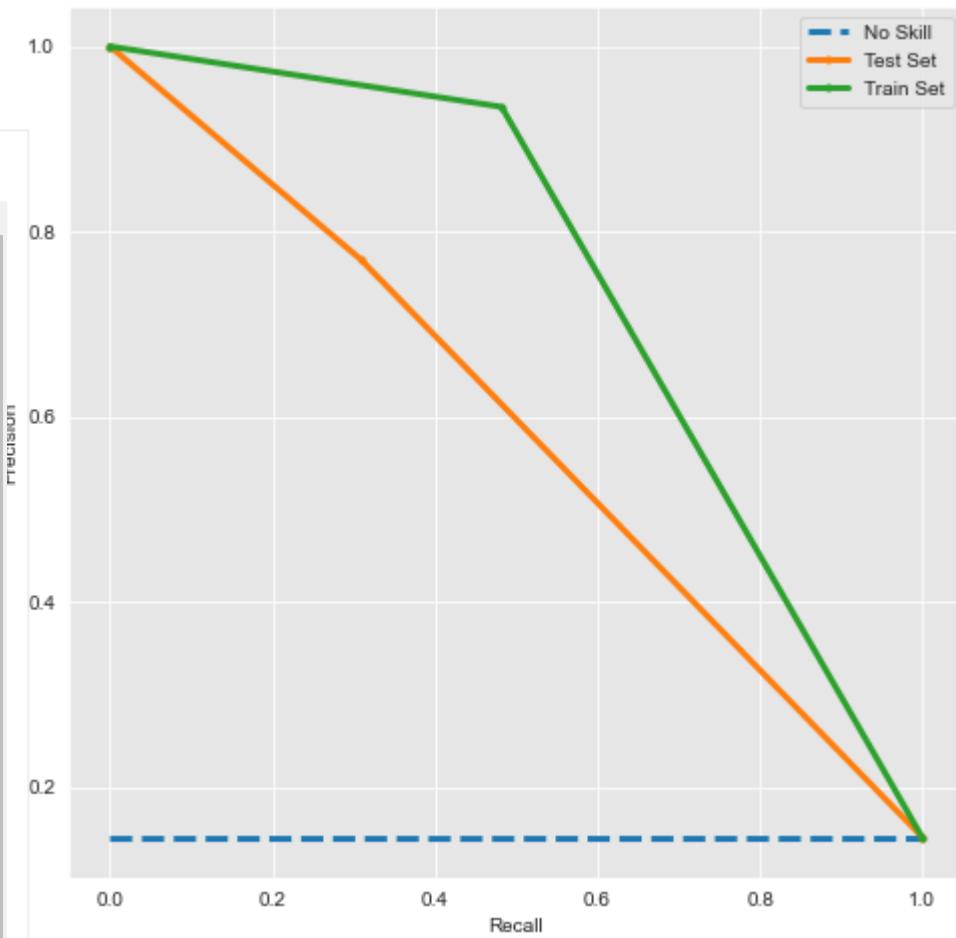
ROC auc score for train and test set:

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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Best Model
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over n values
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

Tuning & Optimization (Iterate over n values)

```

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))
15
16
17

```

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

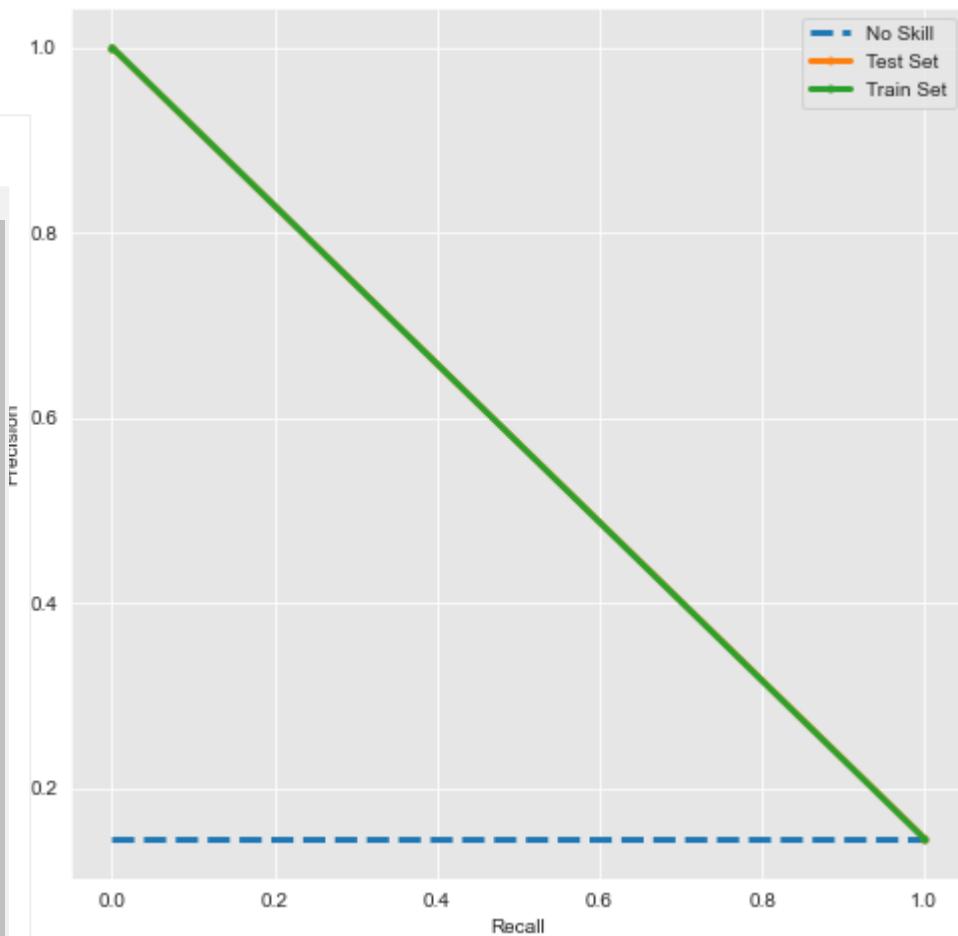
```
Best Value for k: 3
F1-Score: 0.4583333333333334
```

Contents ↻

10.1.1 Select Best K for KNN Machine Learning Model

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state by zip code
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best K
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal K
 - 10.1.1 Select Best K
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over different K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal C
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best C

Mean Cross Validation Score:
0.8437



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - In [88]:
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over n
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

```

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

```

Mean Cross Validation Score: 89.05%

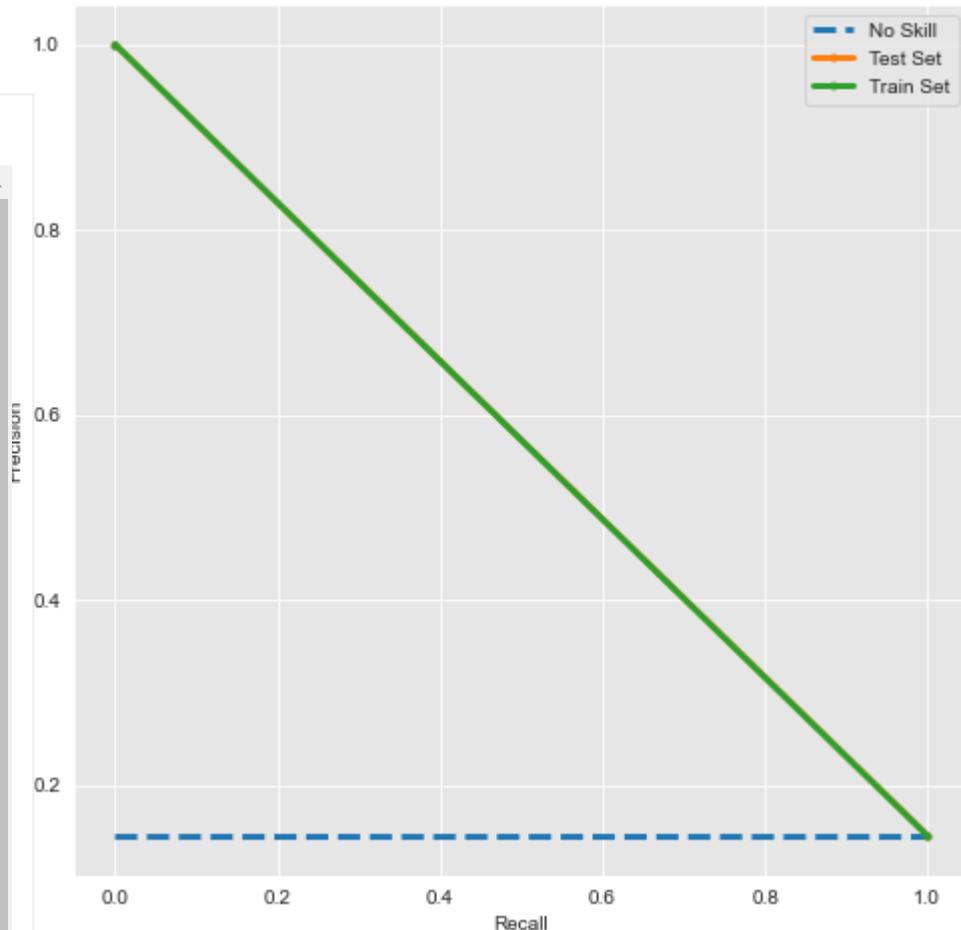
2 Iterate over $1 < n < 26$

```
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 ⚙️

1 Abstract
2 Business Problem
3 Import Libaries
4 Load Data
5 Data Exploration
6 Data Engineering
6.1 Group state by zip code
6.1.1 Comments
6.2 Drop unnecessary columns
6.3 Drop highly correlated columns
6.4 One Hot Encoding
6.5 Train test split
6.6 Standardize Features
6.7 Building analysis
7 Logistic Regression
7.1 Optimization
7.2 Best performance
7.3 Results Comparison
8 Descion Tree Classification
8.1 Hyperparameters
8.1.1 Select Best Model
8.2 Results Comparison
9 Random Forest Classification
9.1 Hyperparameters
9.2 Best Parameters
9.3 Results & Comparison
10 KNN Classifier
10.1 Tuning & Optimization
10.1.1 Select Error Metric
10.1.1.1 K=3
10.1.2 Iterate cross-validation
10.1.2.1 K=1
10.2 Results & Comparison
11 XGBoost Classification
11.1 Results & Comparison
12 Gaussian NB
12.1 Results & Comparison
13 AdaBoostClassifier
13.1 Results & Comparison
14 Gradient Boosting
14.1 Results & Comparison
15 SVM Classifier
15.1 Tuning & Optimization
15.1.1 Hyperparameters
15.1.1.1 Best Model

Mean Cross Validation Score:
0.7942



Contents ⚙️

1	Abstract
2	Business Problem
3	Import Libraries
4	Load Data
5	Data Exploration
6	Data Engineering
6.1	Group state columns
6.1.1	Comments
6.2	Drop unnecessary columns
6.3	Drop highly correlated columns
6.4	One Hot Encoding
6.5	Train test split
6.6	Standarize Features
6.7	Building analysis
7	Logistic Regression
7.1	Optimization
7.2	Best performance
7.3	Results Comparison
8	Decision Tree Classifier
8.1	Hyperparameters
8.1.1	Select Best Model
8.2	Results Comparison
9	Random Forest Classifier
9.1	Hyperparameters
9.2	Best Parameters
9.3	Results & Comparison
10	KNN Classifier
10.1	Tuning & Optimization
10.1.1	Select Best Model
10.1.1.1	K=3
10.1.2	Iterate on K
10.1.2.1	K=1
10.2	Results & Comparison
11	XGBoost
11.1	Results & Comparison
12	Gaussian NB
12.1	Results & Comparison
13	AdaBoostClassifier
13.1	Results & Comparison
14	Gradient Boosting
14.1	Results & Comparison
15	SVM Classifier
15.1	Tuning & Optimization
15.1.1	Hyperparameters
15.1.1.1	Best Model

10 Results & Comments

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

```
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 ↗

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state column
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

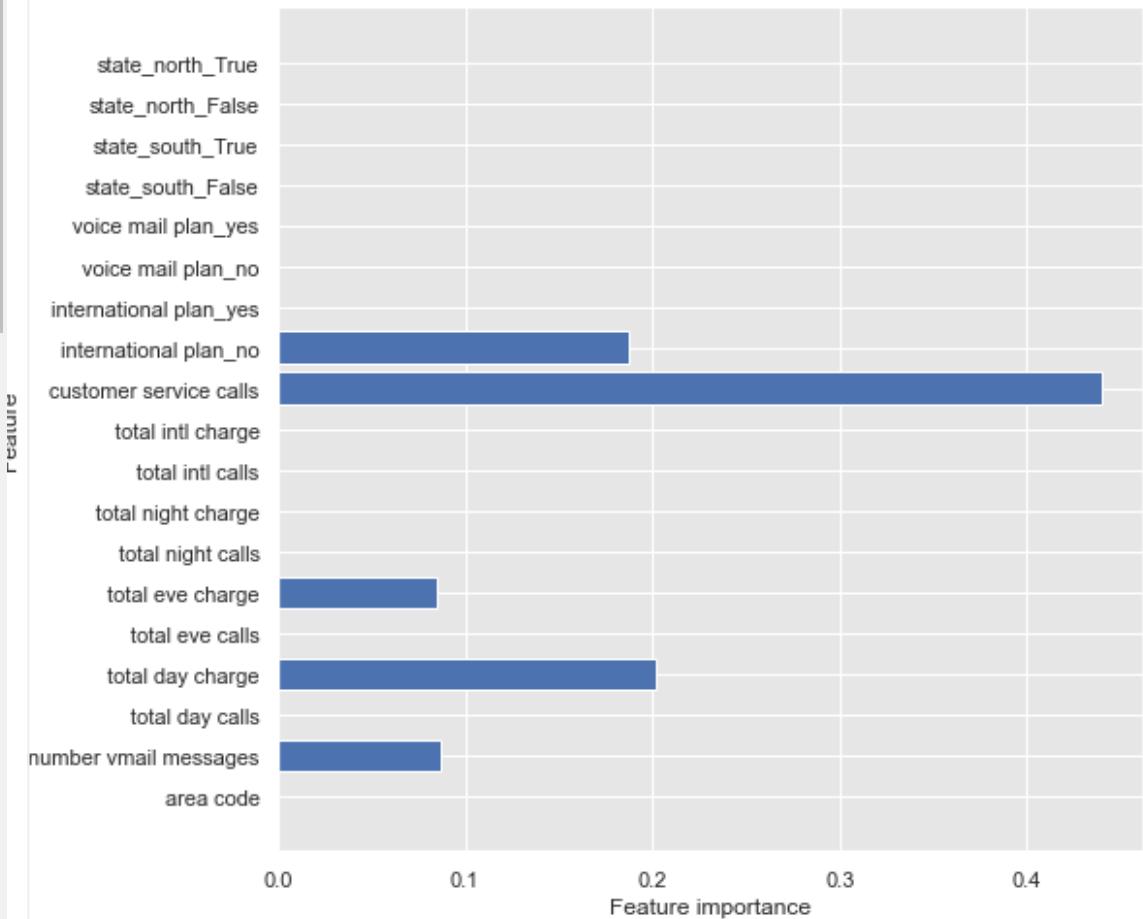
```

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)
```



In [401]:

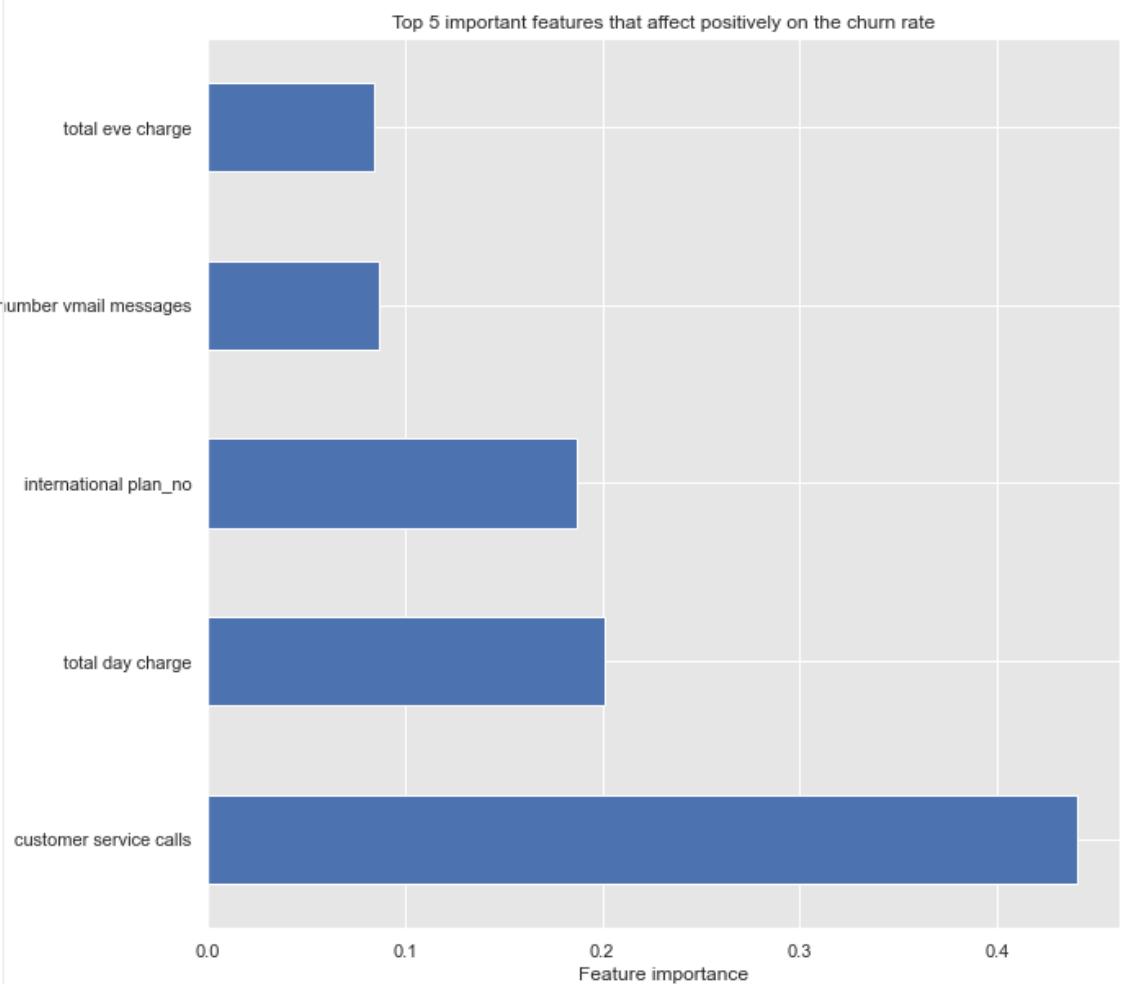
```

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 ↗

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of each customer
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

**11 Results & comments**

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.

The top 5 factors affecting te churn rate:

- 1) Customer service calls
- 2) Total day charge
- 3) International plan_no

Contents ⚙⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration

6 Data Engineering

- 6.1 Group state
- 6.1.1 Commencement
- 6.2 Drop unneeded
- 6.3 Drop highly co
- 6.4 One Hot Encod
- 6.5 Train test spli
- 6.6 Standardize F
- 6.7 Building anal

7 Logistic Regress

- 7.1 Optimization
- 7.2 Best perform
- 7.3 Results Com

8 Descion Tree Cl

- 8.1 Hyperparamet
- 8.1.1 Select Be
- 8.2 Results Com

9 Random Forest C

- 9.1 Hyperparamet
- 9.2 Best Paramet
- 9.3 Results & Co

10 KNN Classifier

- 10.1 Tuning & Opt
- 10.1.1 Select E
- 10.1.1.1 K=3
- 10.1.2 Iterate c
- 10.1.2.1 K=1
- 10.2 Results & Co

11 XGBoost

- 11.1 Results & co

12 Gaussian NB

- 12.1 Results & Co

13 AdaBoostClassi

- 13.1 Results & Co

14 Gradient Boosti

- 14.1 Results & co

15 SVM Classifier

- 15.1 Tuning & Opt
- 15.1.1 Hypertun

- 15.1.1.1 Best

12 Gaussian NB

```

1 nb = GaussianNB()
2 nb_model=nb.fit(scaled_X_train, y_train)
3 y_pred_train = nb_model.predict(scaled_X_train)
4 y_pred_test = nb_model.predict(scaled_X_test)
5

```

```
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

Confusion Matrix for train & test set:

Contents

- 1 Abstract
 - 2 Business Problem
 - 3 Import Libraries
 - 4 Load Data
 - 5 Data Exploration

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

- 6 Data Engineering
 - 6.1 Group state
 - 6.1.1 Comment
 - 6.2 Drop unnece
 - 6.3 Drop higly co
 - 6.4 One Hot Enc

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	570.0	0.0	570.0
Churn	0.0	97.0	97.0
All	570.0	97.0	667.0

Classification Report for train & test set

- ## ▼ 7 Logistic Regress

rain set

precision	recall	f1-score	support
1.00	1.00	1.00	2280
1.00	1.00	1.00	386
		1.00	2666
1.00	1.00	1.00	2666

- ▼ 9 Random Forest (1)
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Conclusion
 - 10 KNN Classification

est set

		precision	recall	f1-score	support
▼ 10 KNN Classifier					
▼ 10.1 Tuning & Ov					

- #### ▼ 10.1.1 Select E

- 1 -

10.1.1 K=3	True	1.00	1.00	1.00	97
▼ 10.1.2 Iterate c					
10.1.2.1 K=1	accuracy			1.00	667
10.2 Results & C	macro avg	1.00	1.00	1.00	667
▼ 11 XGBoost	weighted avg	1.00	1.00	1.00	667

- ## 11.1 Results & co

ohen's Kappa for train and test set:

1.0 1.0

2 score for train and test set:

1010

roc_auc score for train and test set:

1 0 1 0

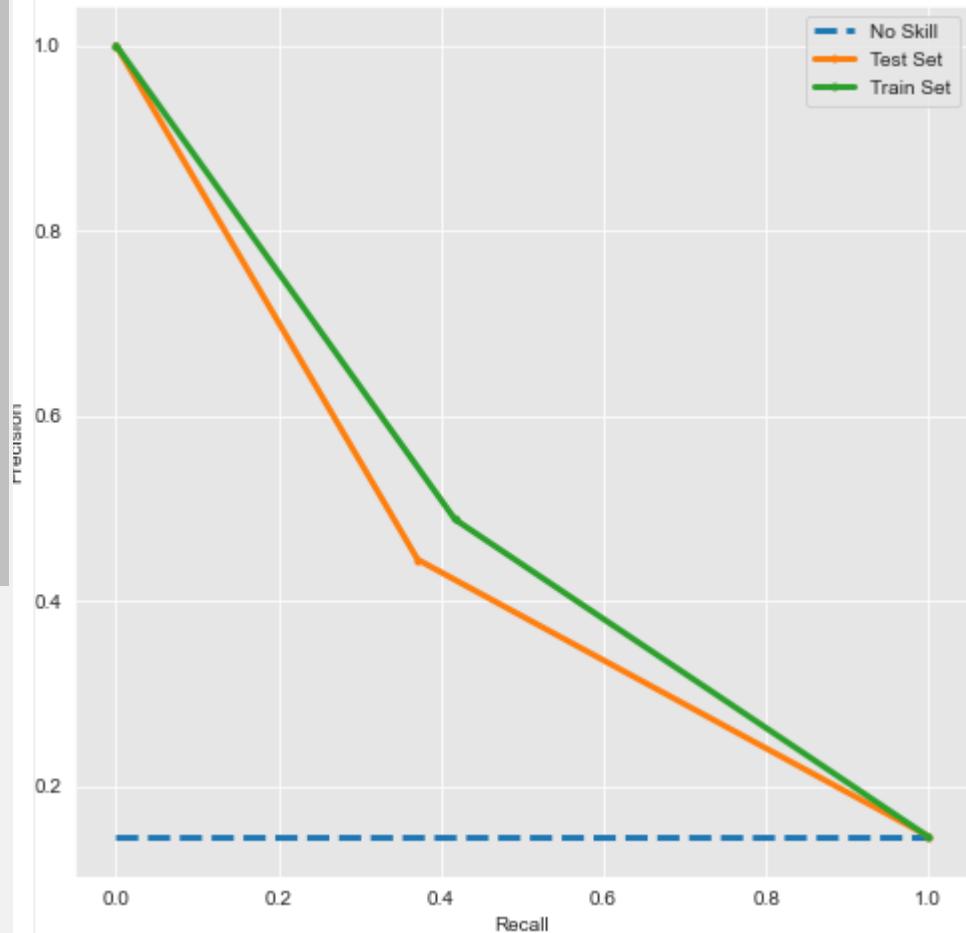
[14:57:27] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

[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.

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state count
 - 6.1.1 Common states
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model



Contents

In [82]:	nb_model.get_params
1 Abstract	out[82]:
2 Business Problem	ocound method BaseEstimator.get_params of GaussianNB()
3 Import Libraries	
4 Load Data	
5 Data Exploration	
6 Data Engineering	The mean cross validation is 88.7%. We have a precision of 1 and recall of 1 for the test set. We 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.
6.1 Group state don't the	
6.2 Drop unneeded	
6.3 Drop highly co	
6.4 One Hot Encod	
6.5 Train test split	
6.6 Standardize Fe	
6.7 Building anal	
7 Logistic Regression	
7.1 Optimization	
7.2 Best performance	
7.3 Results & Comments	
8 Decision Tree Classifier	
8.1 Hyperparameters	
8.1.1 Select Best	
8.2 Results & Comments	
9 Random Forest Classifier	
9.1 Hyperparameters	
9.2 Best Parameters	
9.3 Results & Comments	
10 KNN Classifier	
10.1 Tuning & Optimal	
10.1.1 Select Error	
10.1.1.1 K=3	
10.1.2 Iterate on K	
10.1.2.1 K=1	
10.2 Results & Comments	
11 XGBoost	
11.1 Results & comments	
12 Gaussian NB	
12.1 Results & Comments	
13 AdaBoostClassifier	
13.1 Results & Comments	
14 Gradient Boosting	
14.1 Results & comments	
15 SVM Classifier	
15.1 Tuning & Optimal	
15.1.1 Hyperparameters	
15.1.1.1 Best	

Results & Comments

The mean cross validation is 88.7%. We have a precision of 1 and recall of 1 for the test set. We 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.

13 AdaBoostClassifier

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

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
- 12 Gaussian NB
- 13 AdaBoostClassifier
- 14 Gradient Boosting
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

MODEL EVALUATION METRICS:

Confusion Matrix for train & test set:

Predicted	No Churn	Churn	All
Actual	Nan	Nan	Nan
No Churn	2280.0	0.0	2280.0
Churn	0.0	386.0	386.0
All	2280.0	386.0	2666.0

Predicted	No Churn	Churn	All
Actual	Nan	Nan	Nan
No Churn	570.0	0.0	570.0
Churn	0.0	97.0	97.0
All	570.0	97.0	667.0

Classification Report for train & 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

Results & Comparison

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

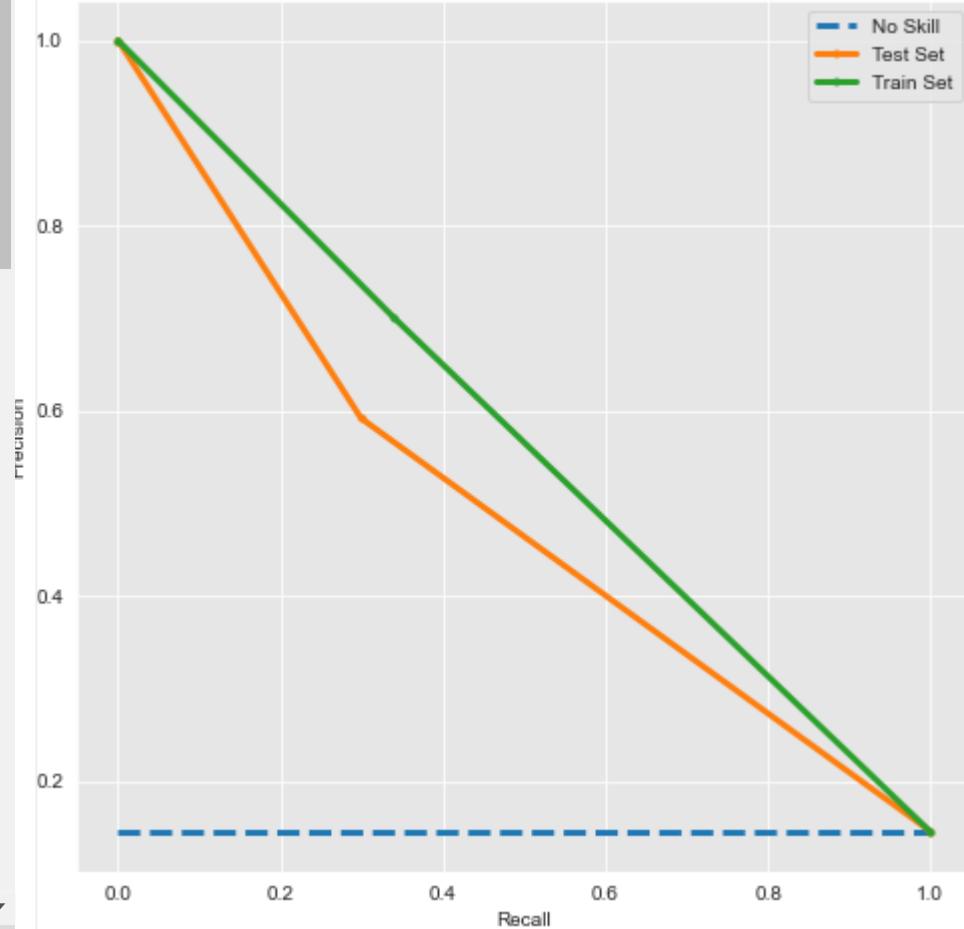
Results & Comparison

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.

Contents ⚙️

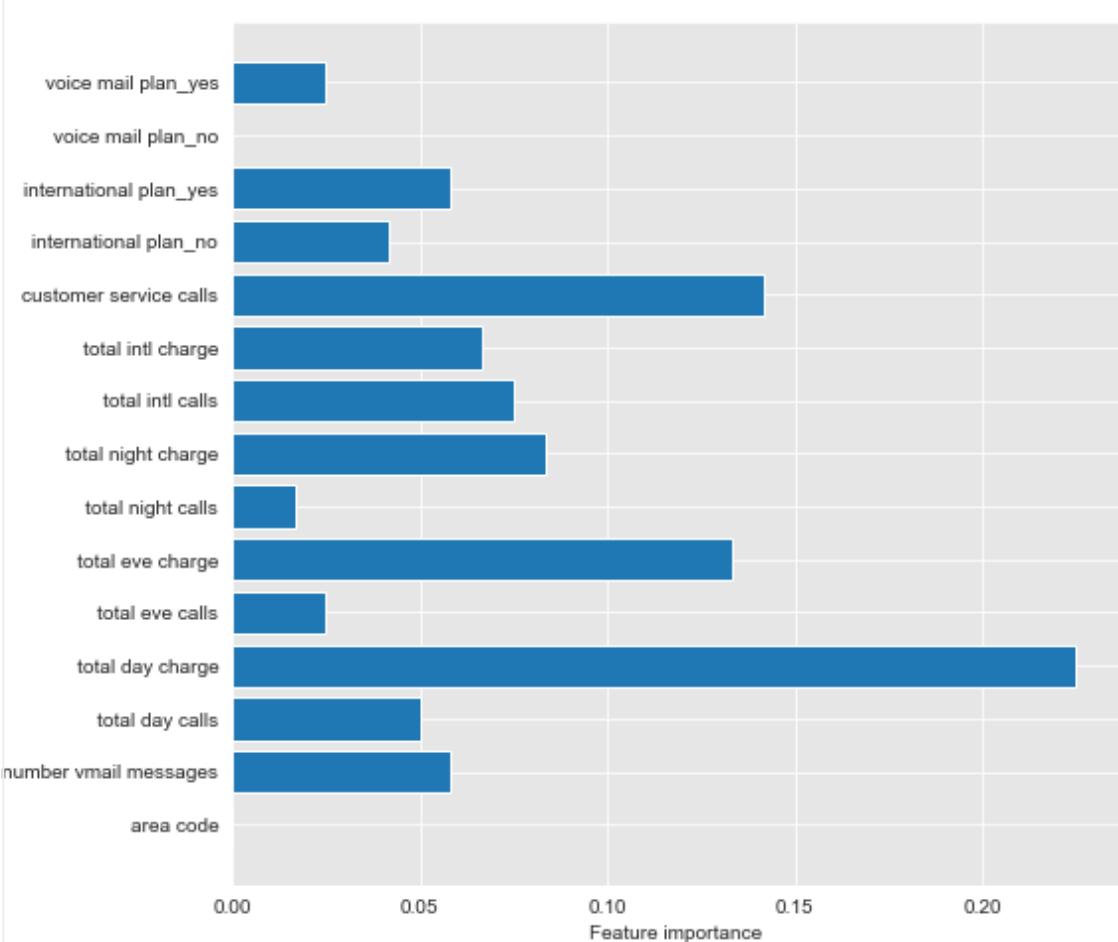
- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Model
 - ▼ 10.1.1 Select Evaluation Metric
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Model
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model



Contents

1	Abstract
2	Business Problem
3	Import Libraries
4	Load Data
5	Data Exploration
6	Data Engineering
6.1	Group state column
6.1.1	Comments
6.2	Drop unnecessary columns
6.3	Drop highly correlated columns
6.4	One Hot Encoding
6.5	Train test split
6.6	Standardize Features
6.7	Building analysis
7	Logistic Regression
7.1	Optimization
7.2	Best performance
7.3	Results Comparison
8	Decision Tree Classifier
8.1	Hyperparameters
8.1.1	Select Best Model
8.2	Results Comparison
9	Random Forest Classifier
9.1	Hyperparameters
9.2	Best Parameters
9.3	Results & Comparison
10	KNN Classifier
10.1	Tuning & Optimization
10.1.1	Select Best Model
10.1.1.1	K=3
10.1.2	Iterate on K
10.1.2.1	K=1
10.2	Results & Comparison
11	XGBoost
11.1	Results & Comparison
12	Gaussian NB
12.1	Results & Comparison
13	AdaBoostClassifier
13.1	Results & Comparison
14	Gradient Boosting
14.1	Results & Comparison
15	SVM Classifier
15.1	Tuning & Optimization
15.1.1	Hyperparameters
15.1.1.1	Best Model

1 plot_feature_importances(ada_model)

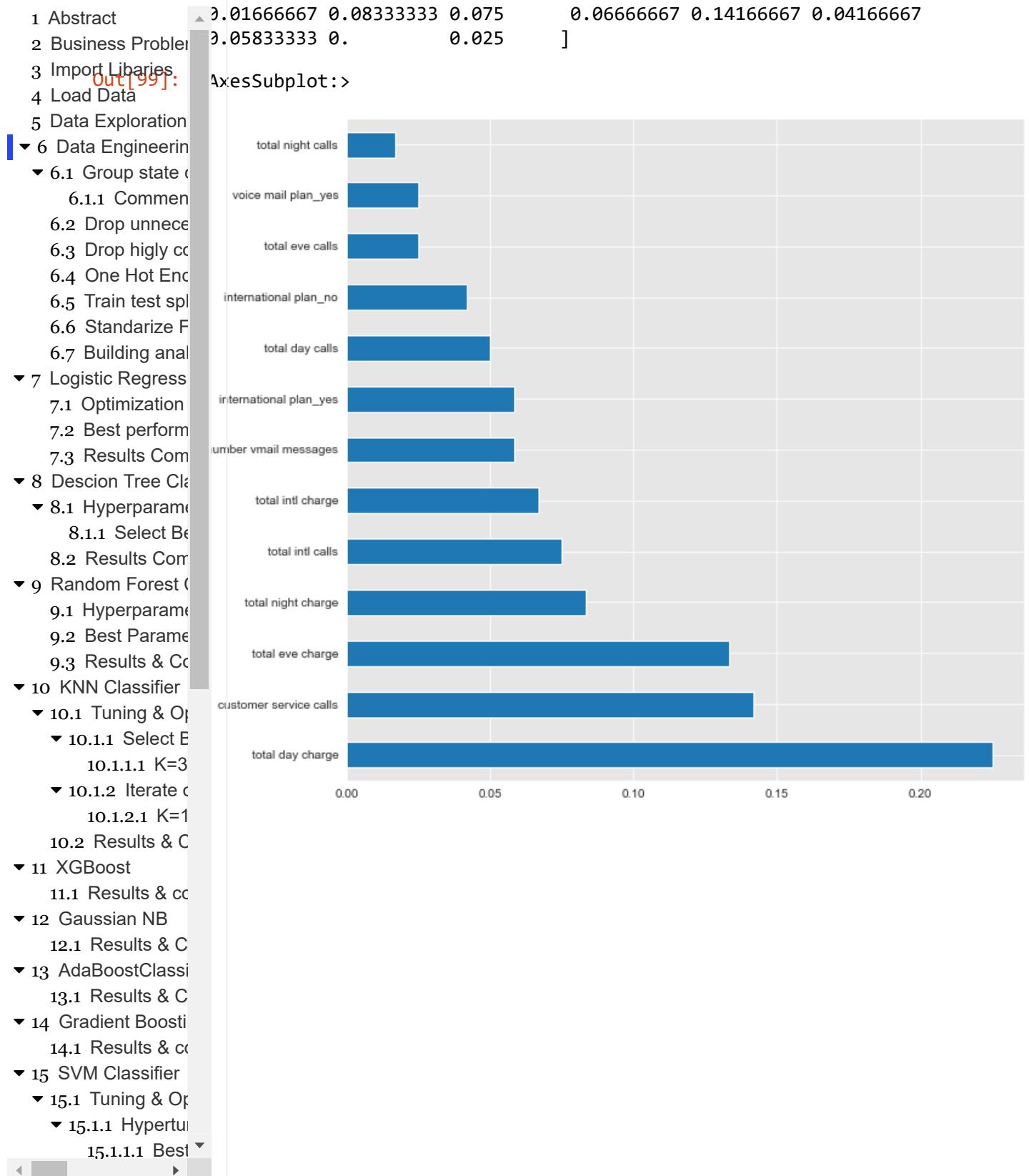


In [99]:

```

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

13.1 Results & Comments

Contents ↻ ⚙

- 1 Abstract 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.
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration The most important churn factors:
 - 1) Total day charge
 - 2) Customer service calls
 - 3) Total evening charge
 - 4) One Hot Encoding
 - 5) International calls
 - 6) Standardize Features
 - 7) Building gradient boosting model night charge
- 6 Data Engineering
 - 6.1 Group state by country
 - 6.1.1 Common prefix
 - 6.1.2 Drop unnecessary columns
 - 6.1.3 Drop highly correlated features
 - 6.1.4 One Hot Encoding
 - 6.1.5 Train test split
 - 6.1.6 Standardize Features
 - 6.1.7 Building gradient boosting model night charge
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classification
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classification
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Number of Neighbors
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.1.2 Iterate over K values
 - 10.1.2.1 K=1
 - 10.1.2.2 K=3
 - 10.1.2.3 K=5
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

12 Gradient Boosting Classifier

In [187]:

```

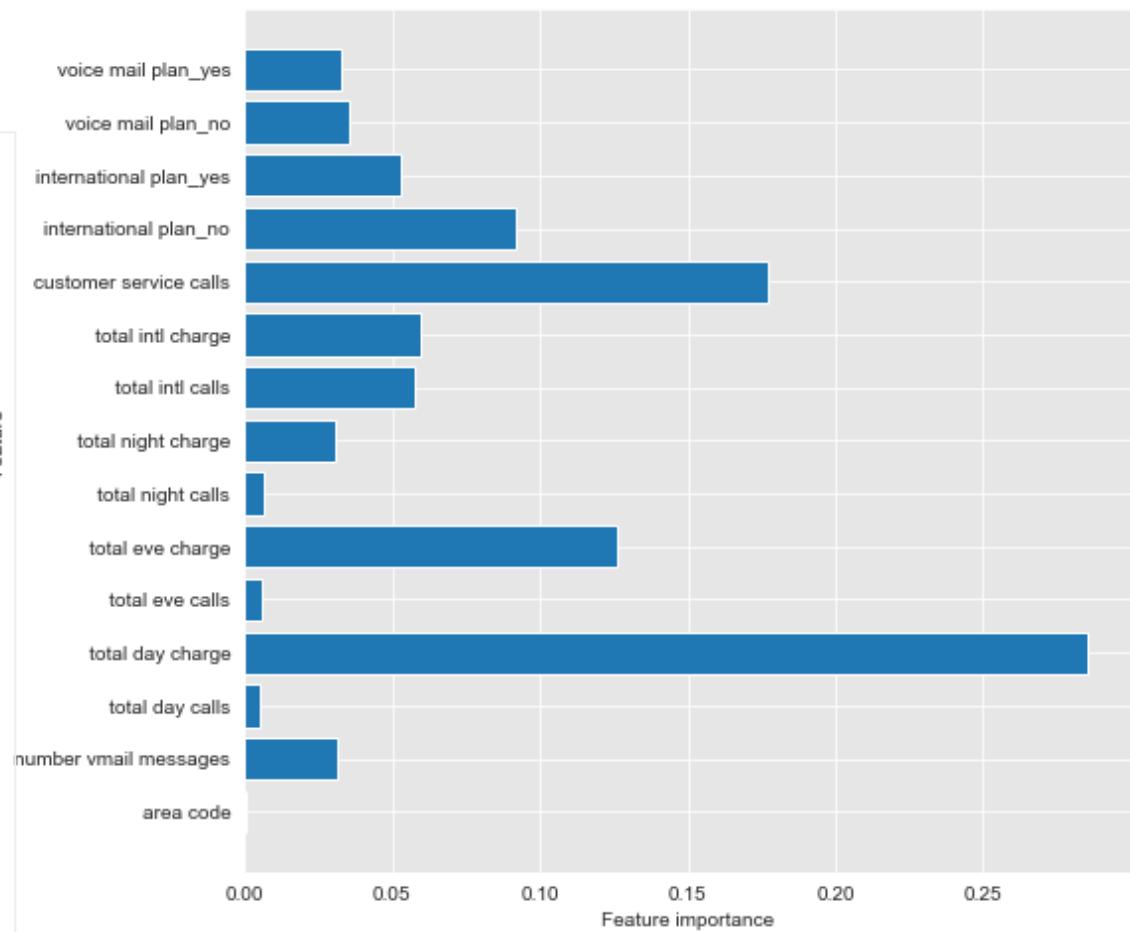
1 # Gradient Boosting
2 from sklearn.ensemble import GradientBoostingClassifier
3
4 gb = GradientBoostingClassifier(learning_rate=0.1, n_estimators=150, random_state=42)
5 gb_model=gb.fit(scaled_X_train,y_train)
6
7
8
9 y_pred_train = gb_model.predict(scaled_X_train)
10 y_pred_test = gb_model.predict(scaled_X_test)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

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state column
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

In [101]: 1 plot_feature_importances(gb_model)



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

```
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))
```

Contents ↗

1 Abstract
2 Business Problem
3 Import Libraries
4 Load Data

5 Data Exploration

6 Data Engineering

6.1 Group state

6.1.1 Comment

6.2 Drop unnecessary

6.3 Drop highly co

6.4 One Hot Encod

6.5 Train test spli

6.6 Standardize Fe

6.7 Building anal

7 Logistic Regression

7.1 Optimization

7.2 Best perform

7.3 Results Com

8 Descion Tree Clas

8.1 Hyperparamet

8.1.1 Select Be

8.2 Results Com

9 Random Forest Cl

9.1 Hyperparamet

9.2 Best Paramet

9.3 Results & Co

10 KNN Classifier

10.1 Tuning & Opt

10.1.1 Select E

10.1.1.1 K=3

10.1.2 Iterate c

10.1.2.1 K=1

10.2 Results & C

11 XGBoost

11.1 Results & Co

12 Gaussian NB

12.1 Results & C

13 AdaBoostClassi

13.1 Results & C

14 Gradient Boosti

14.1 Results & C

15 SVM Classifi

15.1 Tuning & Opt

15.1.1 Hyperpar

15.1.1.1 Best

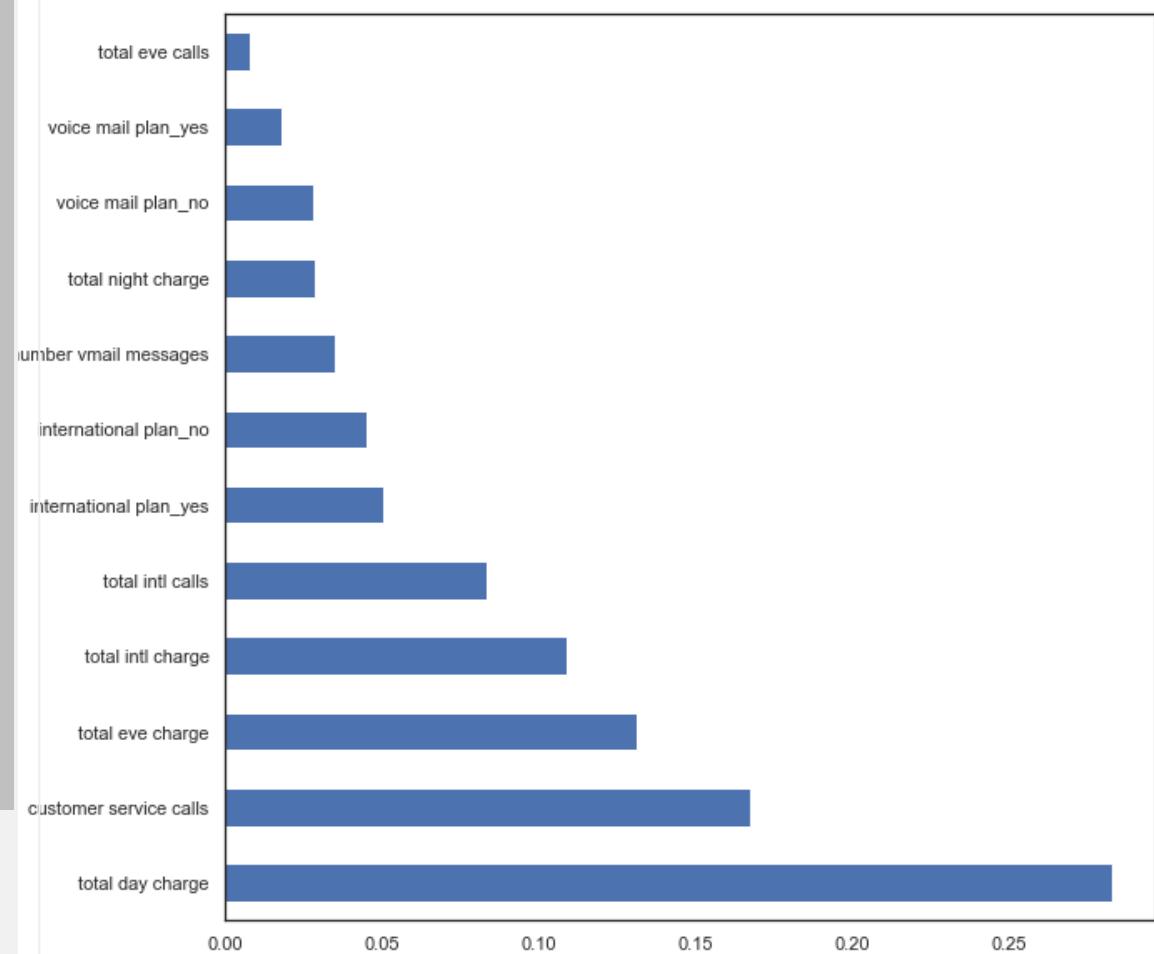
The most important churn factors:

[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:>



Results & comments

The mean cross validation is 88.7%. We have a precision of 1 and recall of 1 for the test set. We have overfitting. It means we are able to detect the 100% of churn customers and 100% of non-churners 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 ⚙⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering

- 16.1 [103]: Standardization

- 6.1.1 Comments
- 6.2 Drop unnecessary columns
- 6.3 Drop highly correlated columns
- 6.4 One Hot Encoding
- 6.5 Train test split

- 16.6 [\\$84]: F

- 6.7 Building analysis

- 7 Logistic Regression

- 7.1 Optimization
- 7.2 Best performance
- 7.3 Results Comparison

- 8 IDes[105]: Tree Classifiers

- 8.1 Hyperparameters

- 8.1.1 Select Best Parameters
- 8.2 Results Comparison

- 9 Random Forest Classifier

- 9.1 Hyperparameters
- 9.2 Best Parameters
- 9.3 Results & Comparison

- 10 KNN Classifier

- 10.1 Tuning & Optimal Parameters
- 10.1.1 Select Error Function
- 10.1.1.1 K=3
- 10.1.2 Iterate cross-validation
- 10.1.2.1 K=1

- 10.2 Results & Comparison

- 11 XGBoost

- 11.1 Results & Comparison

- 12 Gaussian NB

- 12.1 Results & Comparison

- 13 AdaBoostClassifier

- 13.1 Results & Comparison

- 14 Gradient Boosting

- 14.1 Results & Comparison

- 15 SVM Classifier

- 15.1 Tuning & Optimal Parameters

- 15.1.1 Hyperparameters

- 15.1.1.1 Best Parameters

SVM Classifier

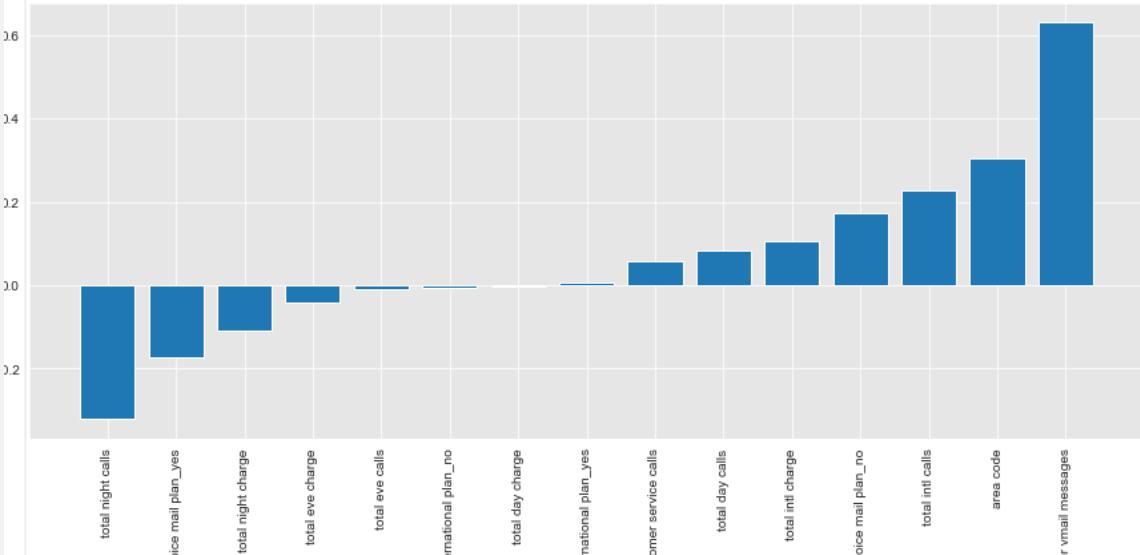
```

1 #SVM M 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 & test set:

Predicted	No Churn	Churn	All
Actual	Nan	Nan	Nan
No Churn	510.0	328.0	838.0
Churn	85.0	77.0	162.0
All	595.0	405.0	1000.0

Predicted	No Churn	Churn	All
Actual	Nan	Nan	Nan
No Churn	347.0	223.0	570.0
Churn	65.0	32.0	97.0
All	412.0	255.0	667.0

Classification Report for train & test set

train set

	precision	recall	f1-score	support
7.1 Optimization	False	0.86	0.61	0.71
	True	0.19	0.48	0.27
8.1 Hyperparameters	accuracy		0.59	1000
	macro avg	0.52	0.54	0.49
	weighted avg	0.75	0.59	0.64

Random Forest

test set

	precision	recall	f1-score	support
10.1 Tuning & Opt	False	0.84	0.61	0.71
	True	0.13	0.33	0.18
10.1.1 Select Est	accuracy		0.57	667
	macro avg	0.48	0.47	0.44
	weighted avg	0.74	0.57	0.63

10.2 Results & C

XGBoost

11.1 Results & C

Gaussian NB

12.1 Results & C

AdaBoostClassi

13.1 Results & C

Gradient Boosti

14.1 Results & C

SVM Classifie

15.1 Tuning & Op

15.1.1 Hypertun

15.1.1.1 Best

ohen's Kappa for train and test set:

0.0523 -0.0366

2 score for train and test set:

0.3656 0.2488

oc 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 cha

nged 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.

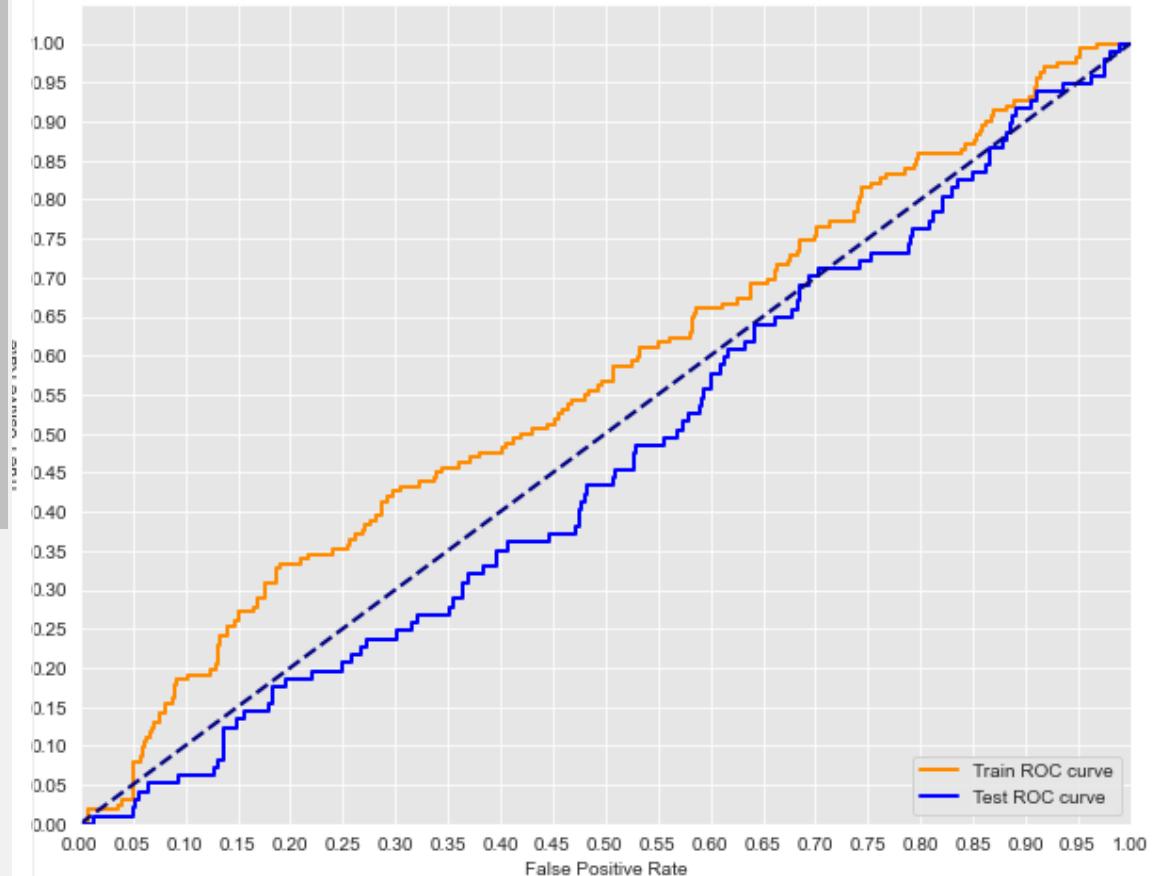
[15:00:58] WARNING: ..\src\learner.cc:1061: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Common states
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

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

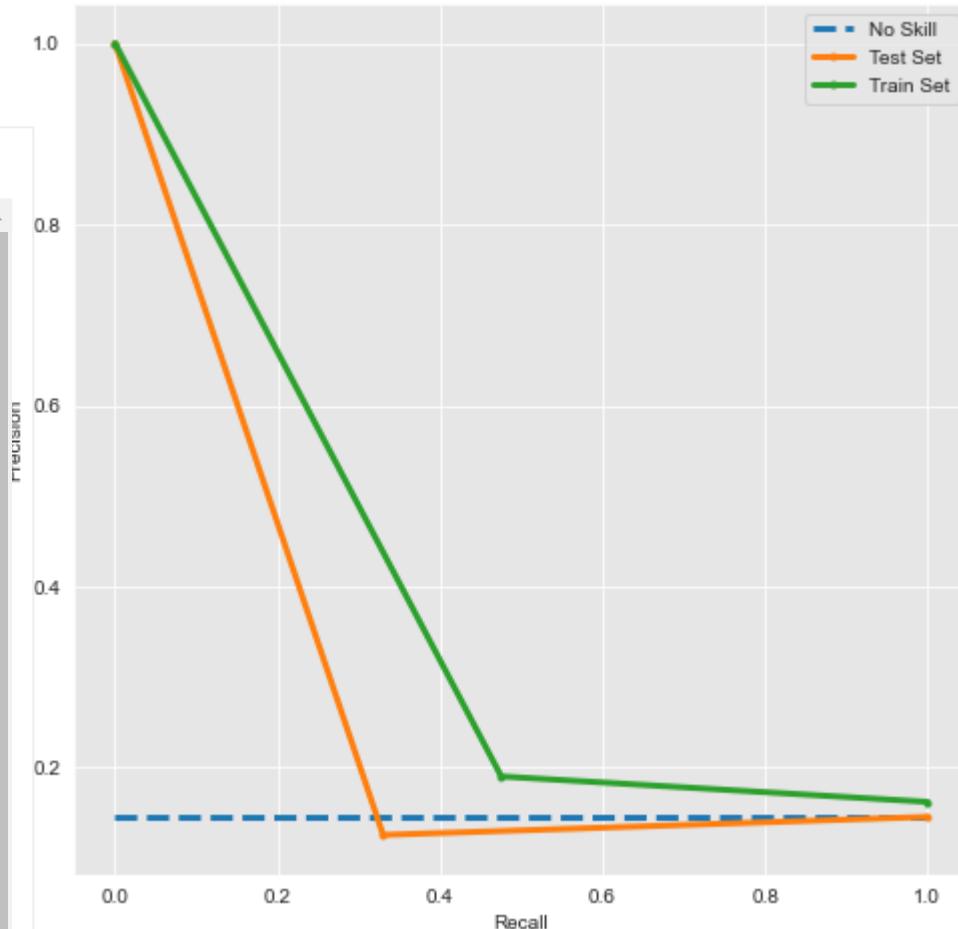
Receiver operating characteristic (ROC) Curve for Training and Testing Sets



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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 7.1 Optimization
- 7.2 Best performance
- 7.3 Results Comparison
- 8 Descion Tree Classifier
- 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
- 9.1 Hyperparameters
- 9.2 Best Parameters
- 9.3 Results & Comparison
- 10 KNN Classifier
- 10.1 Tuning & Optimization
 - 10.1.1 Select Best Parameters
 - 10.1.1.1 K=1
 - 10.1.2 Iterate over K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hypertuning
 - 15.1.1.1 Best Parameters

15 Tuning & Optimization

1 Hypertuning of SVM Classifier throgh GridSearchCV

```
In [109]: #SVM Classifier (balanced class weight and gridsearch)
1
2
3 param_grid = {'C': [0.1, 1, 10, 100, 1000],
4                 'gamma': [1, 0.1, 0.01, 0.001, 0.0001],
5                 'kernel': ['rbf', "linear"]}
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_train)
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_
```

MODEL EVALUATION METRICS:

Confusion Matrix for train & test set:

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	2154.0	126.0	2280.0
Churn	37.0	349.0	386.0
All	2191.0	475.0	2666.0

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	570.0	0.0	570.0
Churn	97.0	0.0	97.0
All	667.0	0.0	667.0

Classification Report for train & test set

```
1 svc2.best_params_
```

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

.1 Best parameters for SVM Classifier

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classification
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classification
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
- 10.2 Results & Comparison
- 11 XGBoost Classification
 - 11.1 Results & Comparison
- 12 Gaussian NB
- 12.1 Results & Comparison
- 13 AdaBoostClassification
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

```
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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Feature
 - 10.1.1.1 K-Nearest Neighbors
 - 10.1.2 Iterate over different K values
 - 10.1.2.1 The best K value
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

MODEL EVALUATION METRICS:

Confusion Matrix for train & test set:

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	1995.0	289.0	2284.0
Churn	340.0	42.0	382.0
All	2335.0	331.0	2666.0

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	502.0	64.0	566.0
Churn	94.0	7.0	101.0
All	596.0	71.0	667.0

Classification Report for train & test set

15 Results & Comments

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.

```
In [459]: 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
```

15.2.1 SVM with RFF

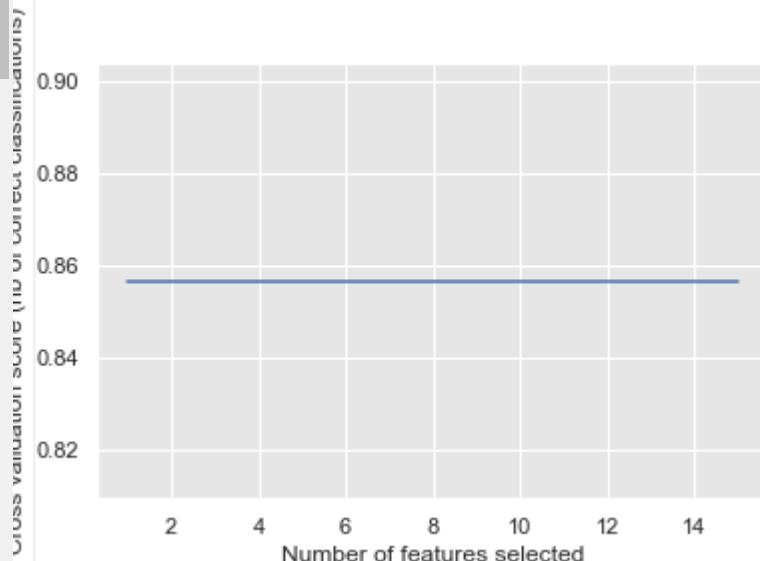
In [153]:

```

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 # Build a classification task using 3 informative features
7
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



In [154]:

```

1 rfecv.support_

```

Out[154]:

```

array([False,  True, False, False, False, False, False, False,
       False, False, False, False, False])

```

In [155]: 1 pd.DataFrame(data=rfecv.support_.reshape(1, -1), columns= X.columns)

Out[155]:

	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	False	True	False	False	False	False	False	False	False	False

Contents ⚙

1 Abstract										
2 Business Problem										
3 Import Libraries										
In [156]: 1 rfevcv.ranking_										
array([10, 1, 8, 7, 9, 4, 13, 15, 11, 14, 5, 12, 6, 3, 2])										
In [157]: 1 pd.DataFrame(data=rfevcv.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	churn	
0	10	1	8	7	9	4	13	15	11	

2 SVM with SelectKBest

15

7 Logistic Regression										
7.1 Optimization										
7.2 Best performance										
7.3 Results & Conclusion										
8 Descion Tree Classifier										
8.1 Hyperparameters										
8.1.1 Select Best Parameters										
8.2 Results Comparison										
9 Random Forest Classifier										
9.1 Hyperparameters										
9.2 Best Parameters										
9.3 Results & Conclusion										
10 KNN Classifier										
10.1 Tuning & Optimization										
10.1.1 Select Error Function										
10.1.1.1 K=3										
10.1.2 Iterate over different K										
10.1.2.1 K=1										
10.2 Results & Conclusion										
11 XGBoost										
11.1 Results & Conclusion										
12 Gaussian NB										
12.1 Results & Conclusion										
13 AdaBoostClassifier										
13.1 Results & Conclusion										
14 Gradient Boosting										
14.1 Results & Conclusion										
15 SVM Classifier										
15.1 Tuning & Optimization										
15.1.1 Hyperparameters										
15.1.1.1 Best Parameters										

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
13 ##### Uivariate feature selection with F-test for feature scoring
14 # We use the default selection function to select the four
15 # most significant features
16 selector = SelectKBest(f_classif, k=4)
17 selector.fit(scaled_X_train, y_train)
18 scores = -np.log10(selector.pvalues_)
19 scores /= scores.max()
20 plt.bar(X_indices - .45, scores, width=.2,
21         label='Univariate score ($-Log(p_{value})$)')
22
23 #####
24 # Compare to the weights of an SVM
25 clf = make_pipeline(LinearSVC())
26 clf.fit(scaled_X_train, y_train)
27 print('Classification accuracy without selecting features: {:.3f}'
28       .format(clf.score(scaled_X_test, y_test)))
29
30
31 svm_weights = np.abs(clf[-1].coef_).sum(axis=0)
32 svm_weights /= svm_weights.sum()
33
34 plt.bar(X_indices - .25, svm_weights, width=.2, label='SVM weight')
35
36 clf_selected = make_pipeline(
37     SelectKBest(f_classif, k=4), LinearSVC())
38
39 clf_selected.fit(scaled_X_train, y_train)
40 print('Classification accuracy after univariate feature selection: {:.3f}'
41       .format(clf_selected.score(scaled_X_test, y_test)))
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()

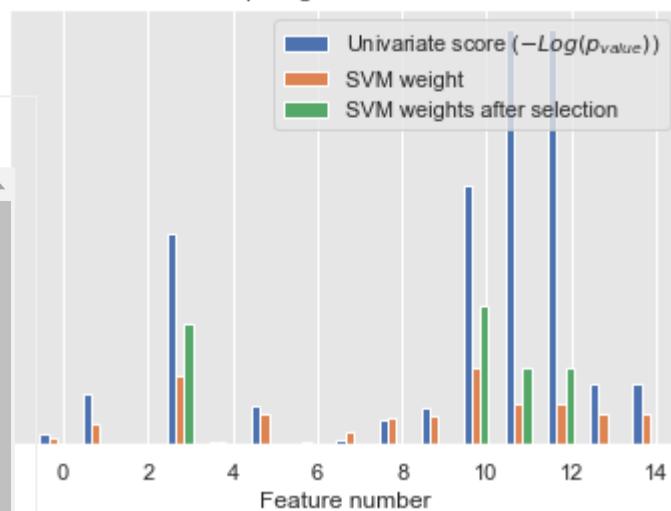
```

Contents ⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Entropy
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hypertuning
 - 15.1.1.1 Best

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

Comparing feature selection



```
In [139]:  
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

	accuracy	macro avg	weighted avg
accuracy	0.84	0.63	0.79
macro avg	0.84	0.53	0.52
weighted avg	0.84	0.84	0.79

15 Results & Comments

VC classifiers performed poorly even with RFE and SELECTK.

```
In [163]:  
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

16 Stacking

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

- 1 Abstract
- 2 Business Problem
- 3 Import Libaries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unneccesary columns
 - 6.3 Drop higly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standarize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Descion Tree Classification
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classification
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Model
 - ▼ 10.1.1 Select Euclidean Distance
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate cross validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
- 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
- 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
- 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Model
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

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

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	2284.0	0.0	2284.0
Churn	55.0	327.0	382.0
All	2339.0	327.0	2666.0

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	560.0	6.0	566.0
Churn	26.0	75.0	101.0
All	586.0	81.0	667.0

Classification Report for train & test set

In [50]: 1 clf.get_params

```
Out[50]: <bound method _BaseHeterogeneousEnsemble.get_params of StackingClassifier(estimators=[('tree', DecisionTreeClassifier(max_depth=10, random_state=5)), ('forest', RandomForestClassifier(max_depth=10, random_state=5)), ('gb', GradientBoostingClassifier(learning_rate=0.1, n_estimators=12, random_state=5))], final_estimator=LogisticRegression())>
```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of customers
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over different K values
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

16 Results & Comments

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

17 SGDClassifier (Lasso & Ridge Regularizations)

```
In [161]: #SCDClassifier (balanced class weight)
1
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 ⚙️

1 Abstract
2 Business Problem
3 Import Libraries
4 Load Data
5 Data Exploration
6 Data Engineering
6.1 Group state & city
6.1.1 Comments
6.2 Drop unnecessary columns
6.3 Drop highly correlated columns
6.4 One Hot Encoding
6.5 Train test split
6.6 Standardize Features
6.7 Building analysis
7 Logistic Regression
7.1 Optimization
7.2 Best performance
7.3 Results Comparison
8 Descion Tree Classifier
8.1 Hyperparameter Tuning
8.1.1 Select Best Model
8.2 Results Comparison
9 Random Forest Classifier
9.1 Hyperparameters Tuning
9.2 Best Parameters
9.3 Results & Comparison
10 KNN Classifier
10.1 Tuning & Optimize
10.1.1 Select Error Metric
10.1.1.1 K=3
10.1.2 Iterate over K
10.1.2.1 K=1
accuracy
10.1.2.2 K=3
macro avg
10.1.2.3 K=5
weighted avg
10.2 Results & Comparison
11 XGBoost Classifier
11.1 Results & Comparison
12 Gaussian NB
12.1 Results & Comparison
13 AdaBoostClassifier
13.1 Results & Comparison
14 Gradient Boosting
14.1 Results & Comparison
15 SVM Classifier
15.1 Tuning & Optimize
15.1.1 Hyperparameter Tuning
15.1.1.1 Best Model
15.2 Results & Comparison
16 Model Evaluation Metrics
16.1 Confusion Matrix for train & test set:
16.2 Classification Report for train & test set
17 Model Comparison
17.1 Confusion Matrix for main set
17.2 Classification Report for main set
18 Final Summary

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:

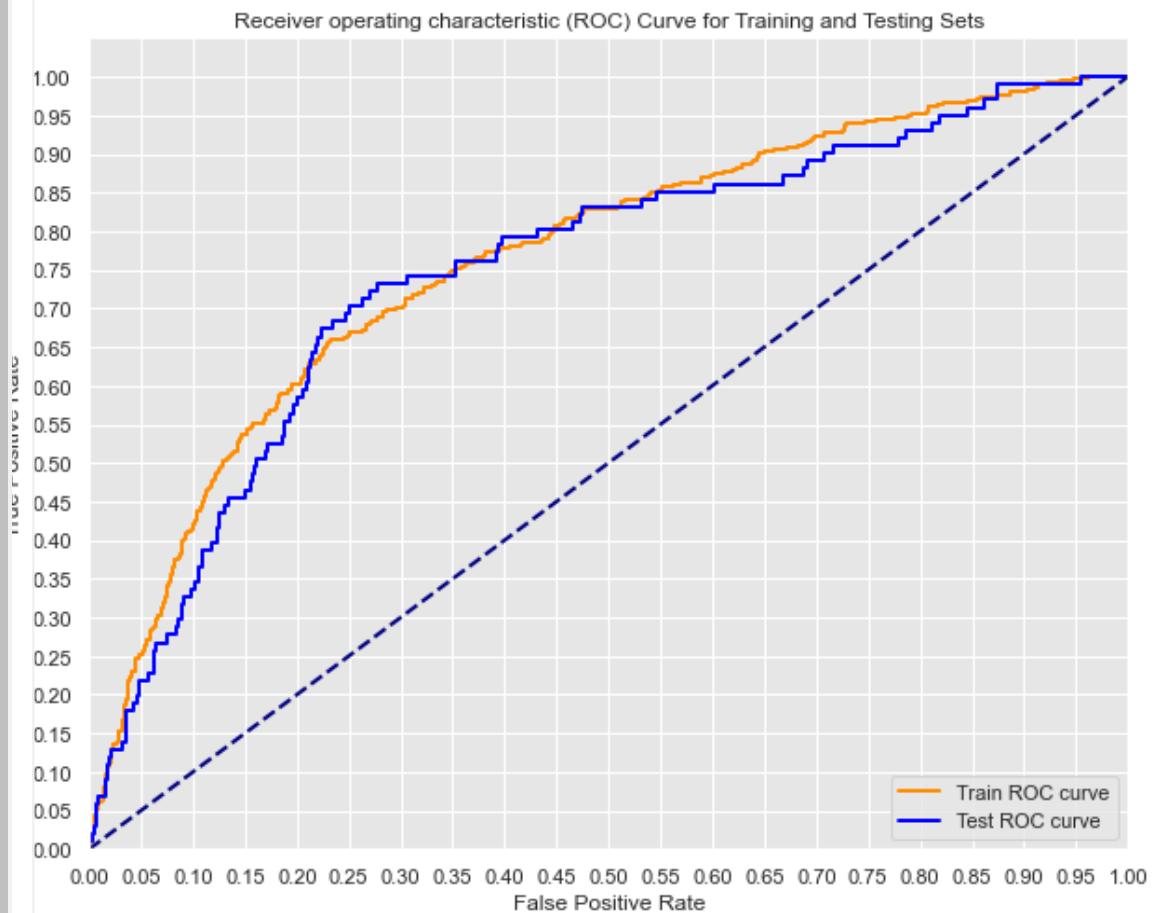
0.6961 0.7187

Mean Cross Validation Score:

0.9502

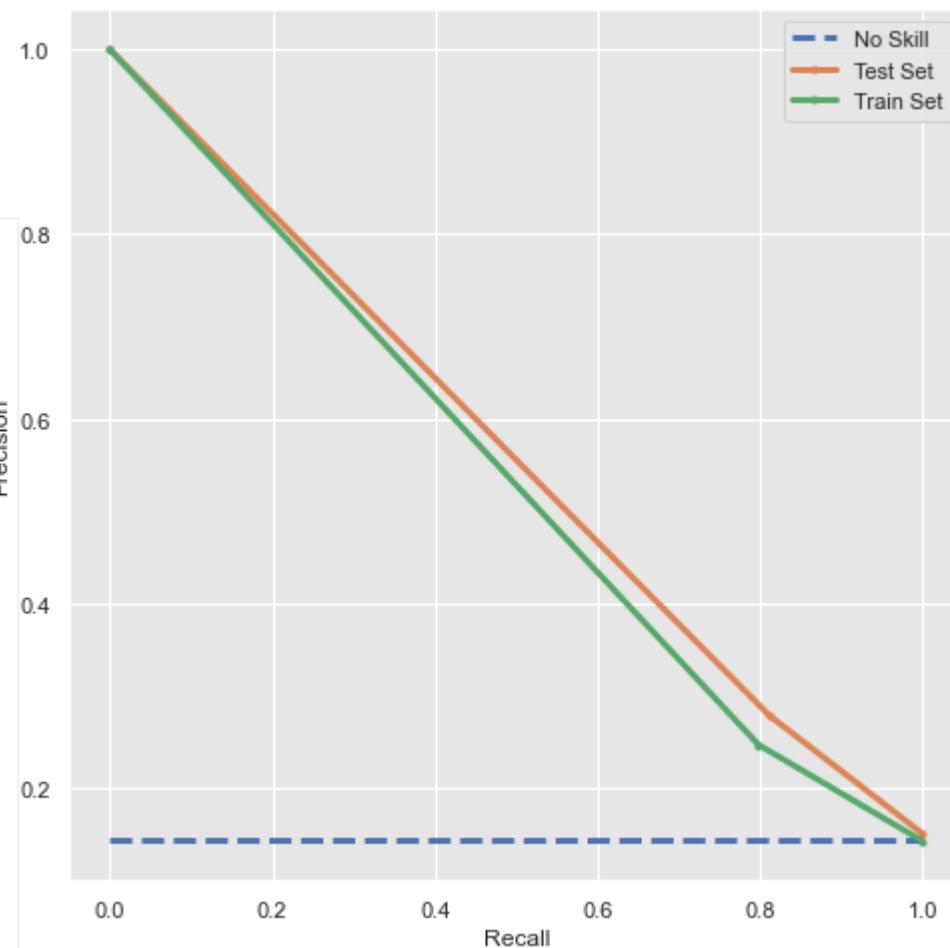
Contents ⚙️⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state column
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model



Training AUC: 0.76274

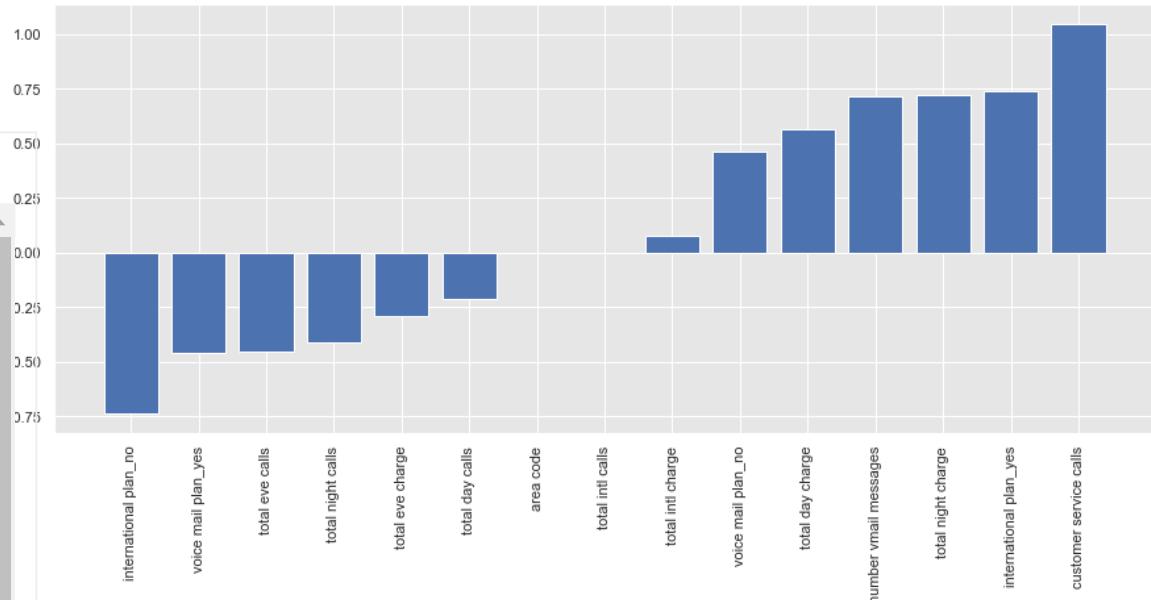
Testing AUC: 0.74873



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimal Model
 - ▼ 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - ▼ 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimal Model
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

In [162]: 1 plot_coefficients(scd)



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classification
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classification
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Number of Neighbors
 - 10.1.1 Select Best K
 - 10.1.1.1 K=3
 - We can see that we don't have overfitting. It means we are able to detect the 81% of churn customers and 28% of non-churn customers we say 'churn' was true.
 - 10.1.2 Iterate over different K values
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian Naive Bayes
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameter Tuning
 - 15.1.1.1 Best Parameters

Results & Comments

The mean cross validation is 95%. We have a precision of 0.28 and recall of 0.81 for the test set. This means we are able to detect the 81% of churn customers and 28% of non-churn 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,
5                                random_state=0)
6 distributions = dict(C=uniform(loc=0, scale=4),
7                       penalty=['l2', 'l1'])
8 clf = RandomizedSearchCV(logistic, distributions, random_state=0)
9 search = clf.fit(scaled_X_train, y_train)
10
11 search.best_params_

```

'C': 1.75034884505077, 'penalty': 'l2'}

Contents ↗

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering

- 6.1 Group state
- 6.1.1 Commence
- 6.2 Drop unnecessary columns
- 6.3 Drop highly correlated columns
- 6.4 One Hot Encoding
- 6.5 Train test split
- 6.6 Standardize Features
- 6.7 Building analysis

- 7 Logistic Regression
- 7.1 Optimization
- 7.2 Best performance
- 7.3 Results Comparison
- 8 Descion Tree Classification
- 8.1 Hyperparameters
- 8.1.1 Select Best Model
- 8.2 Results Comparison

- 9 Random Forest Classification
- 9.1 Hyperparameters
- 9.2 Best Parameters
- 9.3 Results & Comparison

- 10 KNN Classifier
- 10.1 Tuning & Optimization
- 10.1.1 Select Ensemble
- 10.1.1.1 K=3
- 10.1.2 Iterate cross-validation
- 10.1.2.1 K=1
- 10.2 Results & Comparison

- 11 XGBoost
- 11.1 Results & Comparison
- 12 Gaussian NB
- 12.1 Results & Comparison

- 13 AdaBoostClassifier
- 13.1 Results & Comparison
- 14 Gradient Boosting
- 14.1 Results & Comparison

- 15 SVM Classifier
- 15.1 Tuning & Optimization
- 15.1.1 Hyperparameters
- 15.1.1.1 Best Model

```

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)

```

```

1 y_pred_train = clf.predict(scaled_X_train)
2 y_pred_test = clf.predict(scaled_X_test)
3
4
5 model_evaluation(scaled_X_train, scaled_X_test, y_train, y_test, y_pred_train,
6 roc_curve_and_auc(clf, scaled_X_train, scaled_X_test, y_train, y_test)
7
8 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)

```

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

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	2498.0	72.0	2570.0
Churn	337.0	92.0	429.0
All	2835.0	164.0	2999.0

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	273.0	7.0	280.0
Churn	43.0	11.0	54.0
All	316.0	18.0	334.0

Classification Report for train & test set**Results & Comments**

This model had a mean cross validation of 85.66%. Only 20% of the time the churn was true with 85.66% precision.

19 SGDClassifier - RandomizedSearchCV

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.1.2 Iterate over K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

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 setting .

- 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
 - 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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- ▼ 6 Data Engineering
 - ▼ 6.1 Group state by city
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- ▼ 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- ▼ 8 Decision Tree Classifier
 - ▼ 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
 - 8.3 Results & Comparison
- ▼ 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- ▼ 10 KNN Classifier
 - ▼ 10.1 Tuning & Optimization
 - In [168]: `random_search.best_params_`
 - Out[168]: `{'alpha': 0.0006626483822186211, 'average': False, 'l1_ratio': 0.503491101662517}`
 - 10.2 Results & Comparison
- ▼ 11 XGBoost
 - 11.1 Results & Comparison
- ▼ 12 Gaussian NB
 - 12.1 Results & Comparison
- ▼ 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- ▼ 14 Gradient Boosting
 - 14.1 Results & Comparison
- ▼ 15 SVM Classifier
 - ▼ 15.1 Tuning & Optimization
 - ▼ 15.1.1 Hyperparameters
 - 15.1.1 Best Model

19 Tuning & Optimization

```
1 random_search.best_params_
{'alpha': 0.0006626483822186211,
'average': False,
'l1_ratio': 0.503491101662517}
```

In [169]:

```

1 # build a classifier
2 clf = SGDClassifier(loss='hinge', penalty='elasticnet',
3                      fit_intercept=True)
4 param_dist = {'alpha': [0.0006626483822186211], 'average': [False], 'l1_
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 ↗⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Common features
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- In [169]:
- 7.1 Optimization
- 7.2 Best performance
- 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Error Metric
 - 10.1.1.1 K=3
 - 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

```

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_train)
6 roc_curve_and_auc(rs, scaled_X_train, scaled_X_test, y_train, y_test)
7
8 plot_pr_rc_curve(y_test, y_pred_test, y_train, y_pred_train)

```

MODEL EVALUATION METRICS:

Confusion Matrix for train & test set:

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	2284.0	0.0	2284.0
Churn	381.0	1.0	382.0
All	2665.0	1.0	2666.0

Predicted	No Churn	Churn	All
Actual	NaN	NaN	NaN
No Churn	566.0	0.0	566.0
Churn	101.0	0.0	101.0
All	667.0	0.0	667.0

Classification Report for train & test set

19. Results & Comments

This model performed poorly.

20 Data Interpretation

20.1 Most Important features According to DecisionTree Model

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classification
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classification
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over different values of K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters

In [189]:

```

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
28 # Model Performance
29 # Test set predictions
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: 95.15%

Mean Test Score: 93.85%

	precision	recall	f1-score	support
False	0.94	0.99	0.96	566
True	0.89	0.67	0.77	101

	accuracy	0.94	0.96	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

20.1.1 Top 5 Important Features Affecting Churn Rate

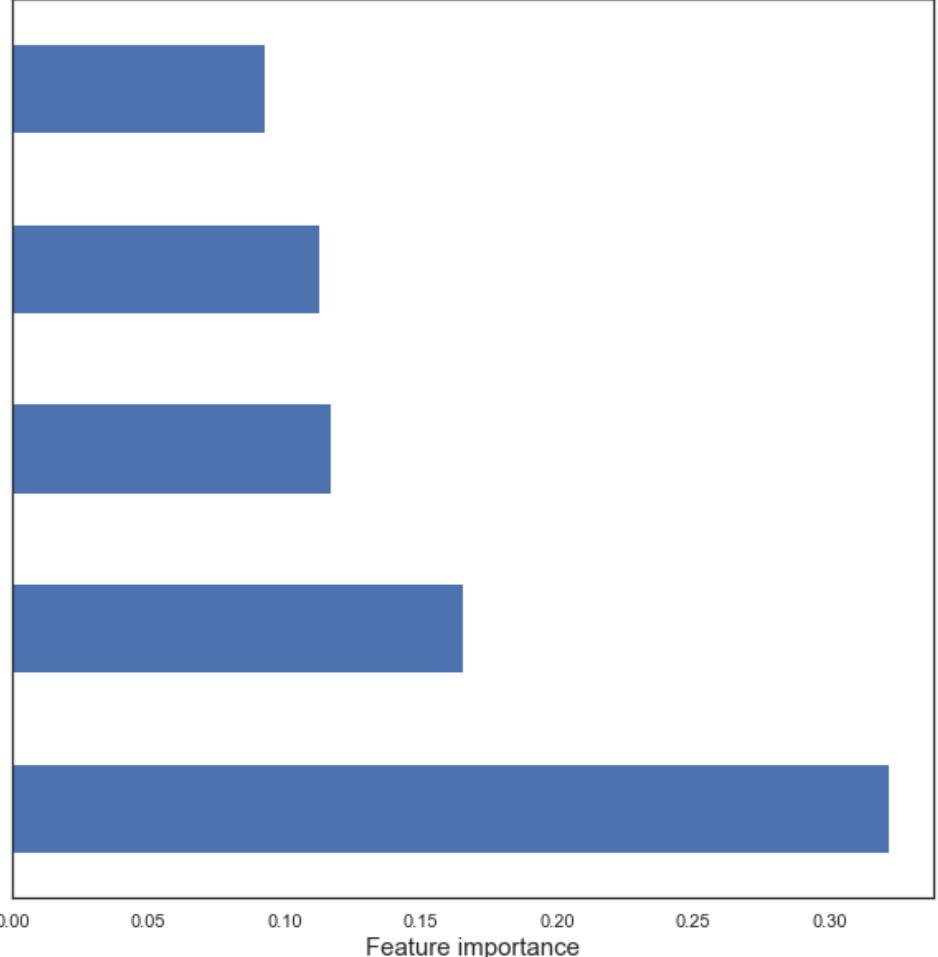
Contents ⚙️

1 Abstract:	1 feature_importance = tree_clf.feature_importances_
2 Business Problem:	2 print(tree_clf.feature_importances_)
3 Import Libaries:	3 feat_importances = pd.Series(tree_clf.feature_importances_, index=X.columns)
4 Load Data:	4 feat_importances = feat_importances.nlargest(5)
5 Data Exploration:	5 feat_importances.plot(kind='barh', figsize=(10,10), color="b")
6 Data Engineering:	6 plt.xlabel('Feature importance', fontsize=15)
6.1 Group state of account:	7 plt.ylabel('Positive Churn Factors', fontsize=15)
6.1.1 Comments:	8 plt.title('Top 5 important features that affect positively on the churn rate')
6.2 Drop unnecessary columns:	9 plt.show()
6.3 Drop highly correlated columns:	0. international plan_no 0. . 0.32193796 0.0134261 0.11687289
6.4 One Hot Encoding:	0.00422496 0.02131784 0.08281552 0.11242512 0.1652558 0.
6.5 Train test split:	0.0928464 0. . 0.0688774]
6.6 Standardize Features:	
6.7 Building analysis:	
7 Logistic Regression:	
7.1 Optimization:	
7.2 Best performance:	
7.3 Results Comparison:	
8 Decision Tree Classifier:	
8.1 Hyperparameters:	
8.1.1 Select Best Parameters:	
8.2 Results Comparison:	
9 Random Forest Classifier:	
9.1 Hyperparameters:	
9.2 Best Parameters:	
9.3 Results & Comparison:	
10 KNN Classifier:	
10.1 Tuning & Optimization:	
10.1.1 Select Error Function:	
10.1.1.1 K=3:	
10.1.2 Iterate over K:	
10.1.2.1 K=1:	
10.2 Results & Comparison:	
11 XGBoost:	
11.1 Results & Comparison:	
12 Gaussian NB:	
12.1 Results & Comparison:	
13 AdaBoostClassifier:	
13.1 Results & Comparison:	
14 Gradient Boosting:	
14.1 Results & Comparison:	
15 SVM Classifier:	
15.1 Tuning & Optimization:	
15.1.1 Hyperparameters:	
15.1.1.1 Best Parameters:	

Out[196]:

```
ext(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



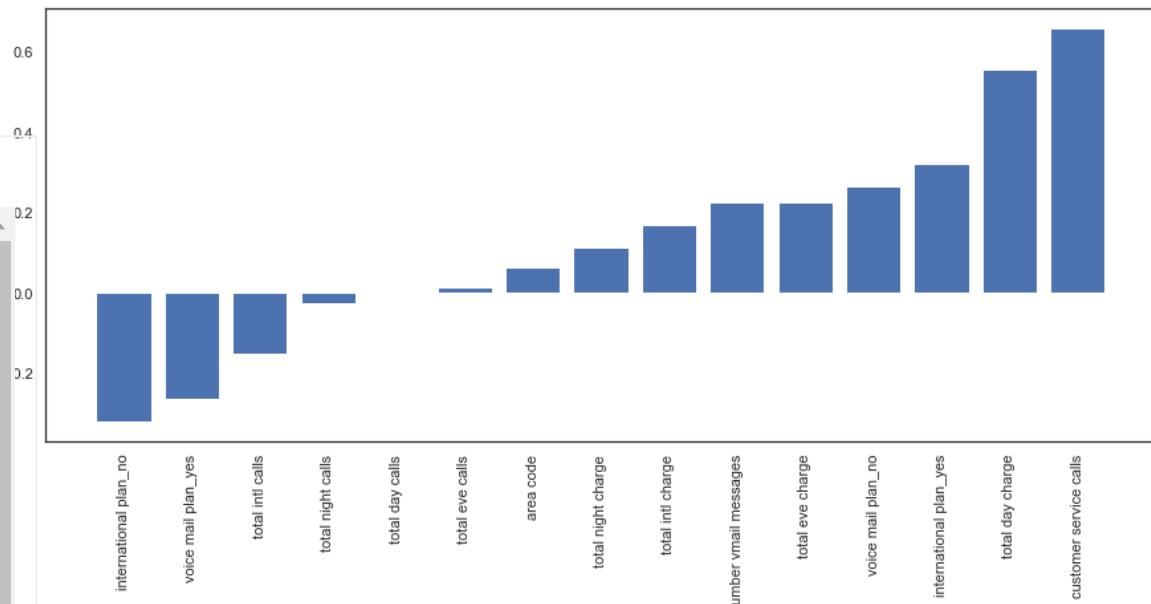
20.2 Top Positive and Negative Features According to Logistic Regression

```
In [178]: 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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state column
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Best Model
 - 10.1.1.1 K=3
 - 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

In [179]: 1 plot_coefficients(lr)



Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
- 7 Logistic Regression
 - 7.1 Building analysis
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimization
 - 10.1.1 Select Best Model
 - 10.1.1.1 K=3
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian Naive Bayes
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimization
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

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]])
```

```
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],
```

1 Top Negative Factors affecting Churn Rate

In [183]:

```

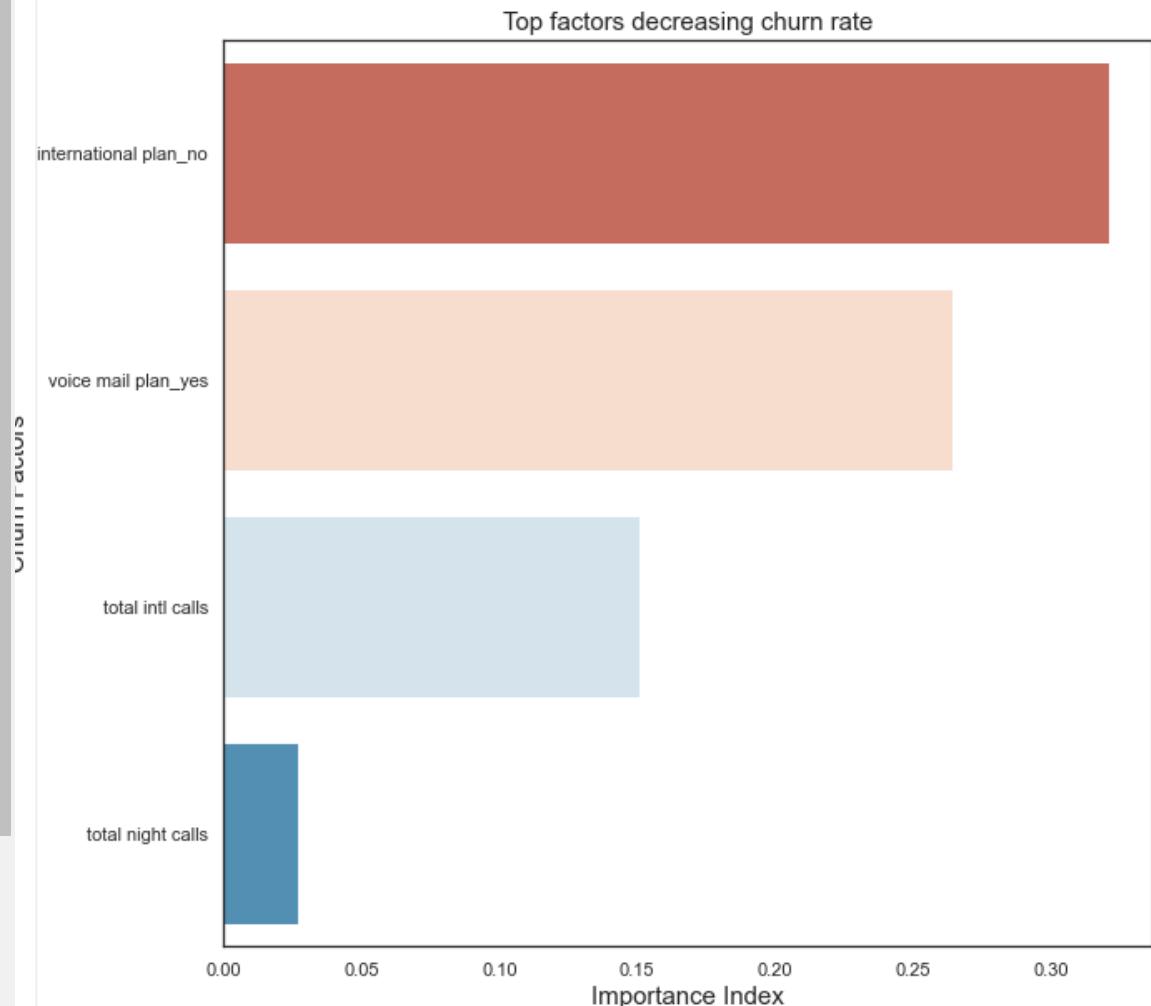
1 sns.set(rc={"figure.figsize":(10,10)})
2 sns.set_style("white")
3 sns.barplot(x=list(coefs_neg.values()), y=list(coefs_neg.keys()), palette="magma")
4 plt.title('Top factors decreasing churn rate', fontsize=15)
5 plt.xlabel("Importance Index", fontsize=15)
6 plt.ylabel('Churn Factors', fontsize=15)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

```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

Biggest Negative Factors on Churn



20.2 Top Positive Factors Affecting Churn Rate

In [184]:

```

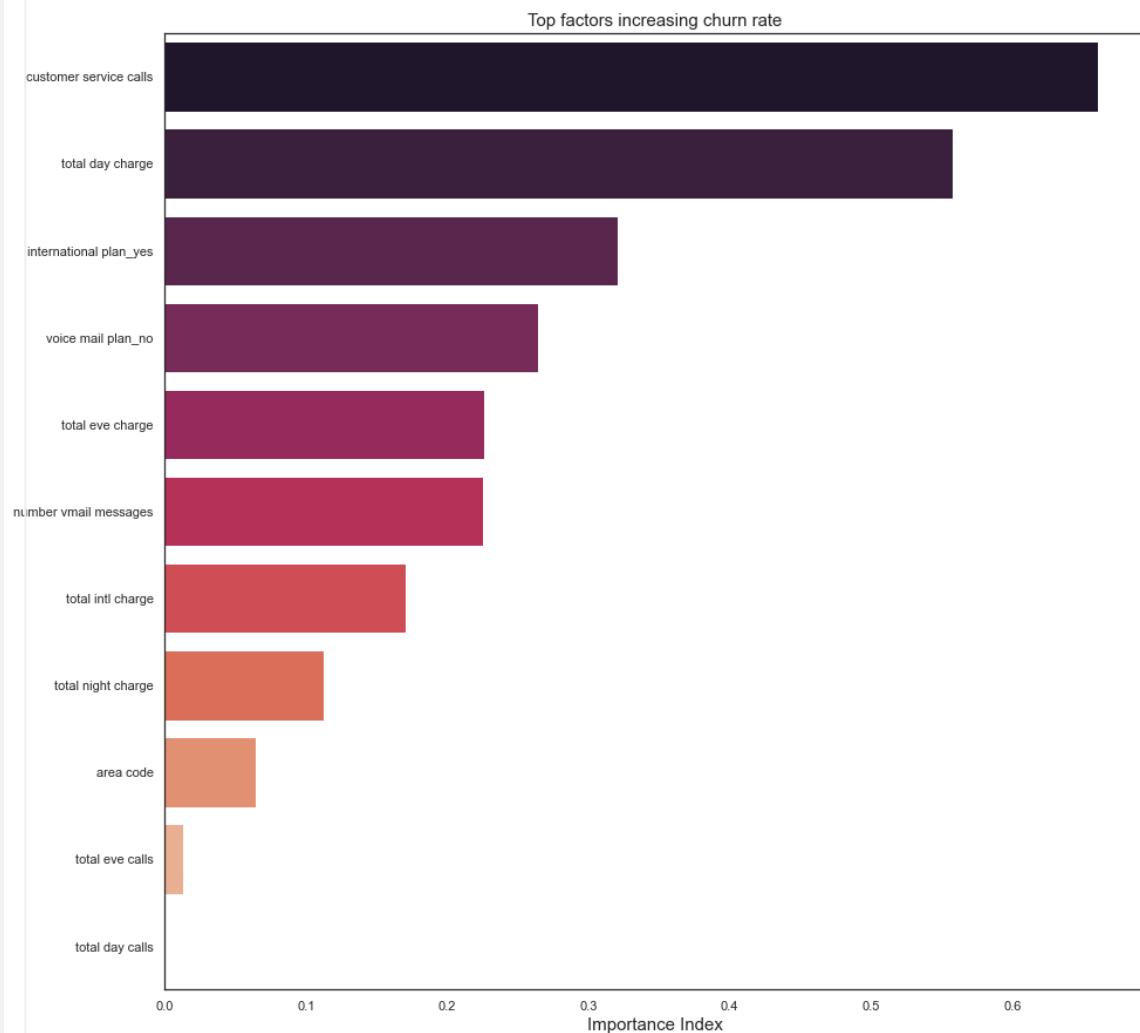
1 sns.set(rc={"figure.figsize":(15,15)})
2 sns.set_style("white")
3 sns.barplot(x=list(coefs_pos.values()), y=list(coefs_pos.keys()), palette="magma")
4 plt.title('Top factors increasing churn rate', fontsize=15)
5 plt.xlabel("Importance Index", fontsize=15)
6 plt.ylabel('Churn Factors', fontsize=15)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

```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comment
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performing model
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Model
 - 10.1.1 Select Ensemble Model
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Model
 - 15.1.1 Hyperparameter Tuning
 - 15.1.1.1 Best Model

Biggest Positive Factors on Churn rate



Conclusion

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.

Gaussian, Adaboost and Gradient Boost Classifier had the optimal recall and precision of

We decided to use DecisionTree classifier which performed well with 65% recall and 87% precision to interpret the data and extract the features.

Adaboost and GradientBoost, DecisionTree and RandoForest intrepreted the data similarly.

Factors that contribute to churn of customers are:

Contents

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state by zip code
 - 6.1.1 Common zip codes
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results
- 8 Descion Tree Classification
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Model
 - 8.2 Results Comparison
- 9 Random Forest Classification
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate on K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Model

Increasing churn factors

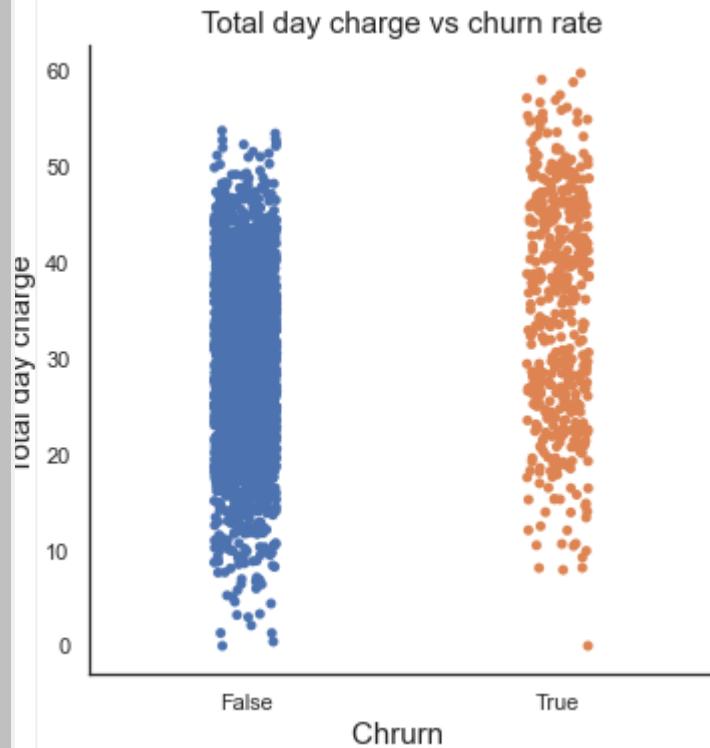
1 Total day charge

In [146]:

```
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 ⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Euclidean Distance
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over K
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



In [214]:

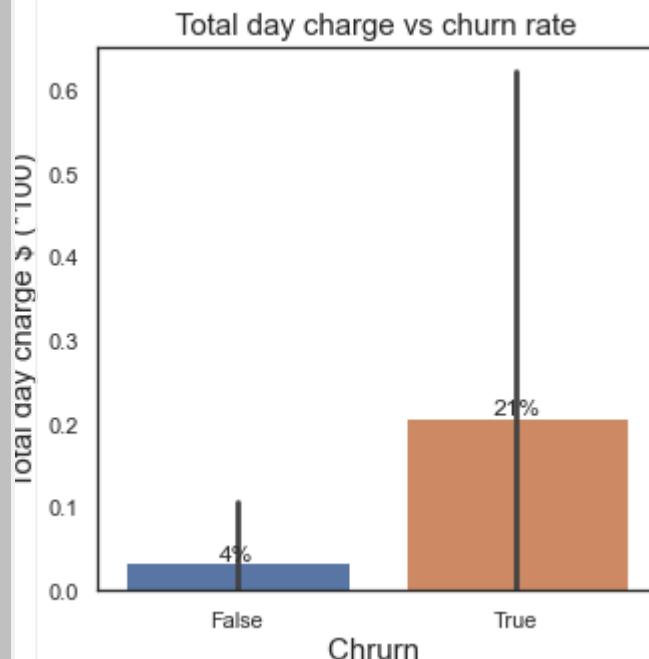
```

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("Churn", 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='ce

```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state column
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Estimator
 - 10.1.1.1 K=3
 - 10.1.2 Iterate cross-validation
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



In [219]:

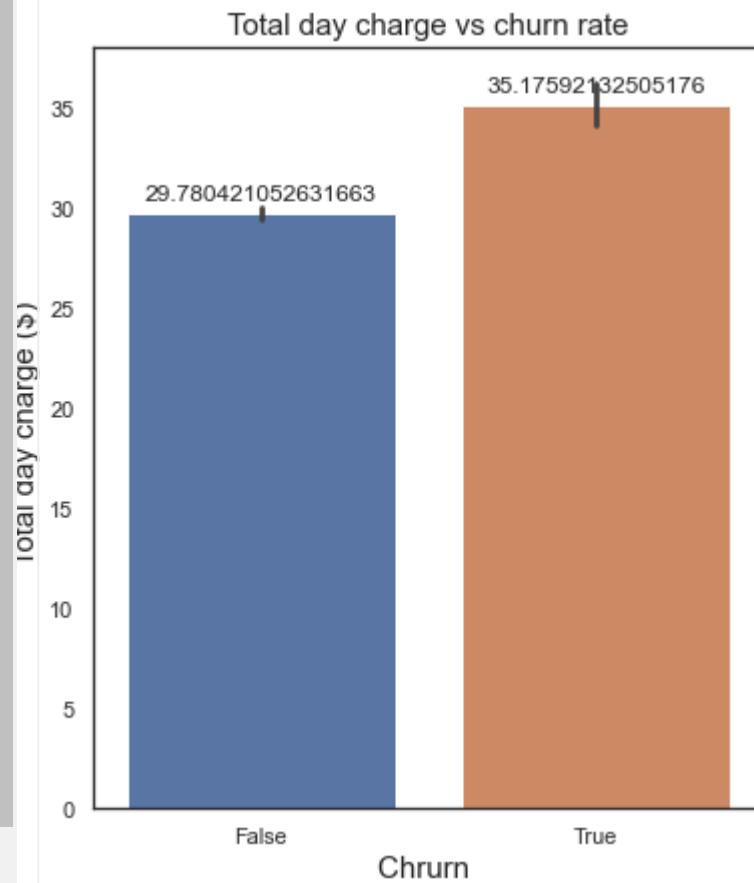
```

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')

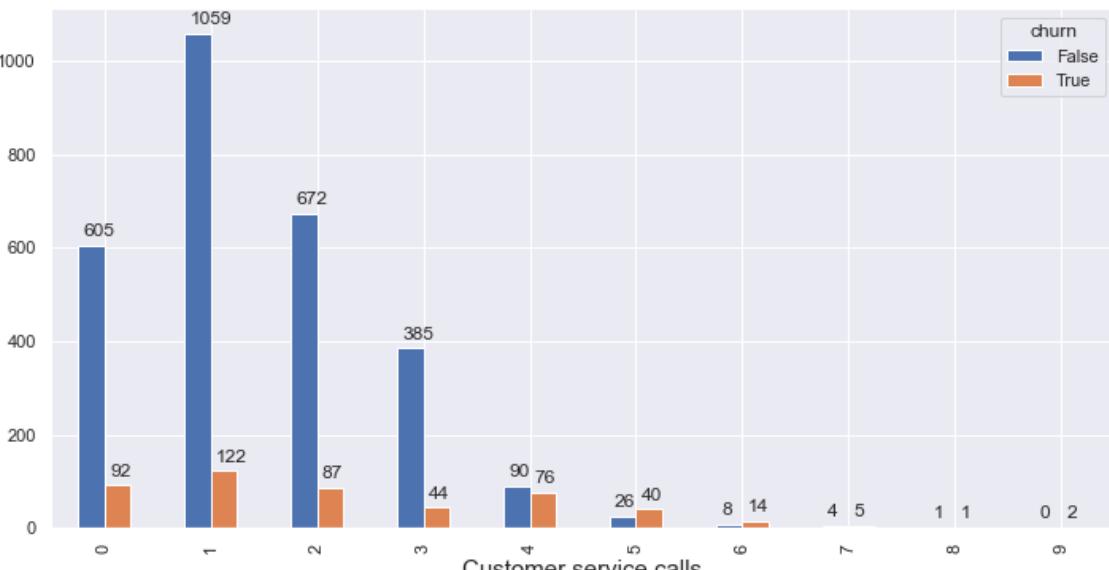
```

Contents ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over different K values
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters
 - 15.1.1.1 Best Parameters



Contents ↻⚙️

1 Abstract	1 cs=df.groupby(["customer service calls", "churn"]).size().unstack().plot
2 Business Problem	In [216]:
3 Import Libaries	for i in cs.patches:
4 Load Data	cs.text(i.get_x() + 0.05, i.get_height() + 20, int(i.get_height()))
5 Data Exploration	4
6 Data Engineering	5 plt.xlabel('Customer service calls', fontsize=15)
6.1 Group state	6 plt.ylabel('Churn number', fontsize=15)
6.1.1 Commen	7 plt.title('Customer service calls vs Churn rate', fontsize=15)
6.2 Drop unnece	
6.3 Drop higly co	ext(0.5, 1.0, 'Customer service calls vs Churn rate')
6.4 One Hot Enc	
6.5 Train test spl	Customer service calls vs Churn rate
6.6 Standarize F	
6.7 Building anal	
7 Logistic Regress	
7.1 Optimization	
7.2 Best perform	
7.3 Results Com	
8 Descion Tree Clas	
8.1 Hyperparamet	
8.1.1 Select Be	
8.2 Results Com	
9 Random Forest C	
9.1 Hyperparamet	
9.2 Best Paramet	
9.3 Results & Co	
10 KNN Classifier	
10.1 Tuning & Op	
10.1.1 Select E	
10.1.1.1 21	
10.1.2 Iterate c	
10.1.2.1 K=1	
10.2 Results & C	
11 XGBoost	
11.1 Results & co	
12 Gaussian NB	
12.1 Results & C	
13 AdaBoostClassi	
13.1 Results & C	
14 Gradient Boosti	
14.1 Results & co	
15 SVM Classifier	
15.1 Tuning & Op	
15.1.1 Hypertu	
15.1.1.1 Best	

In [220]:

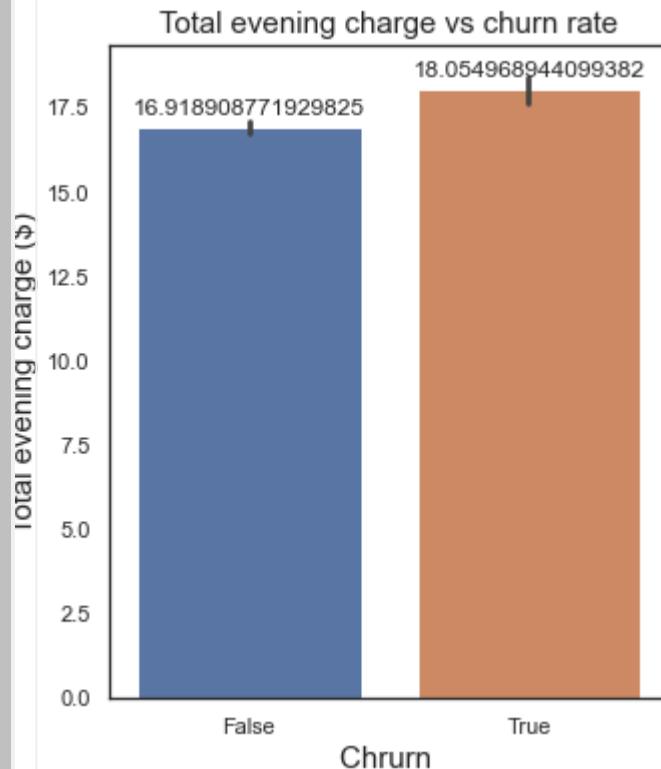
```

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 plt.title('Total evening charge vs churn rate', fontsize=15)
5 plt.xlabel("Churn ", fontsize=15)
6 plt.ylabel('Total evening 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 ⚙

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state of account
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Descion Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Best Parameters
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over K values
 - 10.1.2.1 K=1
 - 10.2 Results & Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameters



4 Total internatioanl charge

In [223]:

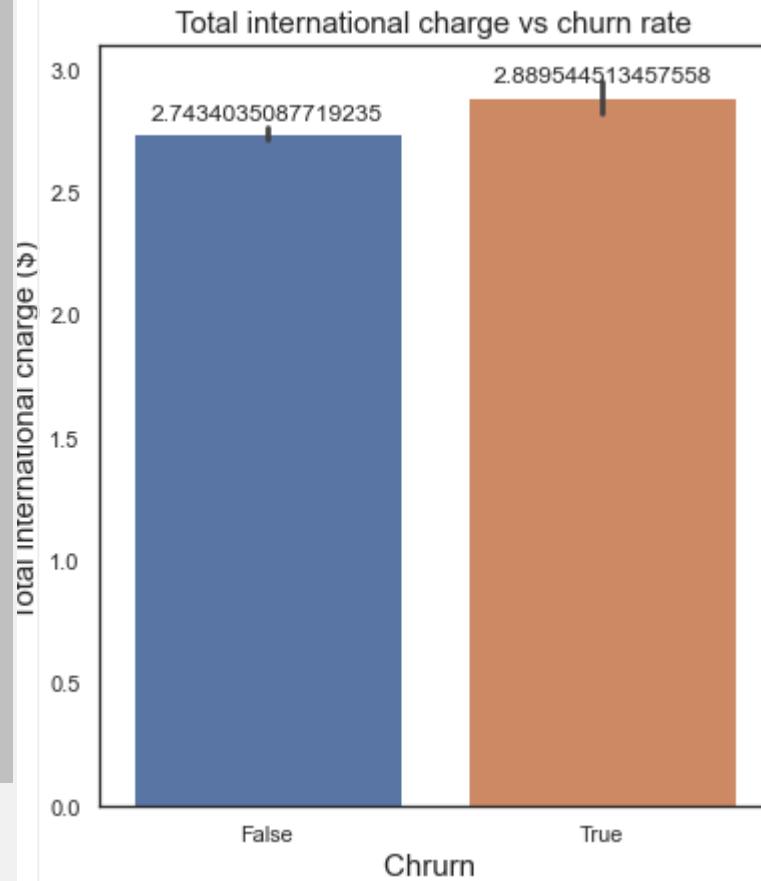
```

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 ⚙️

- 1 Abstract
- 2 Business Problem
- 3 Import Libraries
- 4 Load Data
- 5 Data Exploration
- 6 Data Engineering
 - 6.1 Group state columns
 - 6.1.1 Comments
 - 6.2 Drop unnecessary columns
 - 6.3 Drop highly correlated columns
 - 6.4 One Hot Encoding
 - 6.5 Train test split
 - 6.6 Standardize Features
 - 6.7 Building analysis
- 7 Logistic Regression
 - 7.1 Optimization
 - 7.2 Best performance
 - 7.3 Results Comparison
- 8 Decision Tree Classifier
 - 8.1 Hyperparameters
 - 8.1.1 Select Best Parameters
 - 8.2 Results Comparison
- 9 Random Forest Classifier
 - 9.1 Hyperparameters
 - 9.2 Best Parameters
 - 9.3 Results & Comparison
- 10 KNN Classifier
 - 10.1 Tuning & Optimal Parameters
 - 10.1.1 Select Error Function
 - 10.1.1.1 K=3
 - 10.1.2 Iterate over K values
 - 10.1.2.1 K=1
 - 10.2 Results Comparison
- 11 XGBoost
 - 11.1 Results & Comparison
- 12 Gaussian NB
 - 12.1 Results & Comparison
- 13 AdaBoostClassifier
 - 13.1 Results & Comparison
- 14 Gradient Boosting
 - 14.1 Results & Comparison
- 15 SVM Classifier
 - 15.1 Tuning & Optimal Parameters
 - 15.1.1 Hyperparameter Tuning
 - 15.1.1.1 Best Parameters

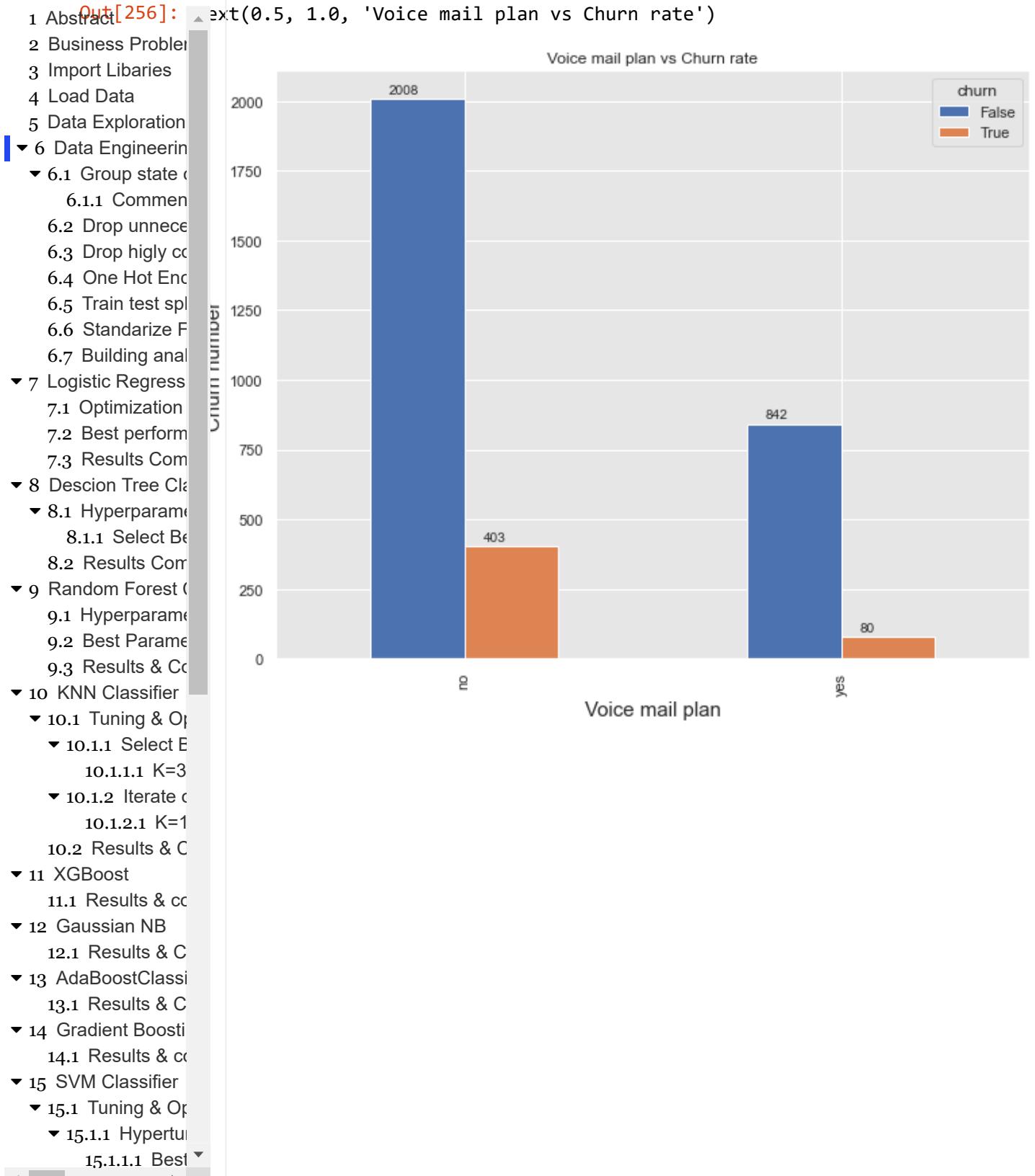


21 Decreasing churn factors

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 ⚙



In [228]: 1 df.columns

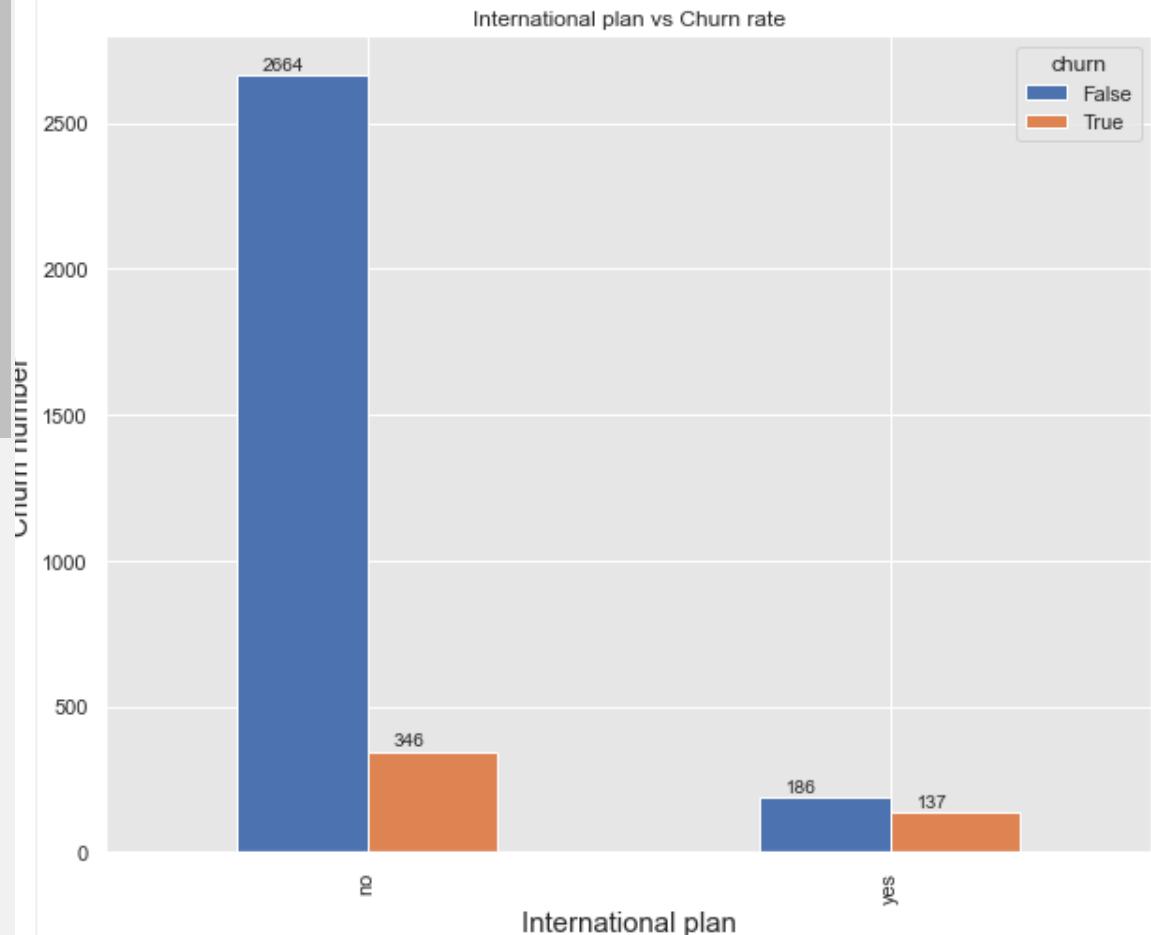
```
Out[228]: Index(['area code', 'number vmail messages', 'total day calls',
       'total day charge', 'total eve calls', 'total eve charge',
       'total night calls', 'total night charge', 'total intl calls',
       'total intl charge', 'customer service calls', 'churn',
       'international plan_no',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

Contents ↗⚙️ no'

1 Abstract
2 Business Problem
3 Import Libraries
4 Load Data
5 Data Exploration
24 6 Data Engineering

6.1 Group state
6.1.1 Comments
6.2 Drop unnecessary columns
6.3 Drop highly correlated columns
6.4 One Hot Encoding
6.5 Train test split
6.6 Standardize Features
6.7 Building analysis

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



24 Recomendations

., increase call charge for day and special charges for evening calls.

- 2) Have a better international plan.
- 3) Customer service calls is a measure usually for bad experience. Look what are the main causes of complaints and improve the quality of phone service.
- 4) Keep up the good work with the voice mail plan which is the best feature that keeps customers in the company.

Contents ↗

1 Abstract	5) C	Compete with other companies on this strongly recommended feature "voice mail plan" which
2 Business Problem	can	tract customers from competing rivals
3 Import Libraries		
4 Load Data		
5 Data Exploration		
6 Data Engineering	2)	Future Work
6.1 Group state of	More	Important data is needed to predict the churn of customers.
6.1.1 Communi		
6.2 Drop unneces		
6.3 Drop highly	1) Co	Competitor Information
6.4 One Hot Enc	2) I	Carrier package information
6.5 Train test spl		
6.6 Standardize	3) C	Contract Information
6.7 Building anal		
7 Logistic Regression	5) F	ment method
7.1 Optimization		
7.2 Best perform	6) C	plaint data
7.3 Results Com		
8 Descion Tree	7) C	ture data
8.1 Hyperparam		
8.1.1 Select Be	In [1]:	
8.2 Results Com		
9 Random Forest		
9.1 Hyperparam		
9.2 Best Paramet		
9.3 Results & Co		
10 KNN Classifier		
10.1 Tuning & Op		
10.1.1 Select E		
10.1.1.1 K=3		
10.1.2 Iterate c		
10.1.2.1 K=1		
10.2 Results & C		
11 XGBoost		
11.1 Results & co		
12 Gaussian NB		
12.1 Results & C		
13 AdaBoostClassi		
13.1 Results & C		
14 Gradient Boosti		
14.1 Results & co		
15 SVM Classifier		
15.1 Tuning & Op		
15.1.1 Hypertu		
15.1.1.1 Best		