

Business Understanding

Introduction

For telecommunications companies such as Syria Tel, the loss of paying subscribers is a major challenge. This leads to a reduction in revenue, increased costs associated with the acquisition of customers and diminished brand loyalty. SyriaTel needs a strong and datadriven approach to proactively tackle customer turnover and increase retention of customers.

Problem Statement

The goal of this project is to use SyriaTel's customer data to create a model for predicting customer churn. SyriaTel will be able to implement targeted retention strategies and reduce customer attrition by using the model to identify customers who are at high risk of churning.

Business Objectives

- Forecasting customer churn to lower client attrition.
- Focusing on high-risk consumers for retention efforts.
- Increasing customer lifetime value by focusing on significant revenue contributors.
- Enhancing client experience to increase satisfaction and loyalty.
- Using model insights for customized marketing campaigns to re-engage potential customers.

Data Description

The project will make use of a dataset from SyriaTel's customer relationship management (CRM) system that contains previous customer data. Features like user demographics, account information, call usage trends, and customer service interactions are anticipated to be included in this data.

Expected Benefits

- Better customer experience: By anticipating problems and knowing how customers behave, businesses can build loyalty and improve the overall customer experience.
- Data-driven decision making: Strategic customer retention initiatives and resource allocation are guided by insights.
- Customized marketing initiatives: Models help with campaign customization for particular clientele groups.
- Proactive customer engagement: By addressing customer complaints and preventing churn, early detection of churn risk helps to prevent it.

Data Understanding

Dataset Overview

The customer relationship management (CRM) system of SyriaTel provides historical customer data in the dataset. Customer demographics, account information, phone usage trends, and customer service exchanges are among the features.

Structure and Features

- **Features:**
 - State of residence
 - Account length
 - Service plans (e.g., international plan, voice mail plan)
 - Call usage metrics (e.g., total minutes, number of calls)
 - Customer service interactions

Potential Issues

- Missing values: Determine appropriate handling procedures after looking for any missing values in the dataset.
- Outliers: Spot any outliers in numerical features and assess whether further action is required.
- Data discrepancies: Check categorical characteristics for errors or inconsistencies that could impact analysis.

Data Preprocessing

To address missing numbers, outliers, and inconsistencies, clean up the dataset. carrying out required transformations, such as scaling numerical characteristics and encoding categorical variables.

Data Exploration

Overview

The first analysis of the SyriaTel customer churn dataset aimed to understand its properties and potential connections to customer attrition. Service plans, consumption patterns, and client demographics are all included in the dataset. Examining feature distribution, spotting anomalies, and figuring out possible connections between variables and the target variable were the key goals.

Importing Libraries

In [13]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import math
6 from scipy import stats
7
8 from sklearn.model_selection import train_test_split, cross_val_score
9 from sklearn.preprocessing import StandardScaler, OneHotEncoder
10 from sklearn.compose import ColumnTransformer
11 from sklearn.pipeline import Pipeline
12 from sklearn.feature_selection import SelectFromModel
13
14 from sklearn.ensemble import RandomForestClassifier
15 from xgboost import XGBClassifier
16 from sklearn.linear_model import LogisticRegression
17 from sklearn.metrics import accuracy_score, precision_score, recall_score,
18 from sklearn.metrics import accuracy_score, classification_report, confusi
19
20 import warnings
21 warnings.filterwarnings("ignore")
```

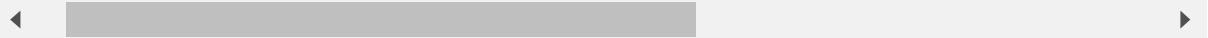
Loading CSV data

```
In [34]: 1 # Load the data from a CSV file
2 df = pd.read_csv("churn.csv")
3
4 # Display first 10 rows
5 df.head(10)
```

Out[34]:

state	account length	area code	phone number	international plan	voice mail plan	number vmail messages	total day minutes	total day calls	total day charge	...	total eve calls	c
KS	128	415	382-4657	no	yes	25	265.1	110	45.07	...	99	
OH	107	415	371-7191	no	yes	26	161.6	123	27.47	...	103	
NJ	137	415	358-1921	no	no	0	243.4	114	41.38	...	110	
OH	84	408	375-9999	yes	no	0	299.4	71	50.90	...	88	
OK	75	415	330-6626	yes	no	0	166.7	113	28.34	...	122	
AL	118	510	391-8027	yes	no	0	223.4	98	37.98	...	101	
MA	121	510	355-9993	no	yes	24	218.2	88	37.09	...	108	
MO	147	415	329-9001	yes	no	0	157.0	79	26.69	...	94	
LA	117	408	335-4719	no	no	0	184.5	97	31.37	...	80	
WV	141	415	330-8173	yes	yes	37	258.6	84	43.96	...	111	

ows × 21 columns



In [28]: 1 df.describe()

Out[28]:

y s	total day calls	total day charge	total eve minutes	total eve calls	total eve charge	total night minutes	total night calls	1
0	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	3333.000000	33
3	100.435644	30.562307	200.980348	100.114311	17.083540	200.872037	100.107711	
9	20.069084	9.259435	50.713844	19.922625	4.310668	50.573847	19.568609	
0	0.000000	0.000000	0.000000	0.000000	0.000000	23.200000	33.000000	
0	87.000000	24.430000	166.600000	87.000000	14.160000	167.000000	87.000000	
0	101.000000	30.500000	201.400000	100.000000	17.120000	201.200000	100.000000	
0	114.000000	36.790000	235.300000	114.000000	20.000000	235.300000	113.000000	
0	165.000000	59.640000	363.700000	170.000000	30.910000	395.000000	175.000000	

◀ ▶

In [17]: 1 df.info()

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

In []: 1

The `info()` method was used to load the dataset and examine its contents. It was found that each column included 3333 non-null entries. This suggests that the dataset is free of missing values.

```
In [20]: 1 df.isnull().sum()
```

```
Out[20]: state          0  
account length        0  
area code             0  
phone number          0  
international plan    0  
voice mail plan       0  
number vmail messages 0  
total day minutes     0  
total day calls       0  
total day charge      0  
total eve minutes     0  
total eve calls       0  
total eve charge      0  
total night minutes   0  
total night calls     0  
total night charge    0  
total intl minutes    0  
total intl calls      0  
total intl charge     0  
customer service calls 0  
churn                  0  
dtype: int64
```

The `df.isnull().sum()` method was used to further verify that there were no missing values, and it returned 0 to indicate this. We can confidently move forward with additional research and modeling now that there are no missing values.

Exploratory Data Analysis

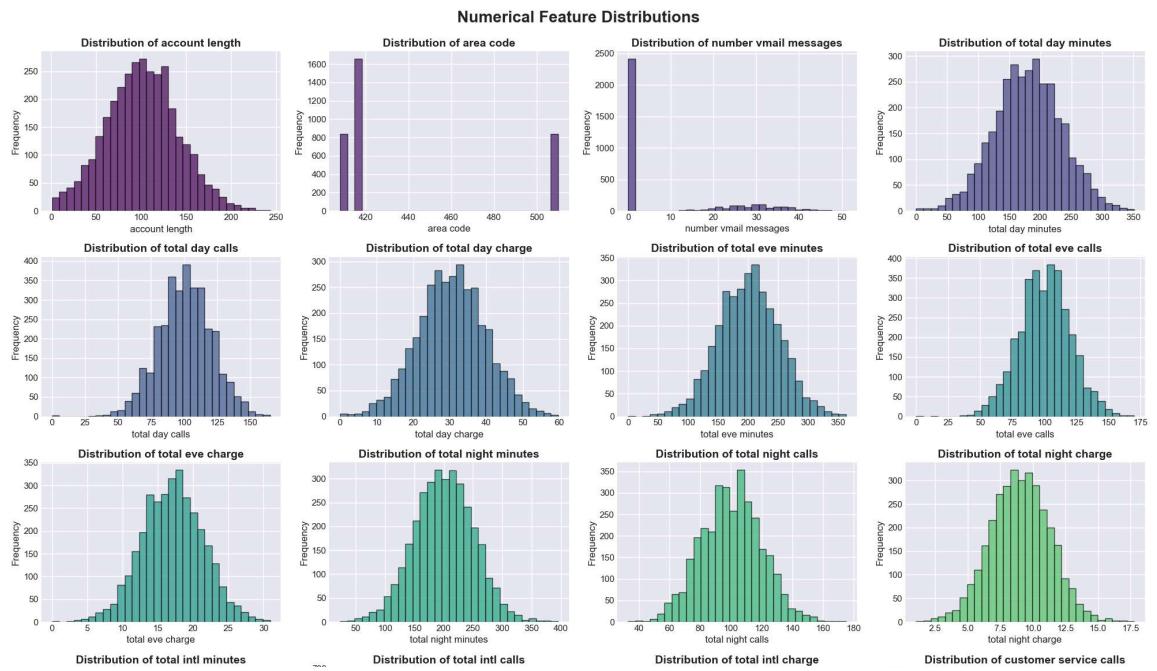
1. Univariate Analysis

Individual variables were analyzed to understand the distribution of the independent variables.

Visualisations

In [35]:

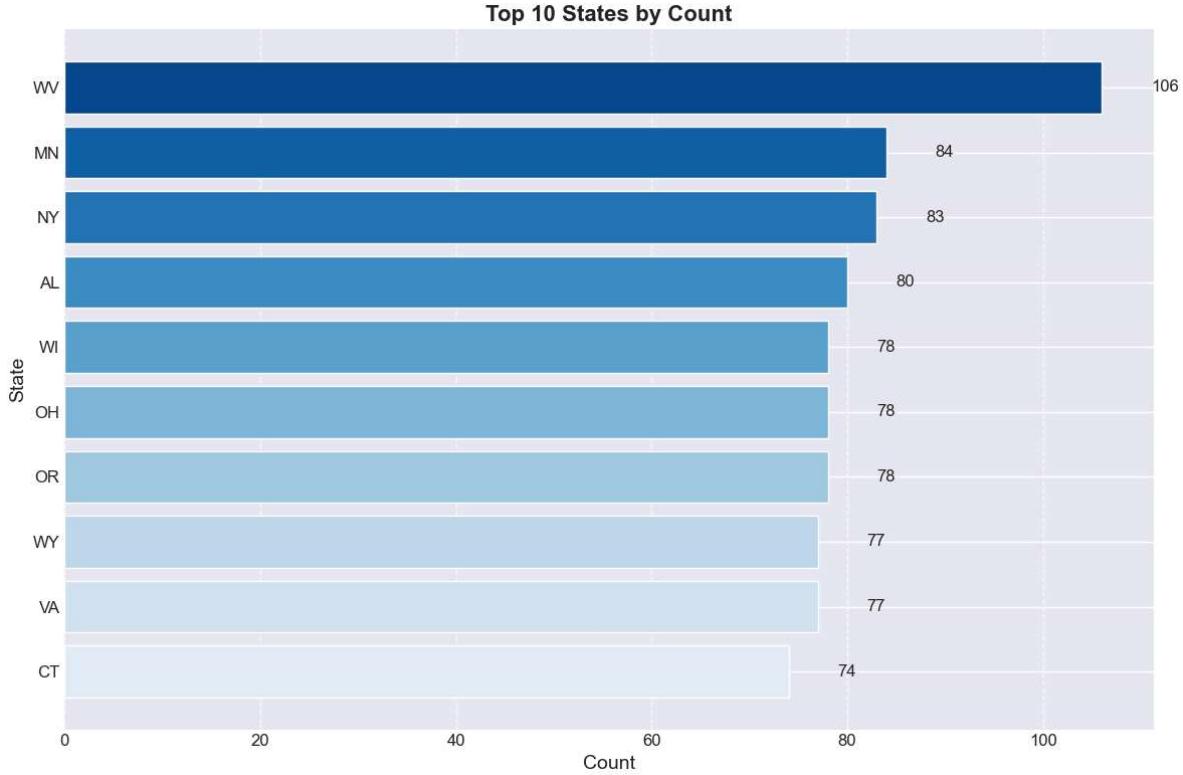
```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 import numpy as np
4 import math
5
6 # Load the data
7 df = pd.read_csv("churn.csv")
8
9 # Select only numerical features
10 numerical_features = df.select_dtypes(include=[np.number]).columns.tolist()
11
12 # Calculate the number of rows and columns needed for subplots
13 num_features = len(numerical_features)
14 num_cols = 4
15 num_rows = math.ceil(num_features / num_cols)
16
17 # Set the style
18 plt.style.use('seaborn-darkgrid')
19
20 # Create a color palette
21 colors = plt.cm.viridis(np.linspace(0, 1, num_features))
22
23 # Visualize numerical features
24 fig, axes = plt.subplots(num_rows, num_cols, figsize=(20, 15))
25 axes = axes.flatten() # Flatten the axes array to make indexing easier
26
27 for i, (column, color) in enumerate(zip(numerical_features, colors)):
28     ax = axes[i]
29     ax.hist(df[column].dropna(), bins=30, color=color, alpha=0.7, edgecolor='black')
30     ax.set_title(f'Distribution of {column}', fontsize=14, fontweight='bold')
31     ax.set_xlabel(column, fontsize=12)
32     ax.set_ylabel('Frequency', fontsize=12)
33     ax.grid(True)
34
35 # Remove empty subplots
36 for j in range(i + 1, num_rows * num_cols):
37     fig.delaxes(axes[j])
38
39 # Adjust layout
40 plt.tight_layout(rect=[0, 0, 1, 0.96])
41 plt.suptitle('Numerical Feature Distributions', fontsize=20, fontweight='bold')
42 plt.show()
43
```



The numerical data showed a right-skewed pattern in customer tenure and call durations across different time periods, suggesting lower usage patterns. International calls may have higher volumes for a smaller customer segment. Customer service calls showed a multimodal distribution, indicating distinct groups with varying frequency of contacting customer service. Understanding these relationships and potential outliers can help predict churn.

In [43]:

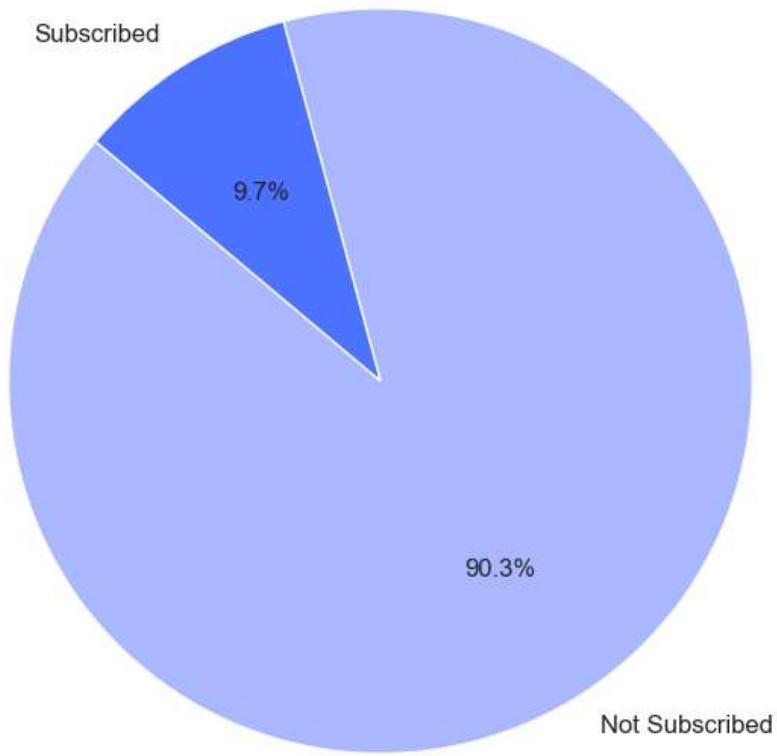
```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Get the top ten states by count
5 top_states = df['state'].value_counts().sort_values(ascending=False).head(10)
6
7 # Plotting
8 plt.figure(figsize=(12, 8))
9 bars = plt.barh(top_states.index, top_states.values, color=sns.color_palette('viridis'))
10
11 # Adding annotations
12 for bar, count in zip(bars, top_states.values):
13     plt.text(bar.get_width() + 5, bar.get_y() + bar.get_height()/2, f'{count}', color='white', va='center', ha='left', fontsize=12)
14
15 # Customizing plot aesthetics
16 plt.title('Top 10 States by Count', fontsize=16, fontweight='bold')
17 plt.xlabel('Count', fontsize=14)
18 plt.ylabel('State', fontsize=14)
19 plt.xticks(fontsize=12)
20 plt.yticks(fontsize=12)
21 plt.grid(axis='x', linestyle='--', alpha=0.7)
22 plt.gca().invert_yaxis() # Invert y-axis to have the highest count at the bottom
23 plt.tight_layout()
24
25
26 plt.show()
27
```



Customers primarily from West Virginia, Minnesota, and New York make up SyriaTel's customer base.

```
In [56]: 1 import matplotlib.pyplot as plt
2
3 # Define custom colors
4 light_blue = '#4c74ff'
5 lighter_blue = '#aabaff'
6
7 # Calculate the percentage of each category
8 international_plan_counts = df['international plan'].value_counts()
9 labels = ['Subscribed' if x == 'yes' else 'Not Subscribed' for x in international_plan_counts.index]
10 sizes = international_plan_counts.values
11 colors = [lighter_blue, light_blue]
12
13 # Plotting
14 plt.figure(figsize=(8, 6))
15 plt.pie(sizes, labels=labels, colors=colors, autopct='%1.1f%%', startangle=90)
16 plt.title('Distribution of International Plan', fontsize=16, fontweight='bold')
17 plt.axis('equal') # Equal aspect ratio ensures that pie is drawn as a circle
18 plt.tight_layout()
19
20 plt.show()
21
```

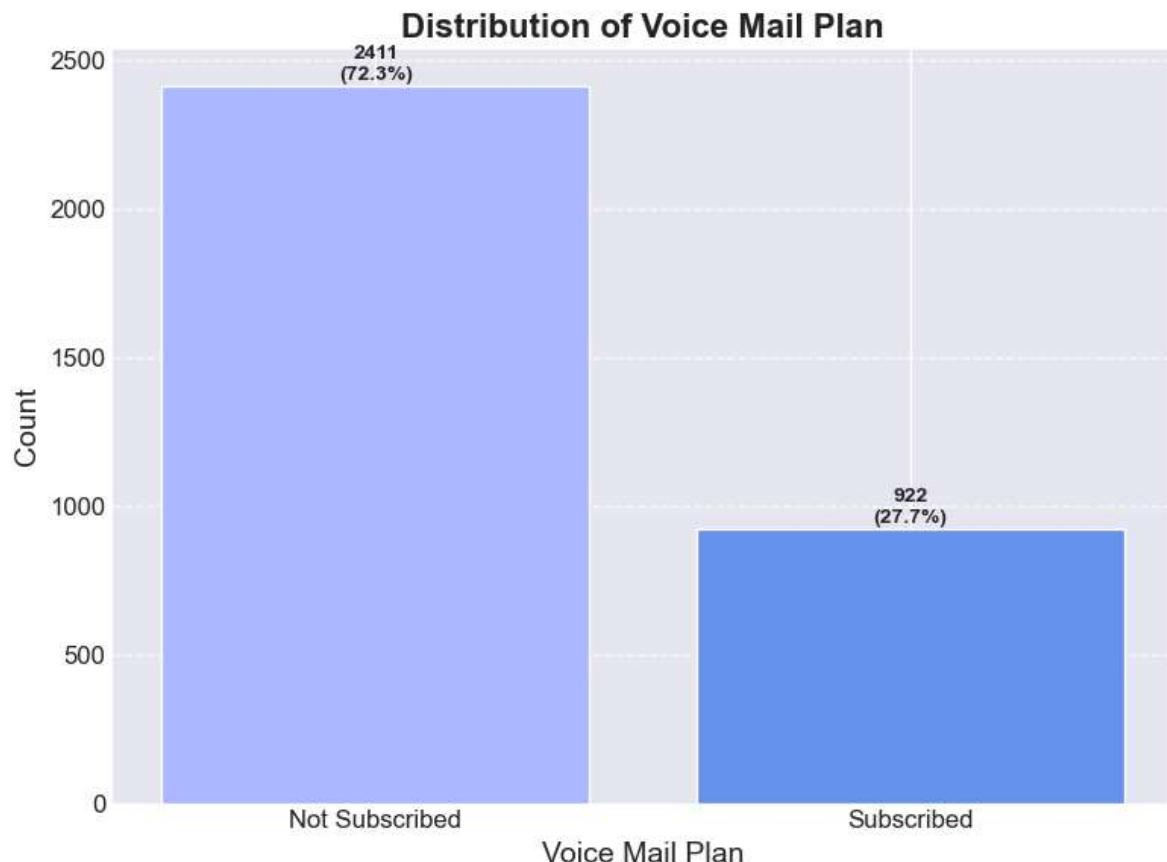
Distribution of International Plan



only 9.7% are subscribed to the international plan.

In [55]:

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 # Replace 'yes' and 'no' with 'Subscribed' and 'Not Subscribed'
5 df['voice mail plan'] = df['voice mail plan'].replace({'yes': 'Subscribed',
6
7     # Calculate counts for each category
8     counts = df['voice mail plan'].value_counts()
9
10    # Plotting a bar plot
11    plt.figure(figsize=(8, 6))
12    bars = plt.bar(counts.index, counts, color=[lighter_blue, light_blue])
13
14    # Adding annotations
15    for bar, count in zip(bars, counts):
16        plt.text(bar.get_x() + bar.get_width()/2, bar.get_height() + 5, f'{count} ({count / sum(counts) * 100:.1f}%)', ha='center', va='bottom', fontsize=10, fontweight='bold')
17
18    # Customizing plot aesthetics
19    plt.title('Distribution of Voice Mail Plan', fontsize=16, fontweight='bold')
20    plt.xlabel('Voice Mail Plan', fontsize=14)
21    plt.ylabel('Count', fontsize=14)
22    plt.xticks(fontsize=12)
23    plt.yticks(fontsize=12)
24    plt.grid(axis='y', linestyle='--', alpha=0.7)
25    plt.tight_layout()
26
27
28    plt.show()
29
```



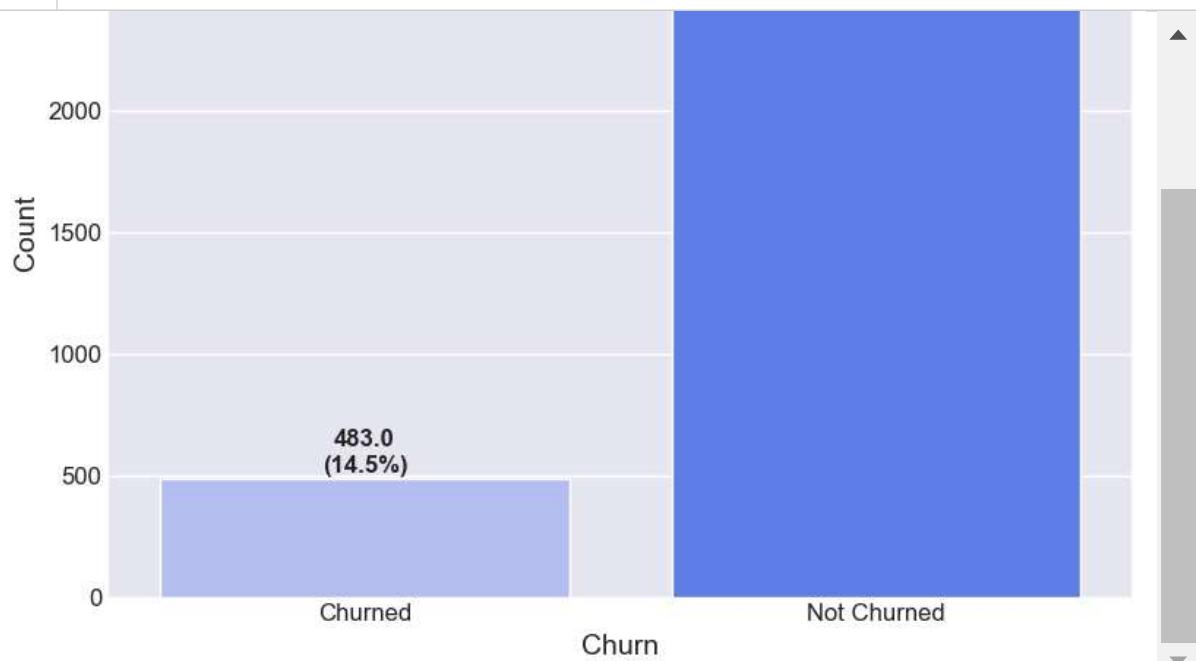
only 27.7% of customers have subscribed to voice mail service

Distribution of Churn

The 'churn' variable delineates whether a customer has churned or not. The endeavor encompassed a univariate analysis to elucidate the distribution of churn within the dataset.

In [61]:

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Define custom colors
5 light_blue = '#4c74ff'
6 lighter_blue = '#aabaff'
7
8 # Replace 'True' and 'False' with 'Churned' and 'Not Churned'
9 df['churn'] = df['churn'].replace({True: 'Churned', False: 'Not Churned'})
10
11 # Calculate churn percentages
12 churn_counts = df['churn'].value_counts()
13 churn_percentages = churn_counts / churn_counts.sum() * 100
14
15 # Plotting
16 plt.figure(figsize=(8, 6))
17 ax = sns.countplot(data=df, x='churn', palette=[lighter_blue, light_blue],
18
19 # Adding annotations
20 total = len(df['churn']) # Total number of records
21 for p in ax.patches:
22     height = p.get_height()
23     ax.text(p.get_x() + p.get_width() / 2., height + 10, f'{height}\n{height:.2f}%', ha='center', va='bottom', fontsize=12, fontweight='bold')
24
25 # Title and Labels
26 plt.title('Distribution of Churn', fontsize=16, fontweight='bold')
27 plt.xlabel('Churn', fontsize=14)
28 plt.ylabel('Count', fontsize=14)
29 plt.xticks(fontsize=12)
30 plt.yticks(fontsize=12)
31
32 # Show plot
33 plt.tight_layout()
34 plt.show()
```



A detailed analysis of churn was conducted by counting the number of consumers who stopped using the service and the number of customers who continued. The following are the outcomes:

Total number of churned customers: 483 There were 2850 clients that remained loyal.

Multivariate analysis

This entailed looking at several variables' relationships at once. This study looked at the relationship between several features and the aim variable (customer churn) when taken into account collectively.

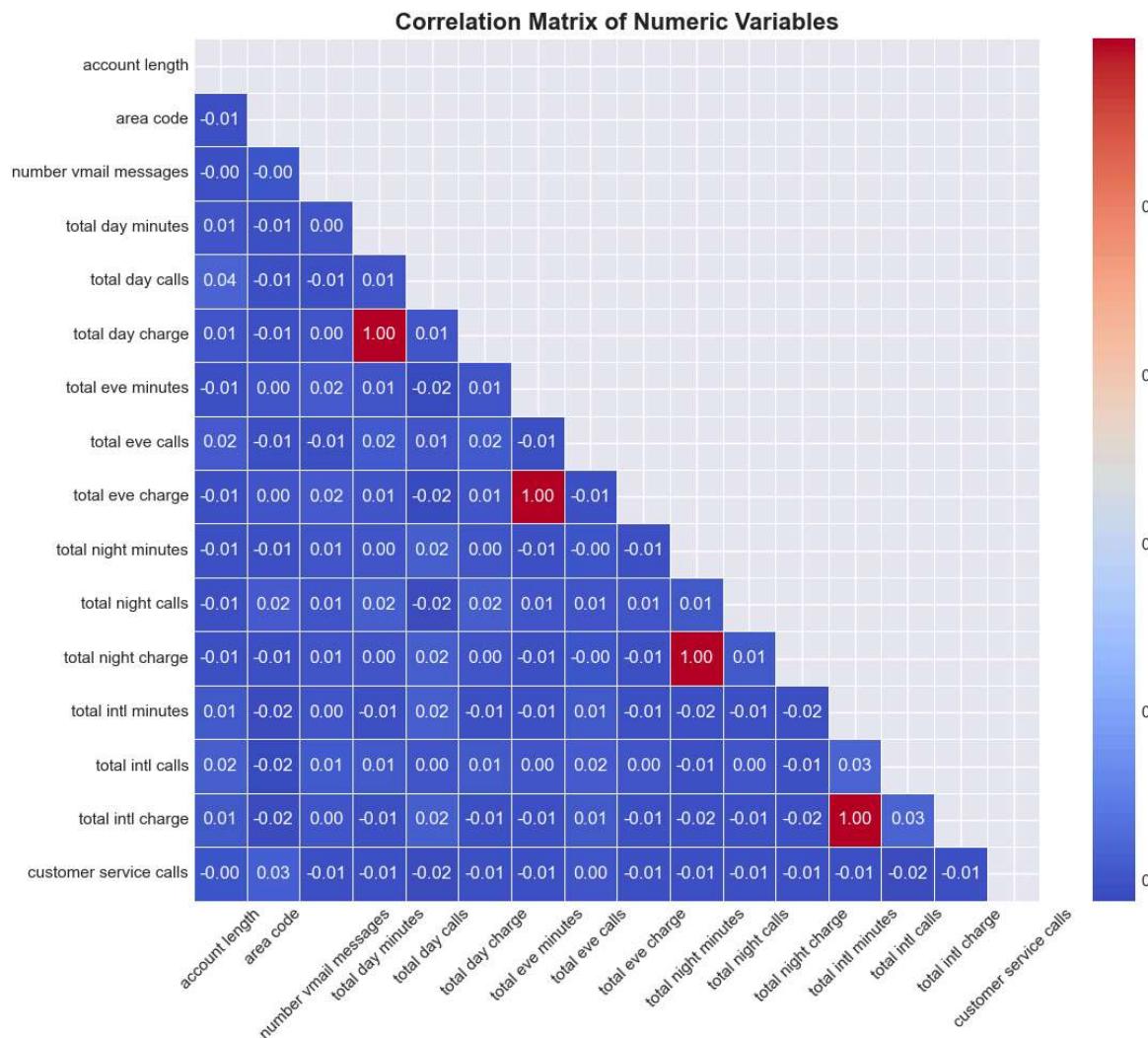
The correlations between the various variables in the dataset were found using a correlation matrix.

In [101]:

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 # Step 1: Select only numeric columns
7 numeric_df = df.select_dtypes(include=['float64', 'int64'])
8
9 # Step 2: Visualize correlation matrix
10 correlation_matrix = numeric_df.corr()
11
12 # Zero out lower triangular part
13 mask = np.triu(np.ones_like(correlation_matrix, dtype=bool), k=0)
14
15 plt.figure(figsize=(12, 10))
16 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=.5)
17 plt.title('Correlation Matrix of Numeric Variables', fontsize=16, fontweight='bold')
18 plt.xticks(rotation=45)
19 plt.yticks(rotation=0)
20 plt.tight_layout()
21 plt.show()
22

```



Correlation Matrix Analysis: Call Numbers, Minutes, Charges

- Significant positive relationships between call numbers, minutes, and charges.
- Higher call spending leads to more calls.
- Weak positive link between call minutes and account length.
- Customer service calls show weak correlations, unrelated to calling habits or account tenure.
- International calls show favorable link, leading to increased spending.

Data Preprocessing

Label Encoding converts label variables in "international plan", "voice mail plan", and "churn" columns to numeric form. The "Yes" and "No" categories are converted to 1 and 0, respectively, representing presence or absence. In the "churn" column, "False" and "True" categories are converted to 0 and 1, respectively, representing customer churn.

```
In [104]:  
1 from sklearn.preprocessing import LabelEncoder  
2  
3 # Define categorical columns to be label encoded  
4 cat_cols = ["international plan", "voice mail plan"]  
5  
6 # Apply label encoding to each categorical column  
7 df[cat_cols] = df[cat_cols].apply(LabelEncoder().fit_transform)  
8  
9 # Display the first ten rows to verify the changes  
10 print(df.head(10))
```

	state	account length	area code	phone number	international plan	\
0	KS	128	415	382-4657	0	
1	OH	107	415	371-7191	0	
2	NJ	137	415	358-1921	0	
3	OH	84	408	375-9999	1	
4	OK	75	415	330-6626	1	
5	AL	118	510	391-8027	1	
6	MA	121	510	355-9993	0	
7	MO	147	415	329-9001	1	
8	LA	117	408	335-4719	0	
9	WV	141	415	330-8173	1	

	voice mail plan	number vmail messages	total day minutes	total day calls	\
0	1	25	265.1	110	
1	1	26	161.6	123	
2	0	0	243.4	114	
3	0	0	299.4	71	
4	0	0	166.7	113	
5	0	0	223.4	98	
6	1	24	218.2	88	
7	0	0	157.0	79	
8	0	0	184.5	97	
9	1	37	258.6	84	

	total day charge	...	total eve calls	total eve charge	\
0	45.07	...	99	16.78	
1	27.47	...	103	16.62	
2	41.38	...	110	10.30	
3	50.90	...	88	5.26	
4	28.34	...	122	12.61	
5	37.98	...	101	18.75	
6	37.09	...	108	29.62	
7	26.69	...	94	8.76	
8	31.37	...	80	29.89	
9	43.96	...	111	18.87	

	total night minutes	total night calls	total night charge	\
0	244.7	91	11.01	
1	254.4	103	11.45	
2	162.6	104	7.32	
3	196.9	89	8.86	
4	186.9	121	8.41	
5	203.9	118	9.18	
6	212.6	118	9.57	
7	211.8	96	9.53	
8	215.8	90	9.71	
9	326.4	97	14.69	

	total intl minutes	total intl calls	total intl charge	\
0	10.0	3	2.70	
1	13.7	3	3.70	
2	12.2	5	3.29	
3	6.6	7	1.78	
4	10.1	3	2.73	
5	6.3	6	1.70	
6	7.5	7	2.03	

7	7.1	6	1.92
8	8.7	4	2.35
9	11.2	5	3.02

```
customer service calls  churn
0                      1  False
1                      1  False
2                      0  False
3                      2  False
4                      3  False
5                      0  False
6                      3  False
7                      0  False
8                      1  False
9                      0  False
```

[10 rows x 21 columns]

One-hot encoding is a method that converts categorical variables into numerical format for machine learning algorithms, especially useful for dealing with variables with multiple categories.

In [105]:

```

1 from sklearn.preprocessing import OneHotEncoder
2
3 # Create an instance of the OneHotEncoder
4 encoder = OneHotEncoder(sparse=False, dtype=np.int64)
5
6 # Encode the "state" column
7 encoded_state = encoder.fit_transform(df[["state"]])
8
9 # Create a DataFrame with the encoded state columns
10 dummy_df_state = pd.DataFrame(encoded_state, columns=encoder.get_feature_names_out())
11
12 # Concatenate the encoded state columns with the original DataFrame
13 one_df = pd.concat([df.drop(columns=["state"]), dummy_df_state], axis=1)
14
15 one_df.head(10)
16

```

Out[105]:

international plan	voice mail plan	number vmail messages	total day minutes	total day calls	total day charge	total eve minutes	...	state_SD	state_TN	state_TX	s
0	1	25	265.1	110	45.07	197.4	...	0	0	0	0
0	1	26	161.6	123	27.47	195.5	...	0	0	0	0
0	0	0	243.4	114	41.38	121.2	...	0	0	0	0
1	0	0	299.4	71	50.90	61.9	...	0	0	0	0
1	0	0	166.7	113	28.34	148.3	...	0	0	0	0
1	0	0	223.4	98	37.98	220.6	...	0	0	0	0
0	1	24	218.2	88	37.09	348.5	...	0	0	0	0
1	0	0	157.0	79	26.69	103.1	...	0	0	0	0
0	0	0	184.5	97	31.37	351.6	...	0	0	0	0
1	1	37	258.6	84	43.96	222.0	...	0	0	0	0



In [107]:

```

1 # Verify the changes
2 print(one_df.head(10))

```

	account	length	area code	phone number	international plan	\
0		128	415	382-4657		0
1		107	415	371-7191		0
2		137	415	358-1921		0
3		84	408	375-9999		1
4		75	415	330-6626		1
5		118	510	391-8027		1
6		121	510	355-9993		0
7		147	415	329-9001		1
8		117	408	335-4719		0
9		141	415	330-8173		1

	voice	mail	plan	number	vmail	messages	total day	minutes	total day	calls
\										
0			1			25		265.1		110
1			1			26		161.6		123
2			0			0		243.4		114
3			0			0		299.4		71
4			0			0		166.7		113
5			0			0		223.4		98
6			1			24		218.2		88
7			0			0		157.0		79
8			0			0		184.5		97
9			1			37		258.6		84

	total day	charge	total eve	minutes	...	state_SD	state_TN	state_TX	\
0		45.07		197.4	...	0	0	0	
1		27.47		195.5	...	0	0	0	
2		41.38		121.2	...	0	0	0	
3		50.90		61.9	...	0	0	0	
4		28.34		148.3	...	0	0	0	
5		37.98		220.6	...	0	0	0	
6		37.09		348.5	...	0	0	0	
7		26.69		103.1	...	0	0	0	
8		31.37		351.6	...	0	0	0	
9		43.96		222.0	...	0	0	0	

	state_UT	state_VA	state_VT	state_WA	state_WI	state_WV	state_WY
0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0
9	0	0	0	0	0	0	1

[10 rows x 71 columns]

```
In [114]: 1 from sklearn.preprocessing import OneHotEncoder
2 import pandas as pd
3
4 # Example DataFrame
5 df = pd.DataFrame({
6     'area code': [408, 415, 408, 510],
7     'other_column': [1, 2, 3, 4]
8 })
9
10 # Check if "area code" column exists in the DataFrame
11 if "area code" in df.columns:
12     # One-hot encoding for "area code" column and updating the DataFrame
13     encoder = OneHotEncoder(sparse_output=False)
14     encoded = encoder.fit_transform(df[['area code']])
15     df = df.join(pd.DataFrame(encoded, columns=encoder.get_feature_names_out()))
16
17 df.head(10)
18
```

Out[114]:

	other_column	area code_408	area code_415	area code_510
0	1	1.0	0.0	0.0
1	2	0.0	1.0	0.0
2	3	1.0	0.0	0.0
3	4	0.0	0.0	1.0

Scaling numerical features is important for algorithms that are sensitive to feature scales, such as gradient descent-based algorithms and distance-based algorithms.

In [115]:

1 newdf = one_df

In [121]:

1 newdf.dtypes

```
Out[121]: account length      int64
area code           int64
international plan    int64
voice mail plan      int64
number vmail messages int64
...
state_VT            int64
state_WA            int64
state_WI            int64
state_WV            int64
state_WY            int64
Length: 70, dtype: object
```

```
In [123]:  
1 from sklearn.preprocessing import MinMaxScaler  
2 import pandas as pd  
3  
4 # Define numeric columns  
5 numeric_columns = newdf.select_dtypes(include=['number']).columns  
6  
7 # Drop rows with NaN values in numeric columns  
8 newdf.dropna(subset=numeric_columns, inplace=True)  
9  
10 # Check if there are numeric columns to scale  
11 if numeric_columns.empty:  
12     print("No numeric columns found")  
13 else:  
14     # Clean and scale numeric columns  
15     scaler = MinMaxScaler()  
16     newdf[numeric_columns] = scaler.fit_transform(newdf[numeric_columns]).  
17  
18     # Print the first few rows of the scaled DataFrame  
19     print(newdf[numeric_columns].head(10))
```

	account length	area code	international plan	voice mail plan	\		
0	0.524793	0.068627	0.0	1.0			
1	0.438017	0.068627	0.0	1.0			
2	0.561983	0.068627	0.0	0.0			
3	0.342975	0.000000	1.0	0.0			
4	0.305785	0.068627	1.0	0.0			
5	0.483471	1.000000	1.0	0.0			
6	0.495868	1.000000	0.0	1.0			
7	0.603306	0.068627	1.0	0.0			
8	0.479339	0.000000	0.0	0.0			
9	0.578512	0.068627	1.0	1.0			
	number vmail messages	total day minutes	total day calls		\		
0	0.490196	0.755701	0.666667				
1	0.509804	0.460661	0.745455				
2	0.000000	0.693843	0.690909				
3	0.000000	0.853478	0.430303				
4	0.000000	0.475200	0.684848				
5	0.000000	0.636830	0.593939				
6	0.470588	0.622007	0.533333				
7	0.000000	0.447548	0.478788				
8	0.000000	0.525941	0.587879				
9	0.725490	0.737172	0.509091				
	total day charge	total eve minutes	total eve calls	...	state_SD \		
0	0.755701	0.542755	0.582353	...	0.0		
1	0.460597	0.537531	0.605882	...	0.0		
2	0.693830	0.333242	0.647059	...	0.0		
3	0.853454	0.170195	0.517647	...	0.0		
4	0.475184	0.407754	0.717647	...	0.0		
5	0.636821	0.606544	0.594118	...	0.0		
6	0.621898	0.958207	0.635294	...	0.0		
7	0.447518	0.283475	0.552941	...	0.0		
8	0.525989	0.966731	0.470588	...	0.0		
9	0.737089	0.610393	0.652941	...	0.0		
	state_TN	state_TX	state_UT	state_VA	state_VT	state_WA	state_WI \
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0
5	0.0	0.0	0.0	0.0	0.0	0.0	0.0
6	0.0	0.0	0.0	0.0	0.0	0.0	0.0
7	0.0	0.0	0.0	0.0	0.0	0.0	0.0
8	0.0	0.0	0.0	0.0	0.0	0.0	0.0
9	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	state_WV	state_WY					
0	0.0	0.0					
1	0.0	0.0					
2	0.0	0.0					
3	0.0	0.0					
4	0.0	0.0					
5	0.0	0.0					
6	0.0	0.0					
7	0.0	0.0					

```
8      0.0      0.0
9      1.0      0.0
```

[10 rows x 69 columns]

splitting the data to guarantee an objective assessment of the models ie: the data must be divided into training and testing sets before modeling can begin.

In [125]:

```
1 from sklearn.model_selection import train_test_split
2
3 # Define the features (X) and target variable (y)
4 X = newdf.drop(columns=['churn']) # Features
5 y = newdf['churn'] # Target variable
6
7 # Split the data into training and testing sets
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, r
9
10 # Print the shapes of the training and testing sets
11 print("X_train shape:", X_train.shape)
12 print("y_train shape:", y_train.shape)
13 print("X_test shape:", X_test.shape)
14 print("y_test shape:", y_test.shape)
```

```
X_train shape: (2666, 69)
y_train shape: (2666,)
X_test shape: (667, 69)
y_test shape: (667,)
```

class imbalance can be averted by the use of random oversampling through the resample function addresses class imbalance, ensuring a fairer representation of data and potentially enhancing the model's performance in handling imbalanced datasets

```
In [134]: 1 from sklearn.utils import resample
2 import numpy as np
3
4 # Combine X_train and y_train to resample
5 train_data = np.column_stack((X_train, y_train))
6
7 # Separate the majority and minority classes
8 majority_class = train_data[train_data[:, -1] == np.bincount(y_train).argmax()]
9 minority_class = train_data[train_data[:, -1] == np.bincount(y_train).argmin()]
10
11 # Perform random oversampling on the minority class
12 minority_class_upsampled = resample(minority_class,
13                                     replace=True, # Sample with replacement
14                                     n_samples=majority_class.shape[0], #
15                                     random_state=42) # Reproducible results
16
17 # Combine majority class with upsampled minority class
18 resampled_data = np.vstack((majority_class, minority_class_upsampled))
19
20 # Separate X and y
21 X_train_resampled = resampled_data[:, :-1]
22 y_train_resampled = resampled_data[:, -1].astype(int) # Convert y to integers
23
24 # Verify the result
25 print("Original class distribution:", np.bincount(y_train))
26 print("Resampled class distribution:", np.bincount(y_train_resampled))
27
```

Original class distribution: [2284 382]
 Resampled class distribution: [2284 2284]

Modelling

Introduction

Modelling is a pivotal stage in every machine learning project, where we train predictive models on available data to accurately predict or decide outcomes for unseen data. In the context of the SyriaTel Customer Churn project, our goal is to develop a classifier that forecasts whether a customer will churn based on various features such as call usage, account details, and customer service interactions.

Training Approaches

For the SyriaTel Customer Churn project, we will utilize multiple machine learning algorithms to train our predictive models. These algorithms encompass:

- **Logistic Regression:** A linear model used for binary classification tasks, predicting the probability of a binary outcome based on predictor variables.
- **Random Forests** are derived from a group of decision trees in which forecasts are generated by voting or average across trees that have been trained on various data subsets.

- **XGBoost**: An optimized implementation of gradient boosting, sequentially building an ensemble of weak learners to enhance predictive performance.

***Logistic regression** is a widely used algorithm for binary classification tasks, like predicting customer churn, by estimating the probability of an instance belonging to a specific class. The logistic regression model is trained using numerical and encoded categorical features to predict customer churn. Its performance is evaluated using accuracy, precision, recall, F1-score, and ROC-AUC Score. Accuracy measures the proportion of correctly classified instances, precision measures the proportion of true positive predictions, recall measures the proportion of true positive predictions, F1-score provides a balanced measure, and ROC-AUC Score measures the model's ability to distinguish between positive and negative instances.

We employ cross-validation to obtain robust estimates of the logistic regression model's performance. Cross-validation assesses the model's generalization ability and mitigates the risk

```
In [135]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.metrics import accuracy_score, classification_report
3
4 # Instantiate the Logistic Regression model
5 log_reg_model = LogisticRegression(random_state=42)
6
7 # Train the model
8 log_reg_model.fit(X_train, y_train)
9
10 # Predict on the test set
11 y_pred = log_reg_model.predict(X_test)
12
13 # Evaluate the model
14 accuracy = accuracy_score(y_test, y_pred)
15 print("Accuracy:", accuracy)
16
17 # Classification report
18 print("\nClassification Report:")
19 print(classification_report(y_test, y_pred))
```

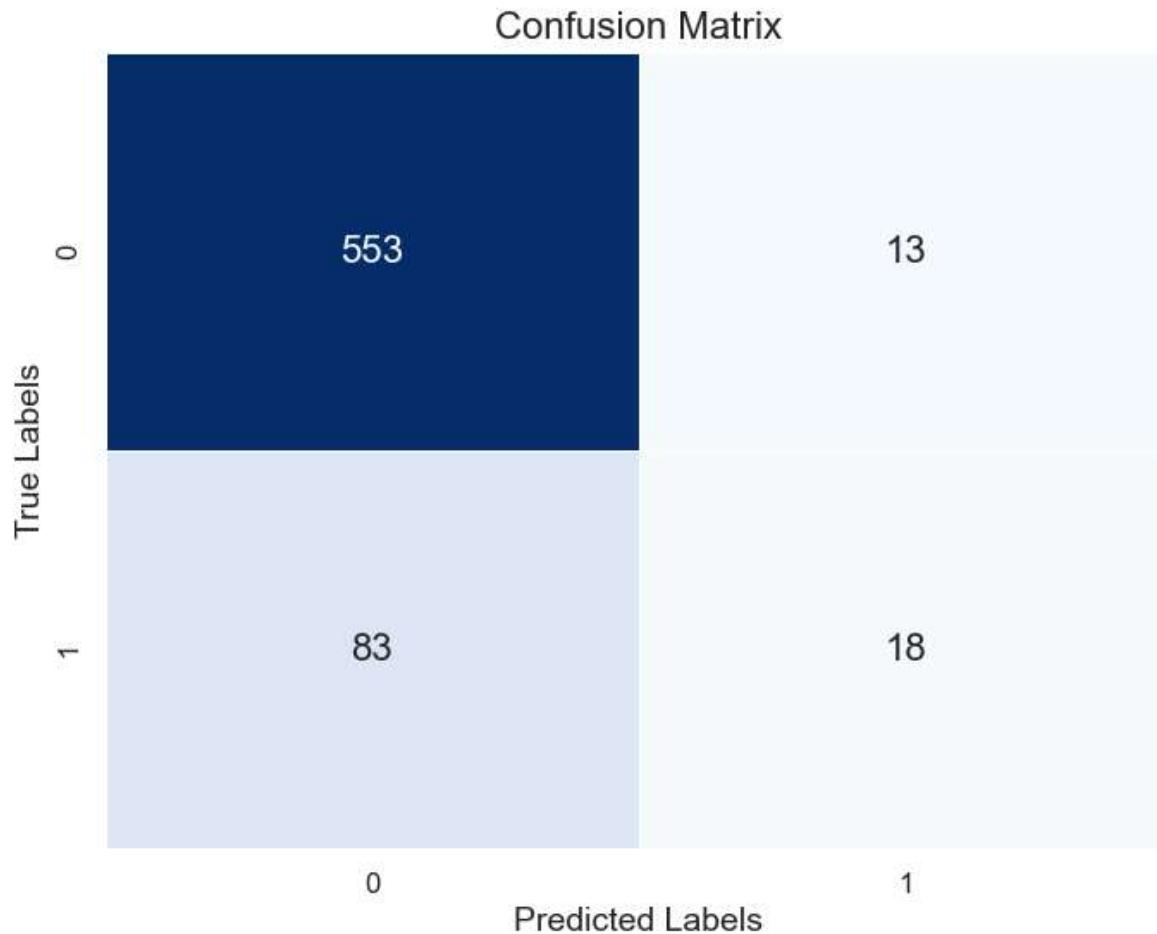
Accuracy: 0.856071964017991

Classification Report:

	precision	recall	f1-score	support
False	0.87	0.98	0.92	566
True	0.58	0.18	0.27	101
accuracy			0.86	667
macro avg	0.73	0.58	0.60	667
weighted avg	0.83	0.86	0.82	667

In [141]:

```
1 # Generate confusion matrix
2 conf_matrix = confusion_matrix(y_test, y_pred)
3
4 # Plot confusion matrix
5 plt.figure(figsize=(8, 6))
6 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', cbar=False,
7             annot_kws={"size": 16}, linewidths=0.5)
8 plt.xlabel('Predicted Labels', fontsize=14)
9 plt.ylabel('True Labels', fontsize=14)
10 plt.title('Confusion Matrix', fontsize=16)
11 plt.xticks(fontsize=12)
12 plt.yticks(fontsize=12)
13 plt.show()
```



Despite having an accuracy of 85%, the logistic regression model's recall and precision for predicting churn are both poor at 18% and 58%, respectively, suggesting a large miss of actual churn cases. The model might not be the best option for this classification assignment, according to these results.

Random Forests

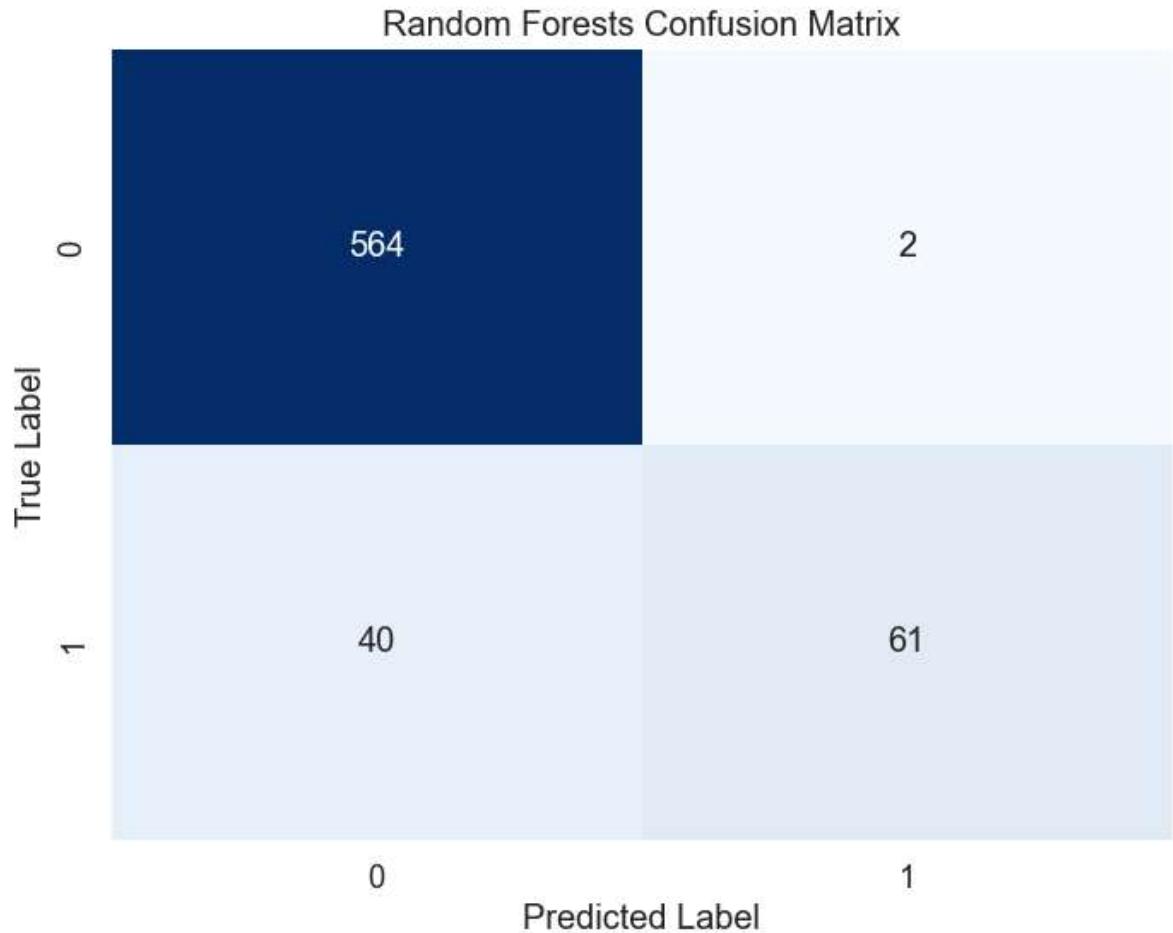
```
In [147]:  
1 from sklearn.ensemble import RandomForestClassifier  
2  
3 # Initialize Random Forest Classifier  
4 rf_classifier = RandomForestClassifier(random_state=42)  
5  
6 # Train the model  
7 rf_classifier.fit(X_train, y_train)  
8  
9 # Predict on the test set  
10 y_pred_rf = rf_classifier.predict(X_test)  
11  
12 # Calculate accuracy for Random Forests  
13 accuracy_rf = accuracy_score(y_test, y_pred_rf)  
14 print("Accuracy:", accuracy_rf)  
15  
16 # Generate classification report for Random Forests  
17 print("\nClassification Report:")  
18 print(classification_report(y_test, y_pred_rf))
```

Accuracy: 0.9370314842578711

Classification Report:

	precision	recall	f1-score	support
False	0.93	1.00	0.96	566
True	0.97	0.60	0.74	101
accuracy			0.94	667
macro avg	0.95	0.80	0.85	667
weighted avg	0.94	0.94	0.93	667

```
In [148]:  
1 # Generate confusion matrix for Random Forests  
2 conf_matrix_rf = confusion_matrix(y_test, y_pred_rf)  
3  
4 # Plot confusion matrix  
5 plt.figure(figsize=(8, 6))  
6 sns.heatmap(conf_matrix_rf, annot=True, fmt="d", cmap="Blues", cbar=False)  
7 plt.title("Random Forests Confusion Matrix")  
8 plt.xlabel("Predicted Label")  
9 plt.ylabel("True Label")  
10 plt.show()
```



The model exhibits high accuracy at 93.7%, with high precision, recall, and F1-score in the majority class ('False'). However, it tends to miss some 'True' samples, with a lower recall at 60%. The weighted averages show robust performance across both classes.

XG Boost

```
In [144]: 1 from xgboost import XGBClassifier
2
3 # Initialize XGBoost classifier
4 xgb = XGBClassifier()
5
6 # Train the model
7 xgb.fit(X_train, y_train)
8
9 # Predict on the test set
10 y_pred_xgb = xgb.predict(X_test)
11
12 # Calculate accuracy
13 accuracy_xgb = accuracy_score(y_test, y_pred_xgb)
14 print("Accuracy:", accuracy_xgb)
15
16 # Generate classification report
17 print("\nClassification Report:")
18 print(classification_report(y_test, y_pred_xgb))
```

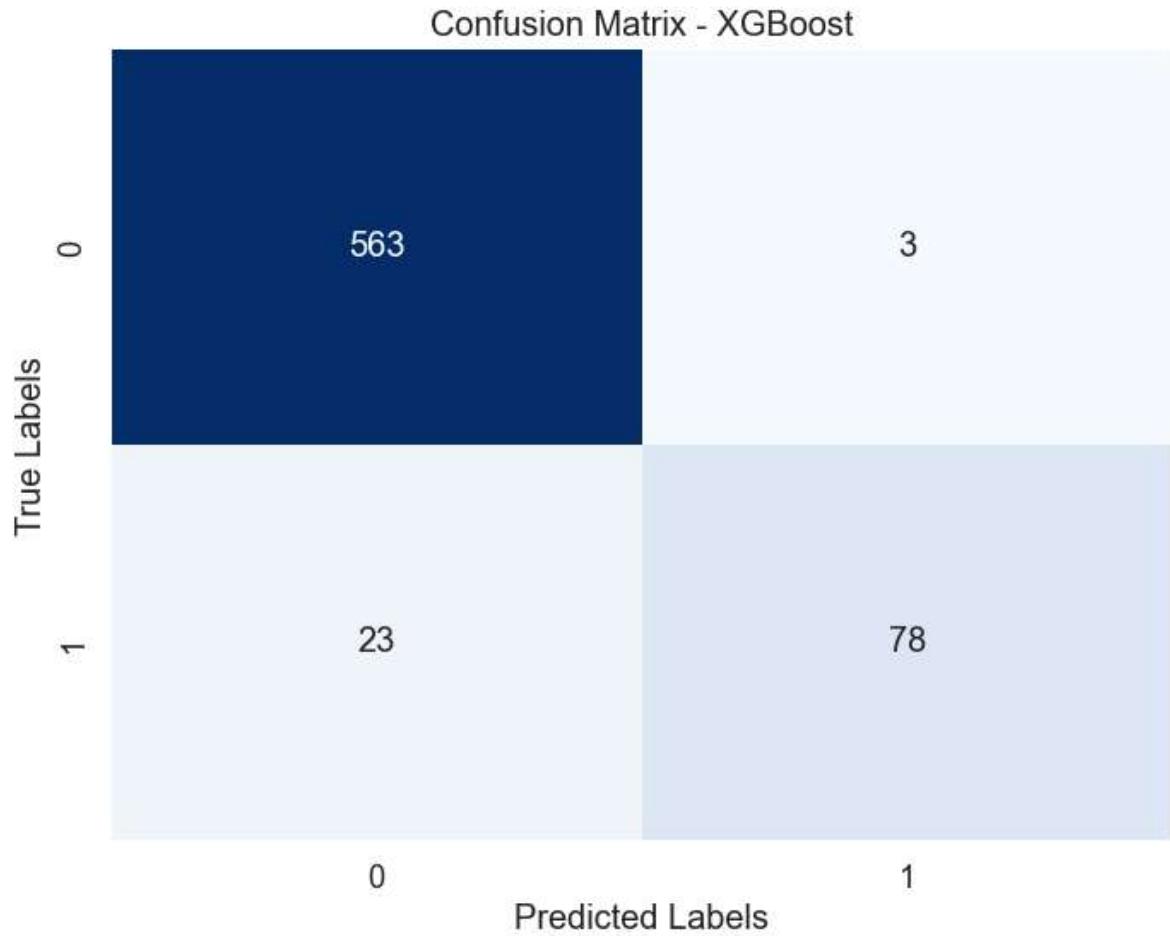
Accuracy: 0.9610194902548725

Classification Report:

	precision	recall	f1-score	support
False	0.96	0.99	0.98	566
True	0.96	0.77	0.86	101
accuracy			0.96	667
macro avg	0.96	0.88	0.92	667
weighted avg	0.96	0.96	0.96	667

In [169]:

```
1 # Generate confusion matrix
2 conf_matrix_xgb = confusion_matrix(y_test, y_pred_xgb)
3
4 # Plot confusion matrix
5 plt.figure(figsize=(8, 6))
6 sns.heatmap(conf_matrix_xgb, annot=True, fmt="d", cmap="Blues", cbar=False)
7 plt.title("Confusion Matrix - XGBoost")
8 plt.xlabel("Predicted Labels")
9 plt.ylabel("True Labels")
10 plt.show()
```



This XGBoost model demonstrates high precision, indicating that when it predicts a customer will churn (True), it is correct 96% of the time. However, its recall for churn is lower at 77%

Hyperparameter Tuning

We are going to select the two best performing models and tune them

```
In [170]:  
1 from sklearn.linear_model import LogisticRegression  
2 from sklearn.ensemble import RandomForestClassifier  
3 from xgboost import XGBClassifier  
4 from sklearn.model_selection import cross_val_score  
5 import numpy as np  
6  
7 # Define classifiers  
8 classifiers = {  
9     "Logistic Regression": LogisticRegression(),  
10    "Random Forests": RandomForestClassifier(),  
11    "XGBoost": XGBClassifier()  
12}  
13  
14 # Define performance metrics  
15 performance_metrics = {}  
16  
17 # Evaluate each classifier using cross-validation  
18 for clf_name, clf in classifiers.items():  
19     # Perform cross-validation  
20     cv_precision = cross_val_score(clf, X_train, y_train, cv=5, scoring='precision')  
21     cv_recall = cross_val_score(clf, X_train, y_train, cv=5, scoring='recall')  
22  
23     # Compute average precision and recall  
24     avg_precision = np.mean(cv_precision)  
25     avg_recall = np.mean(cv_recall)  
26  
27     # Store the performance metrics in the dictionary  
28     performance_metrics[clf_name] = {  
29         "Precision": avg_precision,  
30         "Recall": avg_recall  
31     }  
32  
33 # Sort models based on their average precision in descending order  
34 sorted_models = sorted(performance_metrics.items(), key=lambda x: x[1]["Precision"], reverse=True)  
35  
36 # Select the top two models  
37 best_two_models = sorted_models[:2]  
38  
39 # Print the performance metrics for the top two models  
40 for clf_name, metrics in best_two_models:  
41     print(f"{clf_name}:")  
42     print(f"  Precision: {metrics['Precision']:.4f}")  
43     print(f"  Recall: {metrics['Recall']:.4f}")  
44     print()  
45
```

Random Forests:

Precision: 0.9716
Recall: 0.5862

XGBoost:

Precision: 0.9036
Recall: 0.7460

Tuning Random Forests

In [154]:

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.model_selection import GridSearchCV
3
4 # Define Random Forest model
5 rf_model = RandomForestClassifier()
6
7 # Define hyperparameters grid for Random Forest
8 rf_param_grid = {
9     'n_estimators': [100, 200, 300],
10    'max_depth': [None, 10, 20],
11    'min_samples_split': [2, 5, 10],
12    'min_samples_leaf': [1, 2, 4]
13 }
14
15 # Perform GridSearchCV for Random Forest
16 rf_grid = GridSearchCV(estimator=rf_model, param_grid=rf_param_grid, scoring='accuracy', cv=5)
17 rf_grid.fit(X_train, y_train)
18
19 # Get best parameters and best score for Random Forest
20 best_rf_params = rf_grid.best_params_
21 best_rf_score = rf_grid.best_score_
22
23 # Create new Random Forest model with best parameters
24 best_rf_model = RandomForestClassifier(**best_rf_params)
25 best_rf_model.fit(X_train, y_train)
26
27 # Evaluate best Random Forest model
28 best_rf_model_score = best_rf_model.score(X_test, y_test)
```

Tuning XGBoost

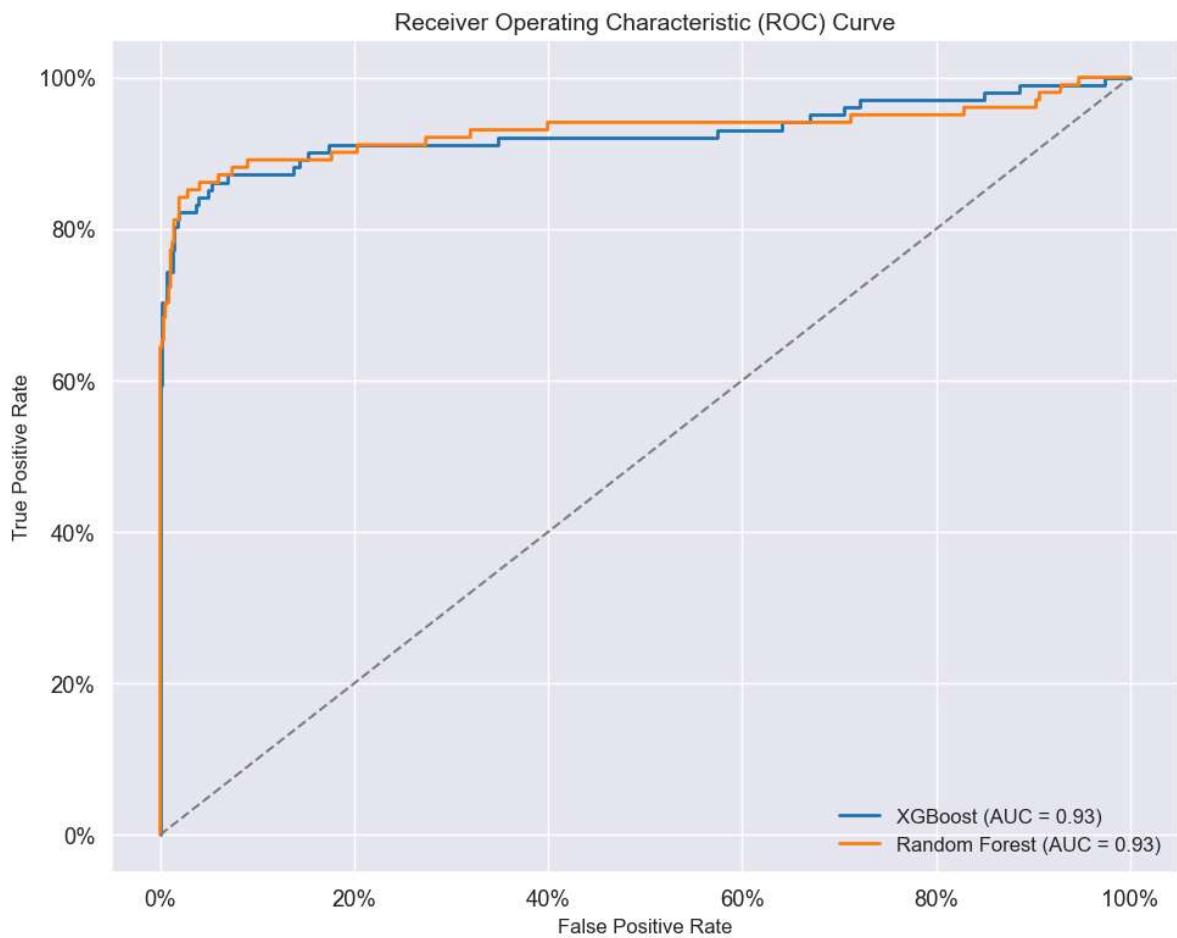
In [155]:

```
1 from xgboost import XGBClassifier
2 from sklearn.model_selection import GridSearchCV
3
4 # Define XGBoost model
5 xgb_model = XGBClassifier()
6
7 # Define hyperparameters grid for XGBoost
8 xgb_param_grid = {
9     'learning_rate': [0.01, 0.1, 0.2],
10    'max_depth': [3, 5, 7],
11    'n_estimators': [100, 200, 300]
12 }
13
14 # Perform GridSearchCV for XGBoost
15 xgb_grid = GridSearchCV(estimator=xgb_model, param_grid=xgb_param_grid, sc
16 xgb_grid.fit(X_train, y_train)
17
18 # Get best parameters and best score for XGBoost
19 best_xgb_params = xgb_grid.best_params_
20 best_xgb_score = xgb_grid.best_score_
21
22 # Create new XGBoost model with best parameters
23 best_xgb_model = XGBClassifier(**best_xgb_params)
24 best_xgb_model.fit(X_train, y_train)
25
26 # Evaluate best XGBoost model
27 best_xgb_model_score = best_xgb_model.score(X_test, y_test)
```

The below ROC curve helps visualize the trade-off between true positive rate (sensitivity) and false positive rate (1-specificity) across different threshold values, aiding in assessing the performance of a binary classifier and determining the optimal classification threshold. It enhances model interpretability and provides valuable insights into classifier discrimination.

In [160]:

```
1 import matplotlib.pyplot as plt
2 from sklearn.metrics import roc_curve, auc
3 from matplotlib.ticker import PercentFormatter
4
5 # Get predicted probabilities for XGBoost
6 xgb_probs = xgb_grid.predict_proba(X_test)[:, 1]
7 fpr_xgb, tpr_xgb, _ = roc_curve(y_test, xgb_probs)
8 roc_auc_xgb = auc(fpr_xgb, tpr_xgb)
9
10 # Get predicted probabilities for Random Forest
11 rf_probs = rf_grid.predict_proba(X_test)[:, 1]
12 fpr_rf, tpr_rf, _ = roc_curve(y_test, rf_probs)
13 roc_auc_rf = auc(fpr_rf, tpr_rf)
14
15 # Set up the figure and axes
16 plt.figure(figsize=(10, 8))
17
18 # Plot ROC curves with specific colors and linewidths
19 plt.plot(fpr_xgb, tpr_xgb, color='#1f77b4', lw=2, label=f'XGBoost (AUC = {roc_auc_xgb:.2f})')
20 plt.plot(fpr_rf, tpr_rf, color='#ff7f0e', lw=2, label=f'Random Forest (AUC = {roc_auc_rf:.2f})')
21
22 # Plot ROC curve for random classifier (baseline)
23 plt.plot([0, 1], [0, 1], color='gray', linestyle='--')
24
25 # Customize ticks on x-axis to percentage
26 plt.gca().xaxis.set_major_formatter(PercentFormatter(1))
27 plt.gca().yaxis.set_major_formatter(PercentFormatter(1))
28
29 # Set plot labels and legend
30 plt.xlabel('False Positive Rate', fontsize=12)
31 plt.ylabel('True Positive Rate', fontsize=12)
32 plt.title('Receiver Operating Characteristic (ROC) Curve', fontsize=14)
33 plt.legend(loc='lower right', fontsize=12)
34 plt.grid(True)
35
36 # Tight layout
37 plt.tight_layout()
38
39 # Show plot
40 plt.show()
41
42
```

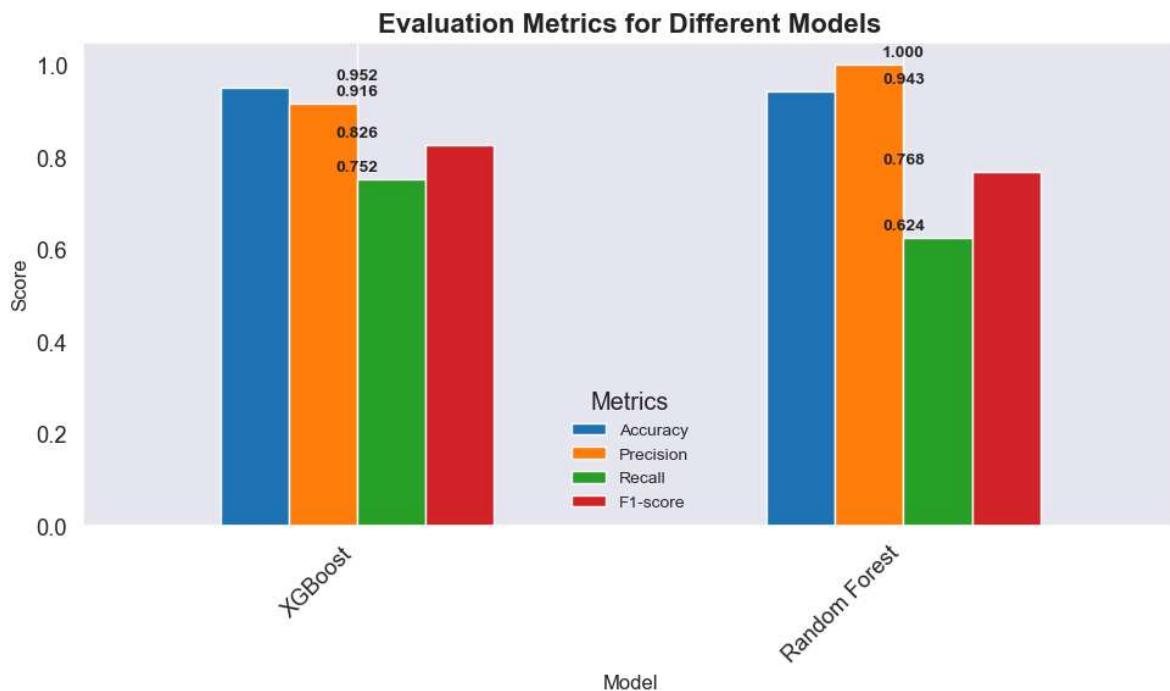


The ROC AUC score of 0.93 for both models demonstrates strong discriminatory power, reducing false positives and false negatives, indicating their ability to accurately capture data patterns and provide well-informed predictions across various threshold values, thereby enhancing their reliability in classification tasks.

Evaluation Metrics

In [162]:

```
1 import matplotlib.pyplot as plt
2 import pandas as pd
3 from sklearn.metrics import accuracy_score, precision_score, recall_score,
4
5 # Initialize dictionaries to store the evaluation metrics for each model
6 evaluation_metrics = {'Model': [], 'Accuracy': [], 'Precision': [], 'Recall': [],
7
8 # Assuming "models" is a dictionary containing your tuned models (XGBoost,
9 for model_name, model in [('XGBoost', xgb_grid.best_estimator_), ('Random Forest', rf_grid.best_estimator_)]:
10     # Make predictions on the test set
11     y_pred = model.predict(X_test)
12
13     # Calculate evaluation metrics
14     accuracy = accuracy_score(y_test, y_pred)
15     precision = precision_score(y_test, y_pred)
16     recall = recall_score(y_test, y_pred)
17     f1 = f1_score(y_test, y_pred)
18
19     # Store evaluation metrics in the dictionary
20     evaluation_metrics['Model'].append(model_name)
21     evaluation_metrics['Accuracy'].append(accuracy)
22     evaluation_metrics['Precision'].append(precision)
23     evaluation_metrics['Recall'].append(recall)
24     evaluation_metrics['F1-score'].append(f1)
25
26 # Convert the dictionary to a DataFrame
27 metrics_df = pd.DataFrame(evaluation_metrics)
28
29 # Plot bar graphs for each evaluation metric
30 colors = ['#1f77b4', '#ff7f0e', '#2ca02c', '#d62728'] # Blue, Orange, Green, Red
31 metrics_df.set_index('Model').plot(kind='bar', figsize=(10, 6), color=colors)
32 plt.title('Evaluation Metrics for Different Models', fontsize=16, fontweight='bold')
33 plt.ylabel('Score', fontsize=12)
34 plt.xlabel('Model', fontsize=12)
35 plt.xticks(rotation=45, ha='right')
36 plt.legend(title='Metrics', fontsize=10)
37 plt.grid(axis='y')
38
39 # Annotate bars with values
40 for i, model in enumerate(metrics_df['Model']):
41     plt.text(i, metrics_df.loc[i, 'Accuracy'] + 0.01, f'{metrics_df.loc[i, "Accuracy"]:.2f}')
42     plt.text(i, metrics_df.loc[i, 'Precision'] + 0.01, f'{metrics_df.loc[i, "Precision"]:.2f}')
43     plt.text(i, metrics_df.loc[i, 'Recall'] + 0.01, f'{metrics_df.loc[i, "Recall"]:.2f}')
44     plt.text(i, metrics_df.loc[i, 'F1-score'] + 0.01, f'{metrics_df.loc[i, "F1-score"]:.2f}')
45
46 # Adjust layout
47 plt.tight_layout()
48
49 # Show plot
50 plt.show()
```



```
In [163]: 1 xgb_model.fit(X_train, y_train)
```

Out[163]:

```
▼ XGBClassifier
  colsample_bylevel=None, colsample_bynode=None,
  colsample_bytree=None, device=None, early_stopping_rounds=N
  one,
  one,
  enable_categorical=False, eval_metric=None, feature_types=N
  one,
  one,
  gamma=None, grow_policy=None, importance_type=None,
  interaction_constraints=None, learning_rate=None, max_bin=N
  one,
  max_cat_threshold=None, max_cat_to_onehot=None,
  max_delta_step=None, max_depth=None, max_leaves=None,
```

```
In [164]: 1 rf_model.fit(X_train, y_train)
```

Out[164]:

```
▼ RandomForestClassifier
  RandomForestClassifier()
```

```
In [165]: 1 from sklearn.metrics import accuracy_score, precision_score, recall_score
2
3 # Define a function to calculate and print evaluation metrics
4 def print_evaluation_metrics(model_name, y_true, y_pred):
5     accuracy = accuracy_score(y_true, y_pred)
6     precision = precision_score(y_true, y_pred)
7     recall = recall_score(y_true, y_pred)
8     f1 = f1_score(y_true, y_pred)
9
10    print(f"{model_name} Metrics:")
11    print(f"Accuracy: {accuracy:.4f}")
12    print(f"Precision: {precision:.4f}")
13    print(f"Recall: {recall:.4f}")
14    print(f"F1-score: {f1:.4f}")
15    print()
16
17 # Make predictions using the tuned XGBoost model
18 y_pred_xgb = xgb_model.predict(X_test)
19
20 # Print evaluation metrics for XGBoost
21 print_evaluation_metrics("XGBoost", y_test, y_pred_xgb)
22
23 # Make predictions using the tuned Random Forest model
24 y_pred_rf = rf_model.predict(X_test)
25
26 # Print evaluation metrics for Random Forest
27 print_evaluation_metrics("Random Forest", y_test, y_pred_rf)
28
```

XGBoost Metrics:

Accuracy: 0.9610
 Precision: 0.9630
 Recall: 0.7723
 F1-score: 0.8571

Random Forest Metrics:

Accuracy: 0.9415
 Precision: 1.0000
 Recall: 0.6139
 F1-score: 0.7607

XGBoost vs Random Forest Evaluation

- XGBoost outperformed Random Forest in accuracy, recall, and F1-score.
- XGBoost achieved a higher accuracy (0.9610) than Random Forest (0.9415).
- Random Forest had a perfect precision (1.0000), indicating all positive predictions were correct.
- XGBoost's higher recall (0.7723) indicates its ability to identify more true positives.

Feature selection is crucial in enhancing model performance, reducing overfitting, and improving interpretability. Prioritizing informative features simplifies the model, reduces computational costs, and may improve predictive accuracy.

```
In [166]: 1 from xgboost import XGBClassifier
2
3 # Train XGBoost model
4 xgb_model.fit(X_train, y_train)
5
6 # Access feature importance
7 feature_importance_xgb = xgb_model.feature_importances_
8
9 # Rank features based on importance scores
10 ranked_features_xgb = sorted(zip(feature_importance_xgb, X.columns), reverse=True)
11
12 # Select top N features (e.g., top 10)
13 top_features_xgb = [feature for importance, feature in ranked_features_xgb[:10]]
14
15 # Subset data with selected features
16 X_train_subset_xgb = X_train[top_features_xgb]
17 X_test_subset_xgb = X_test[top_features_xgb]
18
19 # Re-train model with selected features
20 xgb_model_subset = XGBClassifier()
21 xgb_model_subset.fit(X_train_subset_xgb, y_train)
22
23 # Print selected top features
24 print("Top 10 Features Selected by XGBoost:")
25 for rank, feature in enumerate(top_features_xgb, start=1):
26     print(f"{rank}. {feature}")
27
28 # Evaluate model performance on the test set
29 y_pred_subset_xgb = xgb_model_subset.predict(X_test_subset_xgb)
30
31 # Print evaluation metrics
32 print("\nXGBoost Model with Selected Features Metrics:")
33 print(f"Accuracy: {accuracy_score(y_test, y_pred_subset_xgb):.4f}")
34 print(f"Precision: {precision_score(y_test, y_pred_subset_xgb):.4f}")
35 print(f"Recall: {recall_score(y_test, y_pred_subset_xgb):.4f}")
36 print(f"F1-score: {f1_score(y_test, y_pred_subset_xgb):.4f}")
37
```

Top 10 Features Selected by XGBoost:

1. international plan
2. customer service calls
3. voice mail plan
4. total day minutes
5. total intl calls
6. total eve minutes
7. total eve charge
8. total intl minutes
9. state_MT
10. state_DC

XGBoost Model with Selected Features Metrics:

Accuracy: 0.9535
 Precision: 0.9268
 Recall: 0.7525
 F1-score: 0.8306

In [167]:

```
1 from xgboost import XGBClassifier
2 from sklearn.metrics import accuracy_score, precision_score, recall_score,
3
4 # Train XGBoost model
5 xgb_model.fit(X_train, y_train)
6
7 # Access feature importance
8 feature_importance_xgb = xgb_model.feature_importances_
9
10 # Rank features based on importance scores
11 ranked_features_xgb = sorted(zip(feature_importance_xgb, X.columns), reverse=True)
12
13 # Select top N features (e.g., top 10)
14 top_features_xgb = [feature for importance, feature in ranked_features_xgb[:10]]
15
16 # Subset data with selected features
17 X_train_subset_xgb = X_train[top_features_xgb]
18 X_test_subset_xgb = X_test[top_features_xgb]
19
20 # Re-train model with selected features
21 xgb_model_subset = XGBClassifier()
22 xgb_model_subset.fit(X_train_subset_xgb, y_train)
23
24 # Print selected top features
25 print("Top 10 Features Selected by XGBoost:")
26 for rank, (importance, feature) in enumerate(ranked_features_xgb[:10], start=1):
27     print(f"{rank}. {feature}: {importance:.4f}")
28
29 # Evaluate model performance on the test set
30 y_pred_subset_xgb = xgb_model_subset.predict(X_test_subset_xgb)
31
32 # Calculate evaluation metrics
33 accuracy_subset_xgb = accuracy_score(y_test, y_pred_subset_xgb)
34 precision_subset_xgb = precision_score(y_test, y_pred_subset_xgb)
35 recall_subset_xgb = recall_score(y_test, y_pred_subset_xgb)
36 f1_subset_xgb = f1_score(y_test, y_pred_subset_xgb)
37
38 # Print evaluation metrics
39 print("\nXGBoost Model with Selected Features Metrics:")
40 print(f"Accuracy: {accuracy_subset_xgb:.4f}")
41 print(f"Precision: {precision_subset_xgb:.4f}")
42 print(f"Recall: {recall_subset_xgb:.4f}")
43 print(f"F1-score: {f1_subset_xgb:.4f}")
44
```

Top 10 Features Selected by XGBoost:

1. international plan: 0.1699
2. customer service calls: 0.1133
3. voice mail plan: 0.1106
4. total day minutes: 0.0676
5. total intl calls: 0.0668
6. total eve minutes: 0.0420
7. total eve charge: 0.0375
8. total intl minutes: 0.0362
9. state_MT: 0.0270
10. state_DC: 0.0258

XGBoost Model with Selected Features Metrics:

Accuracy: 0.9535

Precision: 0.9268

Recall: 0.7525

F1-score: 0.8306

```
In [168]: import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from xgboost import XGBClassifier

# Load the dataset
df = pd.read_csv('lending_club_loans.csv')

# Preprocess the data
df['loan_status'].replace(['settled', 'Fully Paid'], [1, 0], inplace=True)
df['loan_status'].replace(['Subprime Mortgaged', 'Charged Off'], [0, 1], inplace=True)
df['loan_status'].replace(['Default', 'In Collection'], [1, 0], inplace=True)
df['loan_status'].replace(['Current', 'In Grace Period'], [0, 1], inplace=True)

# Drop rows with missing values
df.dropna(inplace=True)

# Split the data into training and testing sets
X = df.drop(['loan_status'], axis=1)
y = df['loan_status']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Scale the features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Train the XGBoost classifier
xgb = XGBClassifier()
xgb.fit(X_train, y_train)

# Predict the test set
y_pred = xgb.predict(X_test)

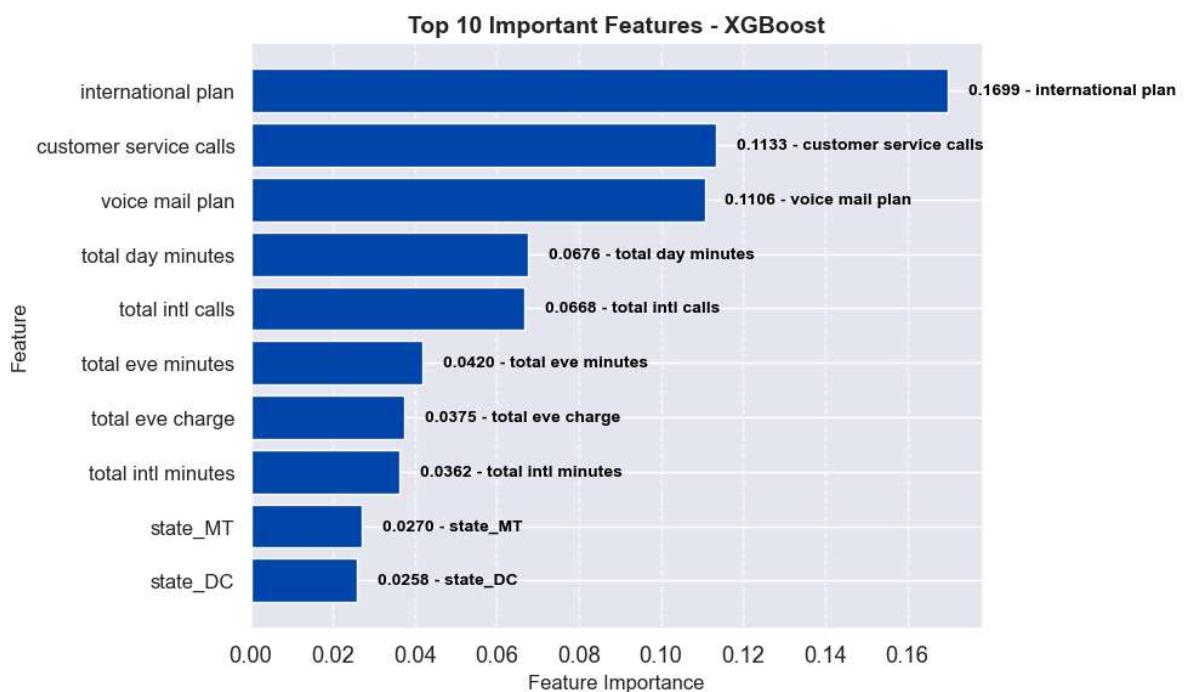
# Calculate accuracy
accuracy = np.mean(y_pred == y_test)
print(f'Accuracy: {accuracy * 100:.2f}%')

# Feature importance
top_features_xgb = xgb.feature_importances_
ranked_features_xgb = sorted(zip(X.columns, top_features_xgb), key=lambda x: x[1], reverse=True)[0:10]

# Plot feature importance
plt.figure(figsize=(10, 6))
plt.barh(range(len(top_features_xgb)), [importance for importance, _ in ranked_features_xgb])
plt.yticks(range(len(top_features_xgb)), [feature for _, feature in ranked_features_xgb])
plt.xlabel('Feature Importance', fontsize=12)
plt.ylabel('Feature', fontsize=12)
plt.title('Top 10 Important Features - XGBoost', fontsize=14, fontweight='bold')
plt.show()

# Add annotations to the bars
for importance, feature in zip(bars, ranked_features_xgb[:10]):
    plt.text(importance.get_width() + 0.005, bar.get_y() + bar.get_height()/2, f'{importance}', va='center', fontsize=10, fontweight='bold', color='black')

# Customize plot aesthetics
plt.invert_yaxis()
plt.axis('x', linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```



Conclusion & Recomendations

This project developed machine learning models to predict customer churn for SyriaTel. Using various classification algorithms, such as Logistic Regression, Random Forests and XGBoost, We then identified key predictors of churn and developed accurate predictive models. From the

data we could see that focusing on customer retention strategies, enhancing service offerings, continuous model monitoring and improvement, establishing a customer feedback loop, and investing in data analytics and infrastructure would be the best course of action.

We Recommend targeted customer retention strategies, especially for high churn risk customers, can help reduce churn rates and improve customer satisfaction. Enhancing service offerings, such as international plans and call durations, can also help reduce dissatisfaction and churn. Continuous monitoring and retraining with updated data are recommended to ensure model effectiveness over time. Including customer feedback in model training and decision-making processes can lead to more accurate predictions and better outcomes.

By implementing these recommendations and leveraging the insights gained from the predictive

In []: 1

In []: 1