

Name

Ogada Thomas Otiende

email: otiende.ogada@student.moringaschool.com (<mailto:otiende.ogada@student.moringaschool.com>)

Business Understanding

Customer churn, the phenomenon where customers cease doing business with a company, is a critical concern for telecommunications companies like SyriaTel. Retaining customers is essential for maintaining revenue and growth in this competitive industry. Identifying factors contributing to churn, such as service dissatisfaction or competitive offers, SyriaTel can take targeted actions to mitigate churn and improve customer retention.

Introduction

This project aims to build a predictive model for SyriaTel, a telecommunications company, to identify customers at risk of churning. Accurately predicting customer churn, SyriaTel can proactively implement retention strategies, thereby reducing financial losses and enhancing customer loyalty.

Background

Syria is a country located in the Middle East, has a telecommunications sector experiencing rapid growth in mobile and internet penetration. SyriaTel, as a key player in this sector, plays a vital role in connecting people and businesses. However, increasing competition and evolving customer preferences pose challenges for customer retention. Understanding and addressing the drivers of churn are crucial for SyriaTel to sustain business success and enhance customer satisfaction.

Business Problem

SyriaTel, a telecommunications company, faces the challenge of customer churn, where customers discontinue their services. This attrition impacts revenue and profitability. The business seeks to proactively identify customers at risk of churning and implement effective retention strategies to mitigate revenue loss and maintain customer loyalty.

Specifically, the project aims to address the following questions:

1. What are the primary factors driving customer churn for SyriaTel?
2. Which machine learning modelling technique to apply in accurately predicting Churn so as to take proactive measures?
3. What actionable insights can SyriaTel derive from the predictive model to improve customer retention efforts?
4. What strategies can SyriaTel put in place to reduce churn rate?

Data Understanding

Importing neccessary Libraries

```
In [1]: 1 import warnings
2 warnings.filterwarnings("ignore")
3 import pandas as pd
4 import numpy as np
5 import datetime as dt
6 from collections import Counter
7 import calendar
8 from dateutil import relativedelta
9 import operator
10 import os
11 import random
12 from functools import reduce
13 import matplotlib.pyplot as plt
14 import seaborn as sns
15 from itertools import combinations
16 import warnings
17 import matplotlib.ticker as ticker
18
19 pd.set_option('display.max_columns', None)
20 pd.set_option('display.max_rows', None)
21 pd.set_option('display.float_format', lambda x: f'{len(str(x%1))-2}f' % x)
22 pd.set_option('display.max_colwidth', None)
23 %matplotlib inline
```

Loading the dataset

```
In [2]: 1 df = pd.read_csv("C:\\Users\\User\\Documents\\Phase_3_Project\\SyriaTel_df.csv")
2 df = df.copy()
```

Understanding the dataframe

```
In [3]: 1 df.head()
```

Out[3]:

	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
0	KS	128	415	382-4657	no	yes	25	265.10000000000002274	110	45.070000000000000028	197.400000
1	OH	107	415	371-7191	no	yes	26	161.59999999999999943	123	27.46999999999999886	
2	NJ	137	415	358-1921	no	no	0	243.40000000000000057	114	41.380000000000000256	121.2000000
3	OH	84	408	375-9999	yes	no	0	299.39999999999997726	71	50.8999999999999986	61.899999
4	OK	75	415	330-6626	yes	no	0	166.6999999999999886	113	28.3399999999999986	148.3000000

In [4]: 1 df.tail()

Out[4]:

	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
3328	AZ	192	415	414-4276	no	yes	36	156.19999999999998863	77	26.550000000000000007	
3329	WV	68	415	370-3271	no	no	0	231.09999999999999432	57	39.28999999999999915	153.4000
3330	RI	28	510	328-8230	no	no	0	180.80000000000000114	109	30.73999999999999984	288.8000
3331	CT	184	510	364-6381	yes	no	0	213.80000000000000114	105	36.35000000000000014	159.5999
3332	TN	74	415	400-4344	no	yes	25	234.40000000000000057	113	39.85000000000000014	265.8999

In [5]: 1 df.columns

Out[5]: Index(['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 minutes', 'total eve calls', 'total eve charge', 'total night minutes', 'total night calls', 'total night charge', 'total intl minutes', 'total intl calls', 'total intl charge', 'customer service calls', 'churn'], dtype='object')

In [6]: 1 df.shape

Out[6]: (3333, 21)

Changing Columns into Title cases

In [7]: 1 df.columns = df.columns.str.title()
2 df.columns

Out[7]: Index(['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 Minutes', 'Total Eve Calls', 'Total Eve Charge', 'Total Night Minutes', 'Total Night Calls', 'Total Night Charge', 'Total Intl Minutes', 'Total Intl Calls', 'Total Intl Charge', 'Customer Service Calls', 'Churn'], dtype='object')

In [8]:

```
1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3333 entries, 0 to 3332
Data columns (total 21 columns):
 #   Column                                Non-Null Count  Dtype  
---  -
 0   State                                3333 non-null   object  
 1   Account Length                      3333 non-null   int64   
 2   Area Code                          3333 non-null   int64   
 3   Phone Number                       3333 non-null   object  
 4   International Plan                  3333 non-null   object  
 5   Voice Mail Plan                    3333 non-null   object  
 6   Number Vmail Messages              3333 non-null   int64   
 7   Total Day Minutes                   3333 non-null   float64  
 8   Total Day Calls                     3333 non-null   int64   
 9   Total Day Charge                    3333 non-null   float64  
10  Total Eve Minutes                   3333 non-null   float64  
11  Total Eve Calls                     3333 non-null   int64   
12  Total Eve Charge                    3333 non-null   float64  
13  Total Night Minutes                 3333 non-null   float64  
14  Total Night Calls                   3333 non-null   int64   
15  Total Night Charge                  3333 non-null   float64  
16  Total Intl Minutes                  3333 non-null   float64  
17  Total Intl Calls                    3333 non-null   int64   
18  Total Intl Charge                   3333 non-null   float64  
19  Customer Service Calls              3333 non-null   int64   
20  Churn                               3333 non-null   bool    
dtypes: bool(1), float64(8), int64(8), object(4)
memory usage: 524.2+ KB
```

Identifying the type of columns

In [9]:

```
1 # Identifying columns
2 df['Churn'] = df['Churn'].astype(bool)
3
4 numerical_cols = df.select_dtypes(include=['int64', 'float64']).columns.tolist()
5 categorical_cols = df.select_dtypes(include=['object', 'bool']).columns.tolist()
6
7 # Removing 'Phone Number' from categorical columns if it exists
8 string_cols = ['Phone Number']
9 for col in string_cols:
10     if col in categorical_cols:
11         categorical_cols.remove(col)
12
13 # Print the identified columns
14 print("Numerical columns:")
15 print(numerical_cols)
16
17 print("\nCategorical columns:")
18 print(categorical_cols)
19
20 print("\nString columns:")
21 print(string_cols)
```

Numerical columns:

```
['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', 'Total Night Minutes', 'Total Night Calls', 'Total Night Charge', 'Total Intl Minutes', 'Total Intl Calls', 'Total Intl Charge', 'Customer Service Calls']
```

Categorical columns:

```
['State', 'International Plan', 'Voice Mail Plan', 'Churn']
```

String columns:

```
['Phone Number']
```

```
In [10]: 1 df.count()
```

```
Out[10]: State                3333
Account Length              3333
Area Code                  3333
Phone Number               3333
International Plan          3333
Voice Mail Plan            3333
Number Vmail Messages      3333
Total Day Minutes          3333
Total Day Calls             3333
Total Day Charge            3333
Total Eve Minutes          3333
Total Eve Calls             3333
Total Eve Charge            3333
Total Night Minutes        3333
Total Night Calls           3333
Total Night Charge          3333
Total Intl Minutes         3333
Total Intl Calls            3333
Total Intl Charge           3333
Customer Service Calls      3333
Churn                      3333
dtype: int64
```

Checking for null values

```
In [11]: 1 df.isnull().sum()
```

```
Out[11]: 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
```

Checking for duplicates

```
In [12]: 1 df.duplicated().sum()
```

```
Out[12]: 0
```

In [13]:

```
1 # function to identify unique values
2 for column in df.select_dtypes(include=['number']):
3     unique_values = df[column].unique()
4     print(f"Unique values in column '{column}': {unique_values}")
5
```

4.94 9.02 11.22 4.97 9.15 5.45 7.27 12.91 7.75 13.46 6.32 12.13

11.97 6.93 11.66 7.42 6.19 11.41 10.33 10.65 11.92 4.77 4.38 7.41

12.1 7.69 8.78 9.36 9.05 12.7 6.16 6.05 10.85 8.93 3.48 10.4

5.05 10.71 9.37 6.75 8.12 11.77 11.49 11.06 11.25 11.03 10.82 8.91

8.57 8.09 10.05 11.7 10.17 8.74 5.51 11.11 3.29 10.13 6.8 8.49

9.55 11.02 9.91 7.84 10.62 9.97 3.44 7.35 9.79 8.89 8.14 6.94

10.49 10.57 10.2 6.29 8.79 10.04 12.41 15.97 9.1 11.78 12.75 11.07

12.56 8.63 8.02 10.42 8.7 9.98 7.62 8.33 6.59 13.12 10.46 6.63

8.32 9.04 9.28 10.76 9.64 11.44 6.48 10.81 12.66 11.34 8.75 13.05

11.48 14.04 13.47 5.63 6.6 9.72 11.68 6.41 9.32 12.95 13.37 9.62

6.03 8.25 8.26 11.96 9.9 9.23 5.58 7.22 6.64 12.29 12.93 11.32

6.85 8.88 7.03 8.48 3.59 5.86 6.23 7.61 7.66 13.63 7.9 11.82

7.47 6.08 8.4 5.74 10.94 10.35 10.68 4.34 8.73 5.14 8.24 9.99

13.93 8.64 11.43 5.79 9.2 10.14 12.11 7.53 12.46 8.46 8.95 9.84

10.8 11.23 10.15 9.21 14.46 6.67 12.83 9.66 9.59 10.48 8.36 4.84

10.54 8.39 7.43 9.06 8.94 11.13 8.87 8.5 7.6 10.73 9.56 10.77

7.73 3.47 11.86 8.11 9.78 9.42 9.65 7. 7.39 9.88 6.56 5.92

6.95 15.71 8.06 4.86 7.8 8.58 10.06 5.21 6.92 6.15 13.49 9.38

12.62 12.26 8.19 11.65 11.62 10.83 7.92 7.33 13.01 13.26 12.22 11.58

5.97 10.99 8.38 9.17 8.08 5.71 3.41 12.63 11.79 12.96 7.64 6.58

Describing the dataframe

In [14]:

```
1 df.describe()
```

Out[14]:

	Account Length	Area Code	Number Vmail Messages	Total Day Minutes	Total Day Calls	Total
count	3333.0	3333.0	3333.0	3333.0	3333.0	
mean	101.06480648064805905	437.18241824182416622	8.0990099009900991	179.7750975097509354	100.4356435643564396	30.562
std	39.8221059285956045	42.3712904856066146	13.6883653720385983	54.46738920237137194	20.0690842073008966	9.2594
min	1.0	408.0	0.0	0.0	0.0	
25%	74.0	408.0	0.0	143.6999999999999886	87.0	24.4299
50%	101.0	415.0	0.0	179.4000000000000057	101.0	
75%	127.0	510.0	20.0	216.4000000000000057	114.0	36.7899
max	243.0	510.0	51.0	350.8000000000000114	165.0	59.6400

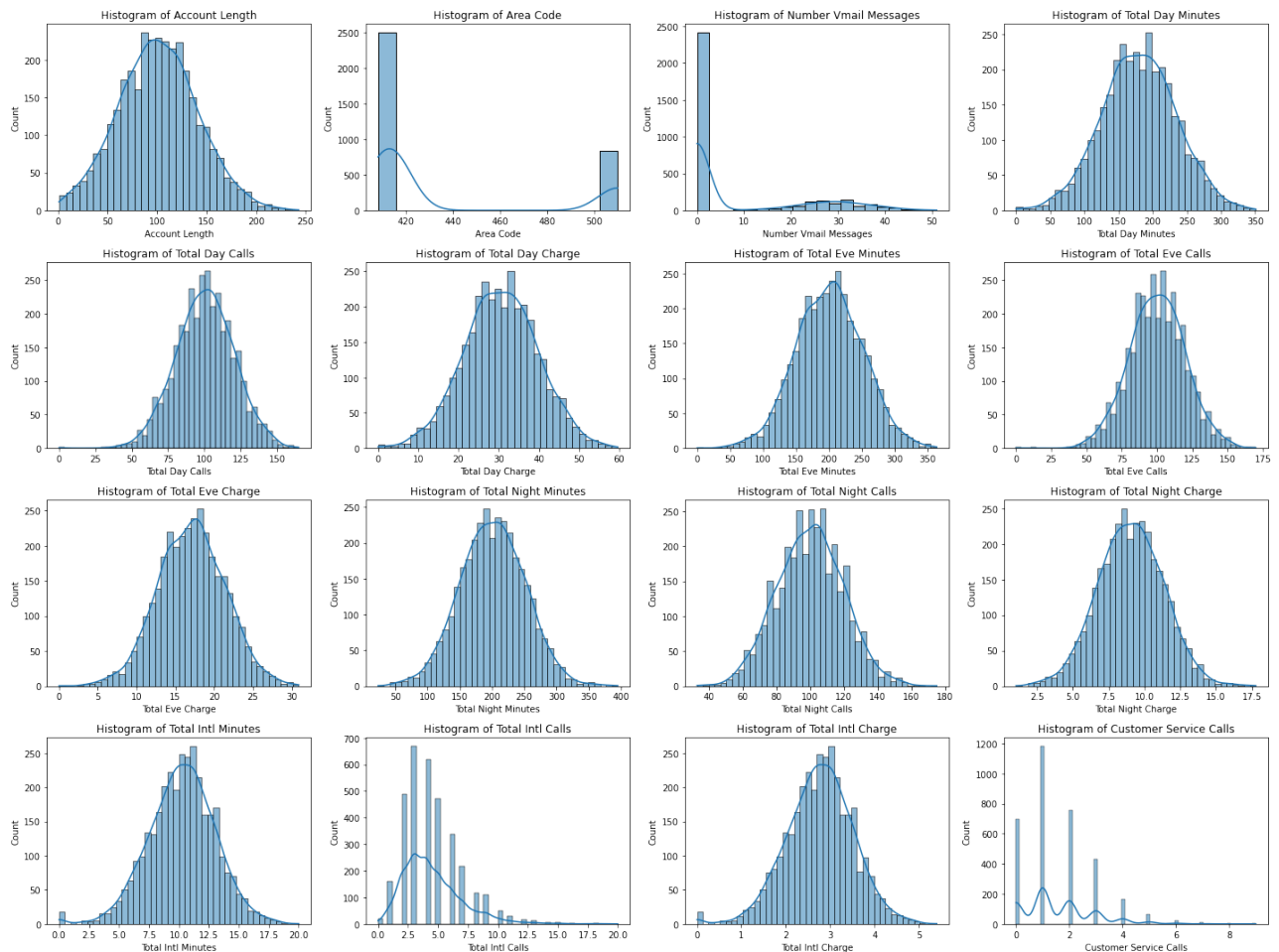
Univariant Analysis

Numerical Columns

```

In [15]: 1 # List of numerical columns
2 numerical_cols = ['Account Length', 'Area Code', 'Number Vmail Messages', 'Total Day Minutes', 'Total
3               'Total Day Charge', 'Total Eve Minutes', 'Total Eve Calls', 'Total Eve Charge',
4               'Total Night Minutes', 'Total Night Calls', 'Total Night Charge', 'Total Intl Minut
5               'Total Intl Calls', 'Total Intl Charge', 'Customer Service Calls']
6
7 # Set the size of the plots
8 plt.figure(figsize=(20, 15))
9
10 # Create histograms for each numerical column
11 for i, col in enumerate(numerical_cols):
12     plt.subplot(4, 4, i + 1) # Adjust subplot layout as needed
13     sns.histplot(df[col], kde=True)
14     plt.title(f'Histogram of {col}')
15
16 # Adjust layout
17 plt.tight_layout()
18 plt.show()

```



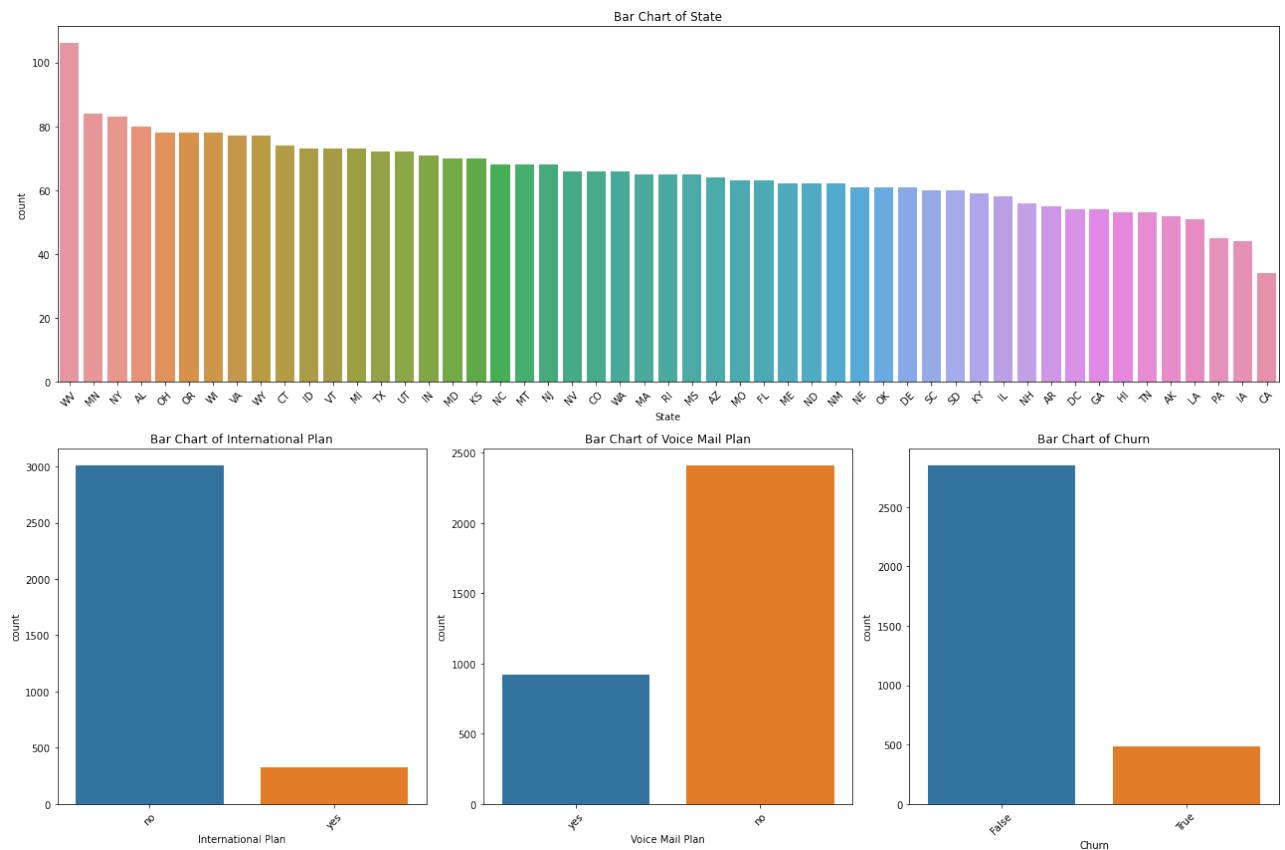
Categorical Columns

In [16]:

```

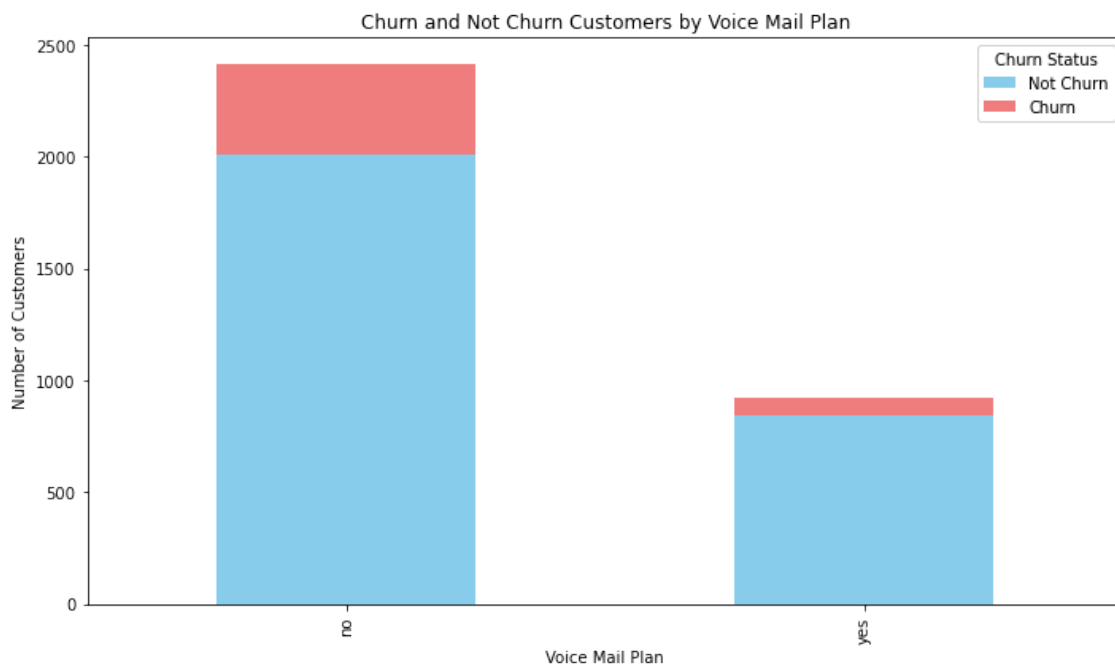
1
2 # List of categorical columns, excluding 'Phone Number'
3 categorical_cols = ['State', 'International Plan', 'Voice Mail Plan', 'Churn']
4
5 # Set the size of the overall figure
6 plt.figure(figsize=(18, 12))
7
8 # Calculate the order of states by their count in descending order
9 state_order = df['State'].value_counts().index
10
11 # Create bar plot for 'State' with a larger size and ordered
12 plt.subplot(2, 1, 1) # 2 rows, 1 column, 1st subplot
13 sns.countplot(data=df, x='State', order=state_order)
14 plt.title('Bar Chart of State')
15 plt.xticks(rotation=45) # Rotate x labels if needed
16
17 # Create smaller bar plots for the other categorical columns
18 for i, col in enumerate(categorical_cols[1:], start=1):
19     plt.subplot(2, 3, i + 3) # 2 rows, 3 columns, starting from 4th subplot
20     sns.countplot(data=df, x=col)
21     plt.title(f'Bar Chart of {col}')
22     plt.xticks(rotation=45) # Rotate x labels if needed
23
24 # Adjust layout
25 plt.tight_layout()
26 plt.show()
27

```



Bivariant Analysis

```
In [17]: 1 # Group by 'Voice Mail Plan' and 'Churn' to get the counts
2 vmail_plan_churn_counts = df.groupby(['Voice Mail Plan', 'Churn']).size().unstack().fillna(0)
3
4 # Sort the categories by the total count of customers in descending order
5 vmail_plan_churn_counts['Total'] = vmail_plan_churn_counts.sum(axis=1)
6 vmail_plan_churn_counts = vmail_plan_churn_counts.sort_values(by='Total', ascending=False).drop('Total', axis=1)
7
8 # Plotting the stacked bar plot
9 vmail_plan_churn_counts.plot(kind='bar', stacked=True, figsize=(10, 6), color=['skyblue', 'lightcoral'])
10
11 # Adding title and labels
12 plt.title('Churn and Not Churn Customers by Voice Mail Plan')
13 plt.xlabel('Voice Mail Plan')
14 plt.ylabel('Number of Customers')
15 plt.legend(['Not Churn', 'Churn'], title='Churn Status')
16
17 # Display the plot
18 plt.tight_layout()
19 plt.show()
20
```



Count of Churn Customers Per State

```
In [18]: 1 count_churn = df.groupby(["Churn", "State"]).agg(churn_count = ("Churn", "count")).reset_index()
2 count_churn.head()
```

Out[18]:

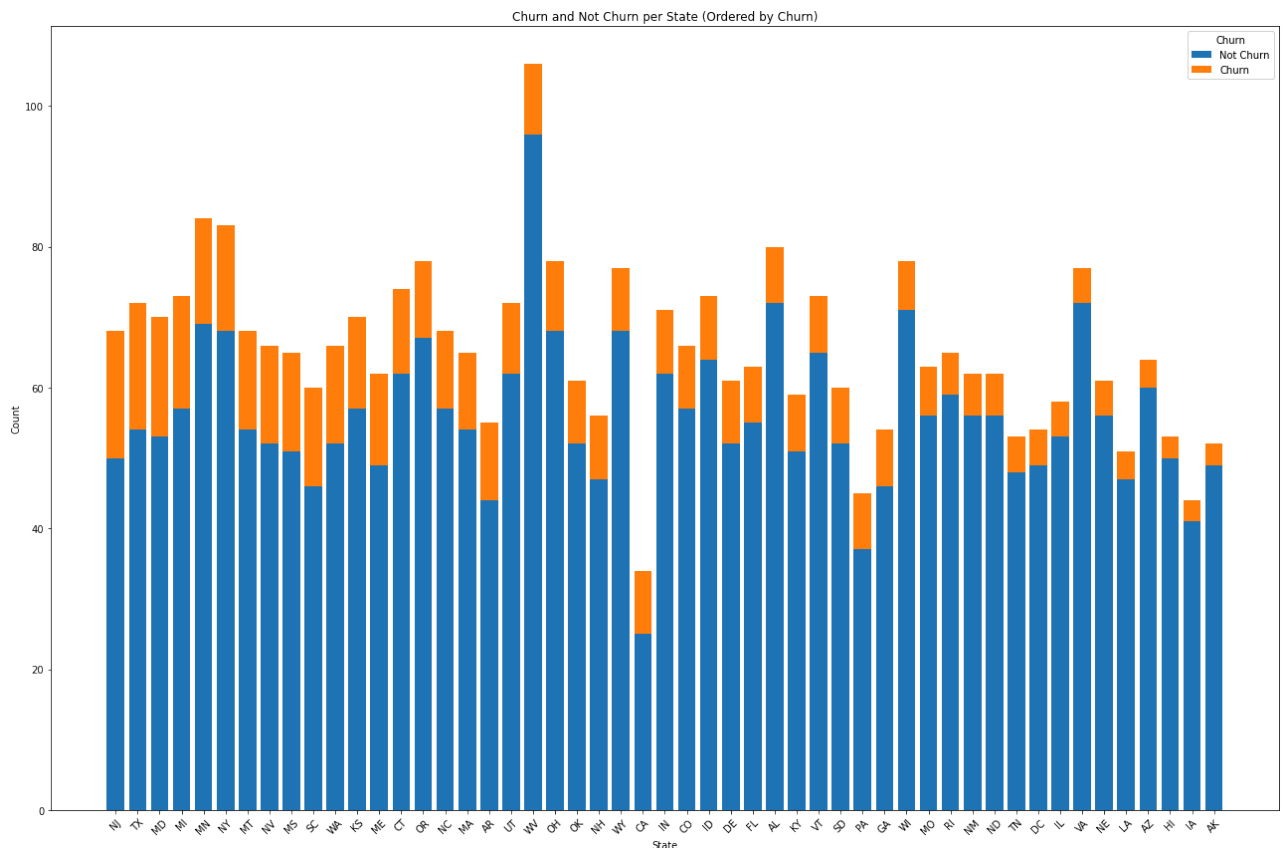
	Churn	State	churn_count
0	False	AK	49
1	False	AL	72
2	False	AR	44
3	False	AZ	60
4	False	CA	25

In [19]:

```

1 # Encode the "Churn" column to 0 and 1
2 df['Churn'] = df['Churn'].map({False: 0, True: 1})
3
4 # Group by "State" and "Churn" and count the occurrences
5 count_churn = df.groupby(['State', 'Churn']).size().unstack(fill_value=0)
6
7 # Sort the DataFrame by the number of churned customers (Churn = 1) in descending order
8 count_churn_sorted = count_churn.sort_values(by=1, ascending=False)
9
10 # Create a new column 'State' from the index to use in seaborn
11 count_churn_sorted = count_churn_sorted.reset_index()
12
13 # Plotting the stacked bar plot
14 plt.figure(figsize=(18, 12))
15 bottom = count_churn_sorted[0]
16 top = count_churn_sorted[1]
17
18 # Plot the not churned customers
19 plt.bar(count_churn_sorted['State'], bottom, label='Not Churn', color='#1f77b4')
20
21 # Plot the churned customers on top of the not churned
22 plt.bar(count_churn_sorted['State'], top, bottom=bottom, label='Churn', color='#ff7f0e')
23
24 plt.title('Churn and Not Churn per State (Ordered by Churn)')
25 plt.ylabel('Count')
26 plt.xlabel('State')
27 plt.xticks(rotation=45) # Rotate x Labels if needed
28 plt.legend(title='Churn')
29
30 # Show plot
31 plt.tight_layout()
32 plt.show()

```



Count of Churn Customers Per Area Code

In [20]:

```
1 count_area_code_sorted = df.groupby(["Churn", "Area Code"]).agg(count_area_code = ("Churn", "count"))
2 count_area_code_sorted.head(10)
```

Out[20]:

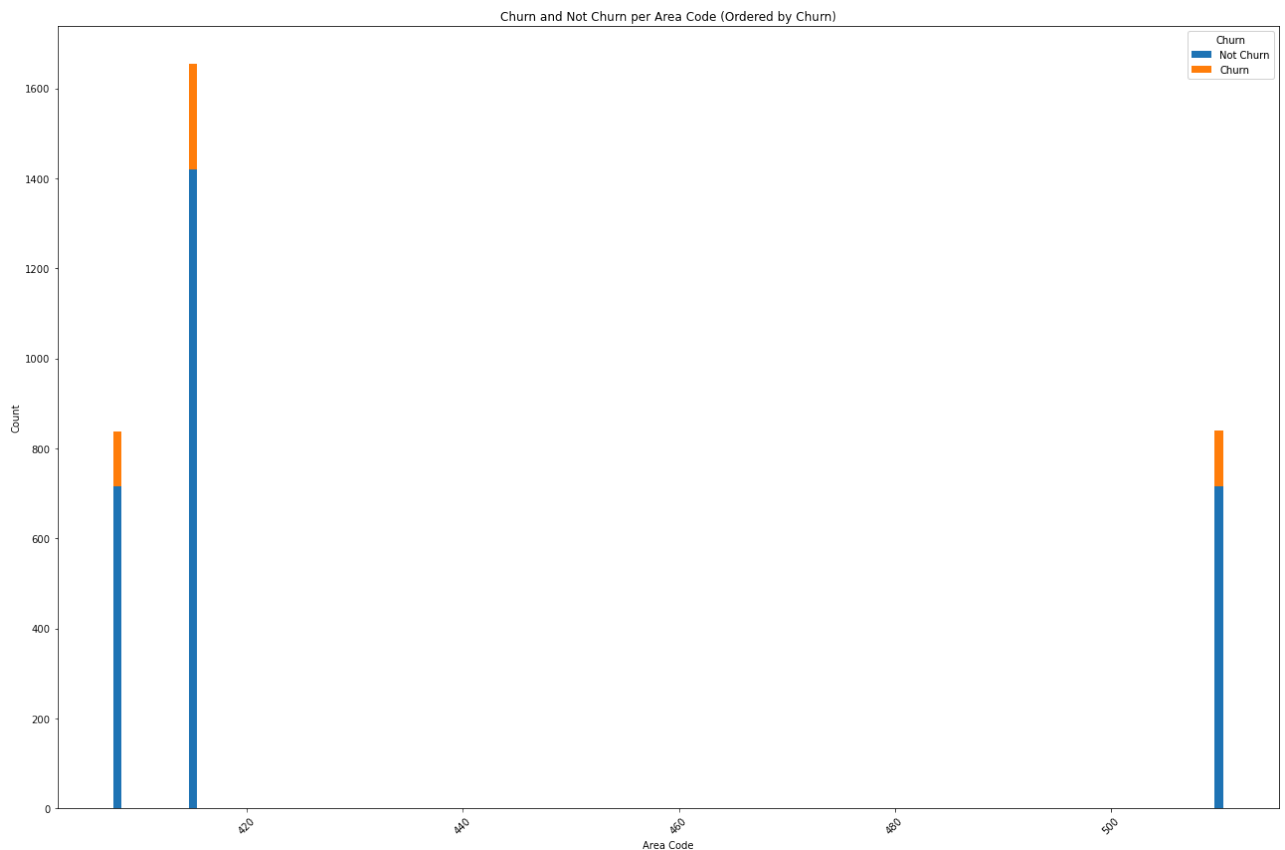
	Churn	Area Code	count_area_code
0	0	408	716
1	0	415	1419
2	0	510	715
3	1	408	122
4	1	415	236
5	1	510	125

In [21]:

```

1 # Encode the "Churn" column to 0 and 1
2 df['Churn'] = df['Churn'].map({False: 0, True: 1})
3
4 # Group by "Area Code" and "Churn" and count the occurrences
5 count_area_code_sorted = df.groupby(['Area Code', 'Churn']).size().unstack(fill_value=0)
6
7 # Sort the DataFrame by the number of churned customers (Churn = 1) in descending order
8 count_area_code_sorted_sorted = count_area_code_sorted.sort_values(by=1, ascending=False)
9
10 # Create a new column 'Area Code' from the index to use in seaborn
11 count_area_code_sorted_sorted = count_area_code_sorted_sorted.reset_index()
12
13 # Plotting the stacked bar plot
14 plt.figure(figsize=(18, 12))
15 bottom = count_area_code_sorted_sorted[0]
16 top = count_area_code_sorted_sorted[1]
17
18 # Plot the not churned customers
19 plt.bar(count_area_code_sorted_sorted['Area Code'], bottom, label='Not Churn', color='#1f77b4')
20
21 # Plot the churned customers on top of the not churned
22 plt.bar(count_area_code_sorted_sorted['Area Code'], top, bottom=bottom, label='Churn', color='#ff7f0e')
23
24 plt.title('Churn and Not Churn per Area Code (Ordered by Churn)')
25 plt.ylabel('Count')
26 plt.xlabel('Area Code')
27 plt.xticks(rotation=45) # Rotate x Labels if needed
28 plt.legend(title='Churn')
29
30 # Show plot
31 plt.tight_layout()
32 plt.show()

```



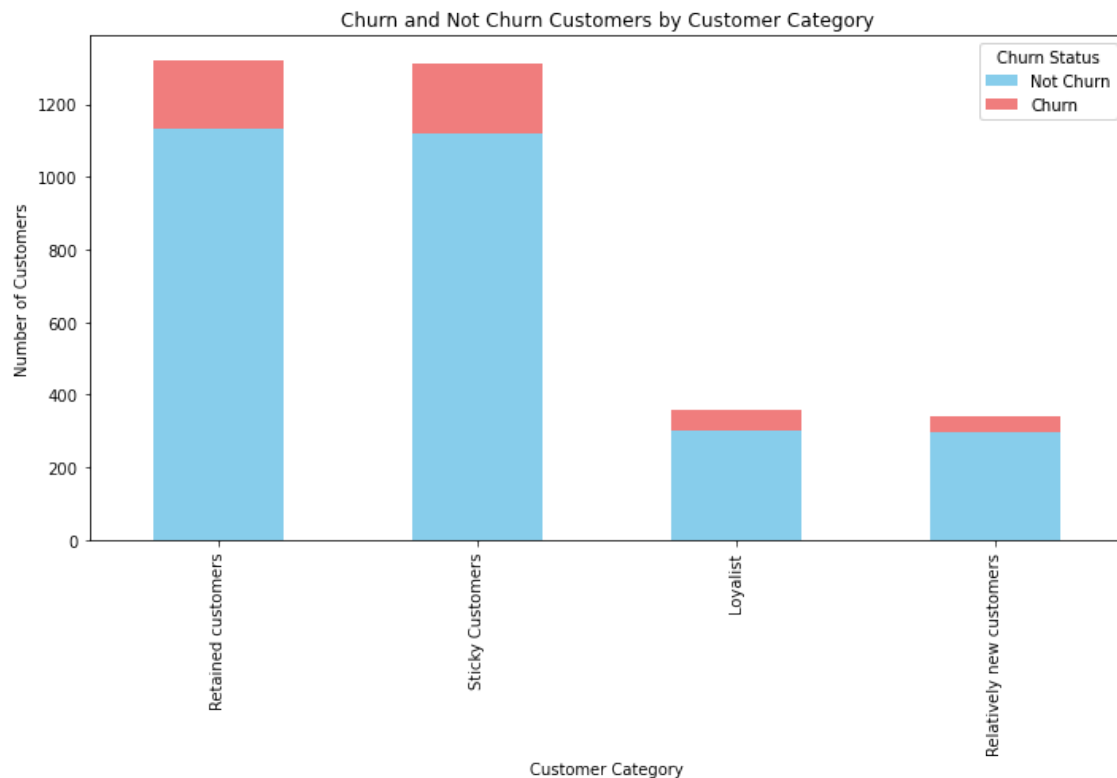
Creating Account Category Column

```
In [22]: 1 import pandas as pd
2
3 # Function to categorize customers
4 def categorize_customer(account_length):
5     if account_length <= 50:
6         return 'Relatively new customers'
7     elif 50 < account_length <= 100:
8         return 'Retained customers'
9     elif 100 < account_length <= 150:
10        return 'Sticky Customers'
11    else:
12        return 'Loyalist'
13
14 # Create the Customer Category column in df
15 df['Customer Category'] = df['Account Length'].apply(categorize_customer)
16
17 # Display the first few rows to verify
18 print(df[['Account Length', 'Customer Category']].head())
19
```

	Account Length	Customer Category
0	128	Sticky Customers
1	107	Sticky Customers
2	137	Sticky Customers
3	84	Retained customers
4	75	Retained customers

Distribution of Customer Category vs Churn

```
In [23]: 1 # Group by 'Customer Category' and 'Churn' to get the counts
2 category_churn_counts = df.groupby(['Customer Category', 'Churn']).size().unstack().fillna(0)
3
4 # Sort the categories by the total count of customers in descending order
5 category_churn_counts['Total'] = category_churn_counts.sum(axis=1)
6 category_churn_counts = category_churn_counts.sort_values(by='Total', ascending=False).drop(columns='Total')
7
8 # Plotting the stacked bar plot
9 category_churn_counts.plot(kind='bar', stacked=True, figsize=(10, 7), color=['skyblue', 'lightcoral'])
10
11 # Adding title and labels
12 plt.title('Churn and Not Churn Customers by Customer Category')
13 plt.xlabel('Customer Category')
14 plt.ylabel('Number of Customers')
15 plt.legend(['Not Churn', 'Churn'], title='Churn Status')
16
17 # Display the plot
18 plt.tight_layout()
19 plt.show()
```



Plot of Syria using the Area_Code

```
In [24]: 1 # Extract unique area codes
2 unique_area_codes = df['Area Code'].unique()
3
4 # Print all unique area codes
5 print("Unique Area Codes:")
6 for area_code in unique_area_codes:
7     print(area_code)
```

Unique Area Codes:

415
408
510

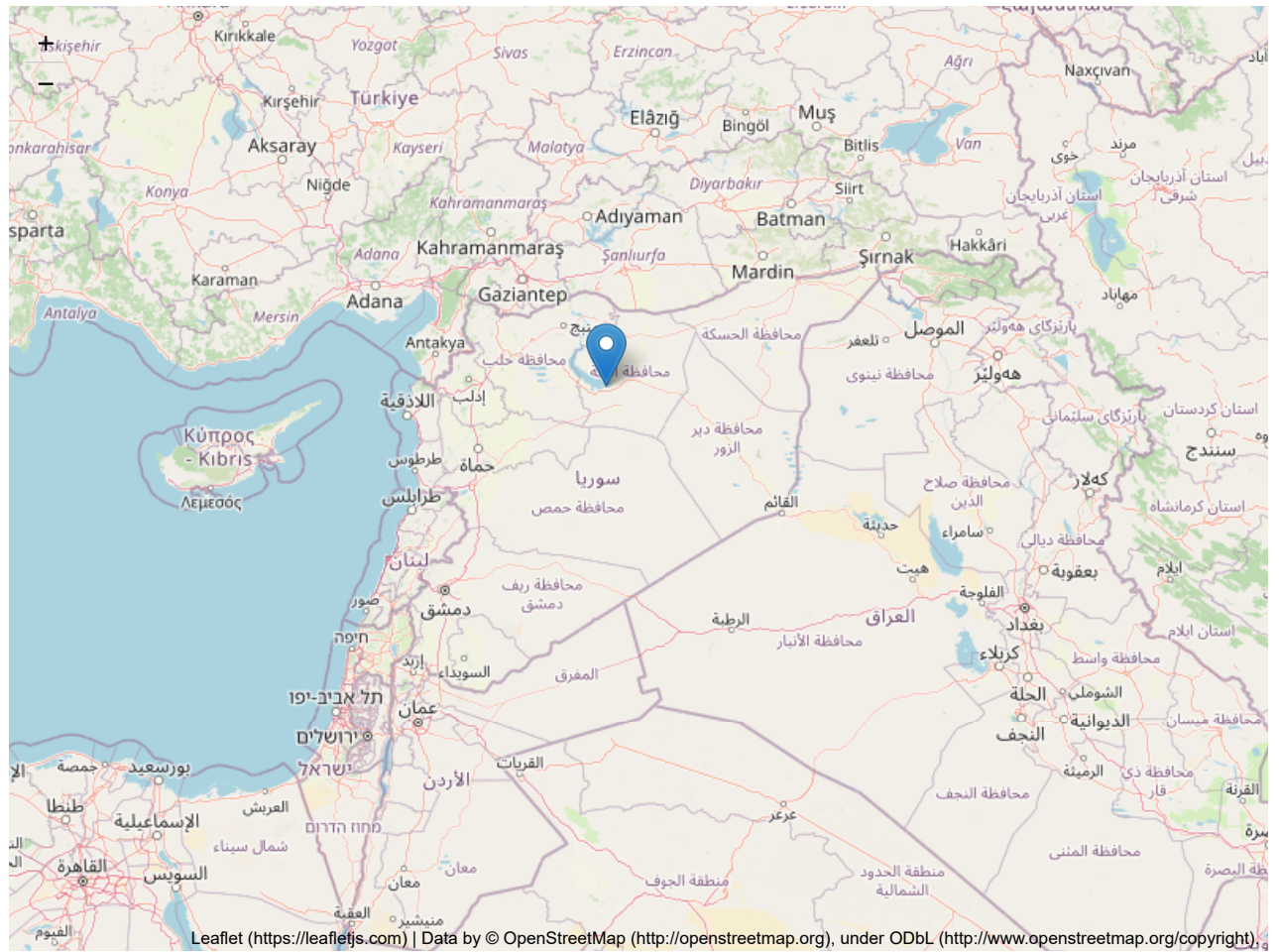
```

In [25]: 1 !pip install geopy folium
2 from geopy.geocoders import Nominatim
3 import folium
4
5 # Extract unique area codes
6 unique_area_codes = df['Area Code'].unique()
7
8 # Initialize geocoder with a longer timeout
9 geolocator = Nominatim(user_agent="area_locator", timeout=10)
10
11 # Create a map centered on Syria
12 map_syria = folium.Map(location=[34.802075, 38.996815], zoom_start=6)
13
14 # Loop through unique area codes and get coordinates
15 for area_code in unique_area_codes:
16     try:
17         location = geolocator.geocode(f"{area_code}, Syria")
18         if location:
19             # Add marker to the map
20             folium.Marker([location.latitude, location.longitude], popup=f"Area Code: {area_code}").a
21     except Exception as e:
22         print(f"Error fetching location for area code {area_code}: {e}")
23
24 # Save the map as an HTML file
25 map_syria.save("area_codes_map.html")
26
27 # Display the map in a Jupyter notebook
28 from IPython.display import IFrame
29 IFrame("area_codes_map.html", width=800, height=600)
30

```

Requirement already satisfied: geopy in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (2.4.1)
Requirement already satisfied: folium in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (0.11.0)
Requirement already satisfied: geographiclib<3,>=1.52 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from geopy) (2.0)
Requirement already satisfied: numpy in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from folium) (1.18.5)
Requirement already satisfied: branca>=0.3.0 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from folium) (0.4.1)
Requirement already satisfied: Jinja2>=2.9 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from folium) (2.11.2)
Requirement already satisfied: requests in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from folium) (2.24.0)
Requirement already satisfied: MarkupSafe>=0.23 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from Jinja2>=2.9->folium) (1.1.1)
Requirement already satisfied: urllib3!=1.25.0,!>=1.25.1,<1.26,>=1.21.1 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from requests->folium) (1.25.10)
Requirement already satisfied: chardet<4,>=3.0.2 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from requests->folium) (3.0.4)
Requirement already satisfied: certifi>=2017.4.17 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from requests->folium) (2020.6.20)
Requirement already satisfied: idna<3,>=2.5 in c:\users\user\anaconda3\envs\learn-env\lib\site-packages (from requests->folium) (2.10)

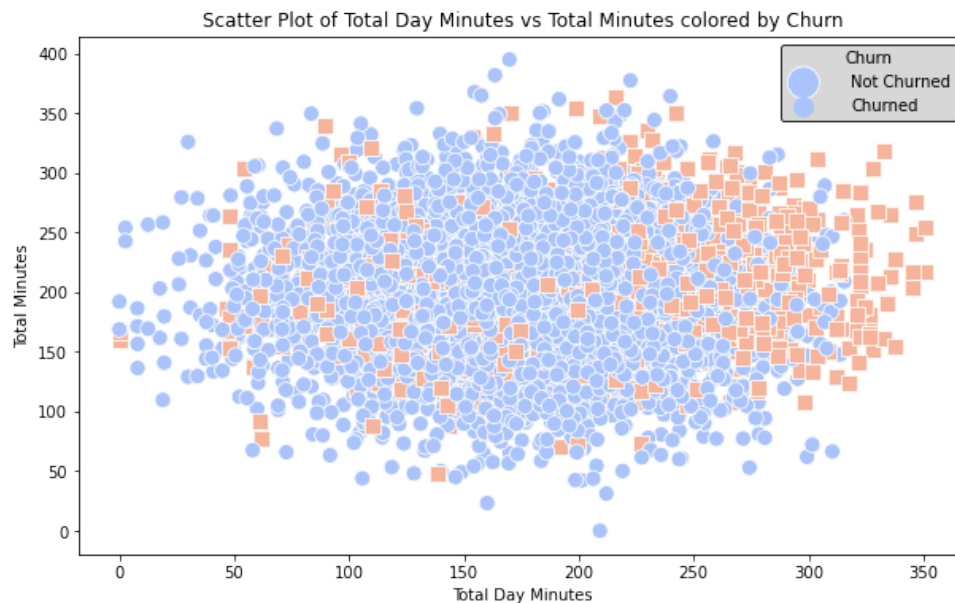
Out[25]:



Multivariant Analysis

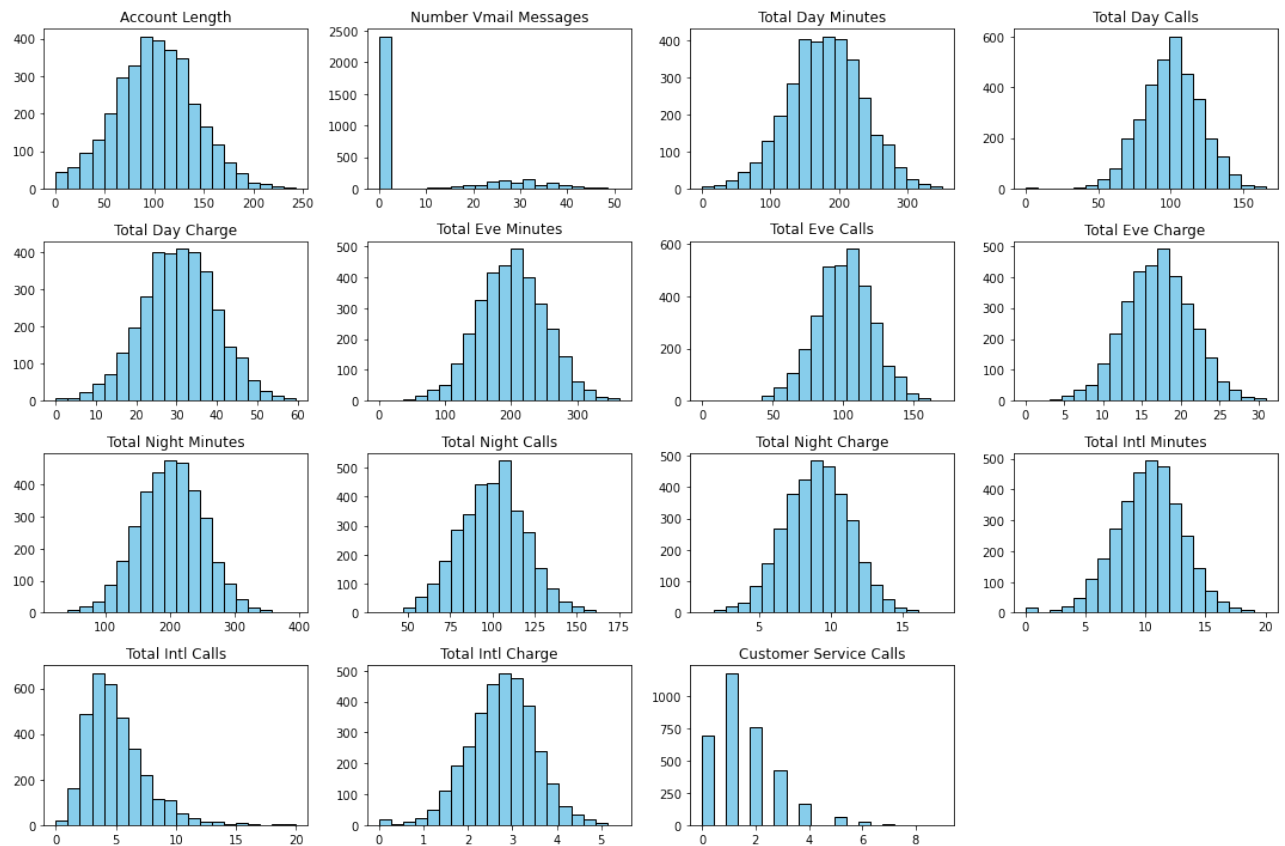
In [26]:

```
1 # Scatter plot of Total Day Minutes vs Total Eve Minutes vs Total Night Minutes colored by Churn
2 plt.figure(figsize=(10, 6))
3 sns.scatterplot(x='Total Day Minutes', y='Total Eve Minutes', data=df, hue='Churn', style='Churn', palette='coolwarm')
4 sns.scatterplot(x='Total Day Minutes', y='Total Night Minutes', data=df, hue='Churn', style='Churn', palette='coolwarm')
5 plt.xlabel('Total Day Minutes')
6 plt.ylabel('Total Minutes')
7 plt.title('Scatter Plot of Total Day Minutes vs Total Minutes colored by Churn')
8 legend_colors = {'Churned': sns.color_palette('coolwarm')[1], 'Not Churned': sns.color_palette('coolwarm')[0]}
9 plt.legend(title='Churn', loc='upper right', labels=['Not Churned', 'Churned'], facecolor='lightgrey')
10 plt.show()
11
```



Data Preparation

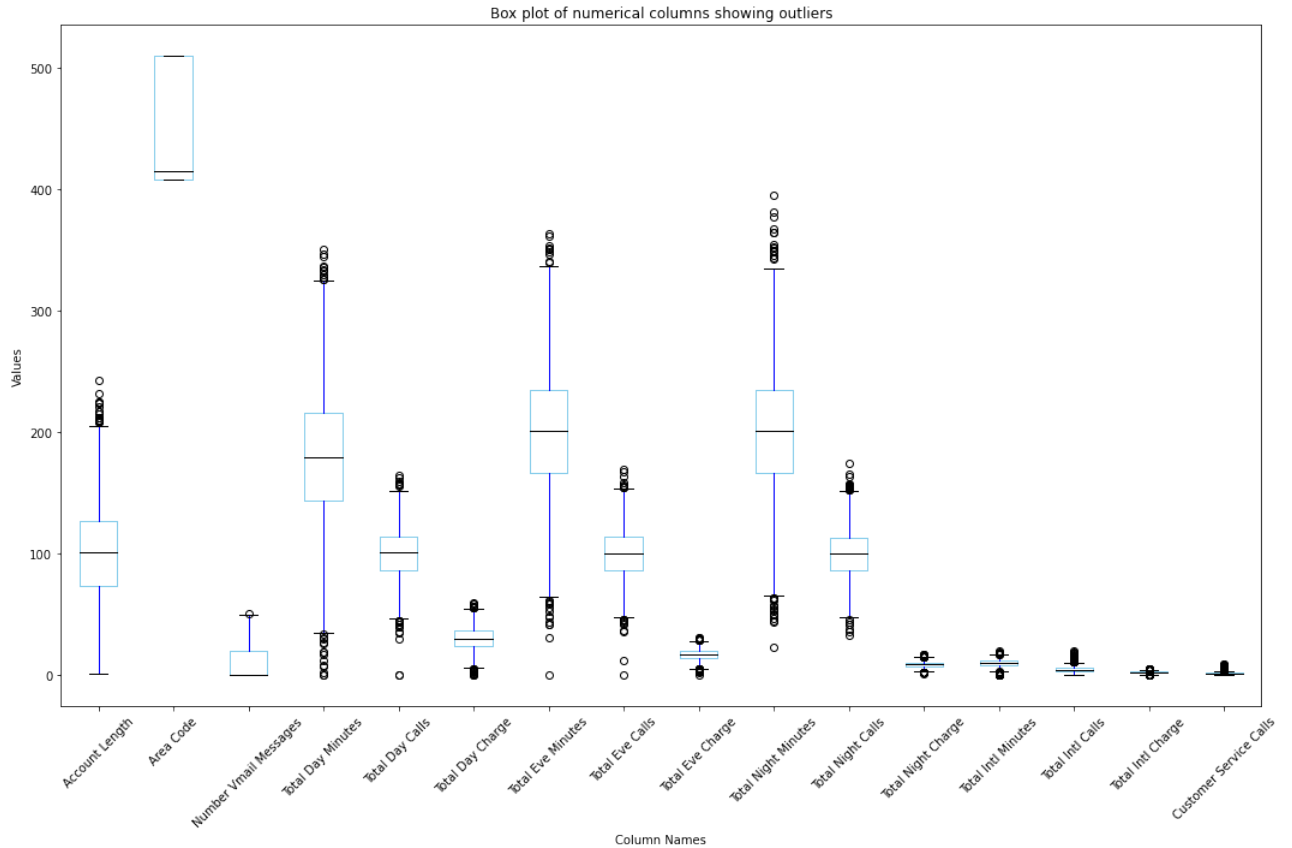
```
In [27]: 1 # Select numerical columns excluding 'Area Code'
2 numerical_cols = ['Account Length', 'Number Vmail Messages', 'Total Day Minutes',
3                 'Total Day Calls', 'Total Day Charge', 'Total Eve Minutes', 'Total Eve Calls',
4                 'Total Eve Charge', 'Total Night Minutes', 'Total Night Calls',
5                 'Total Night Charge', 'Total Intl Minutes', 'Total Intl Calls',
6                 'Total Intl Charge', 'Customer Service Calls']
7
8 # Plot histograms before handling outliers
9 plt.figure(figsize=(15, 10))
10 for i, col in enumerate(numerical_cols, 1):
11     plt.subplot(4, 4, i)
12     plt.hist(df[col], bins=20, color='skyblue', edgecolor='black')
13     plt.title(col)
14 plt.tight_layout()
15 plt.show()
```



```

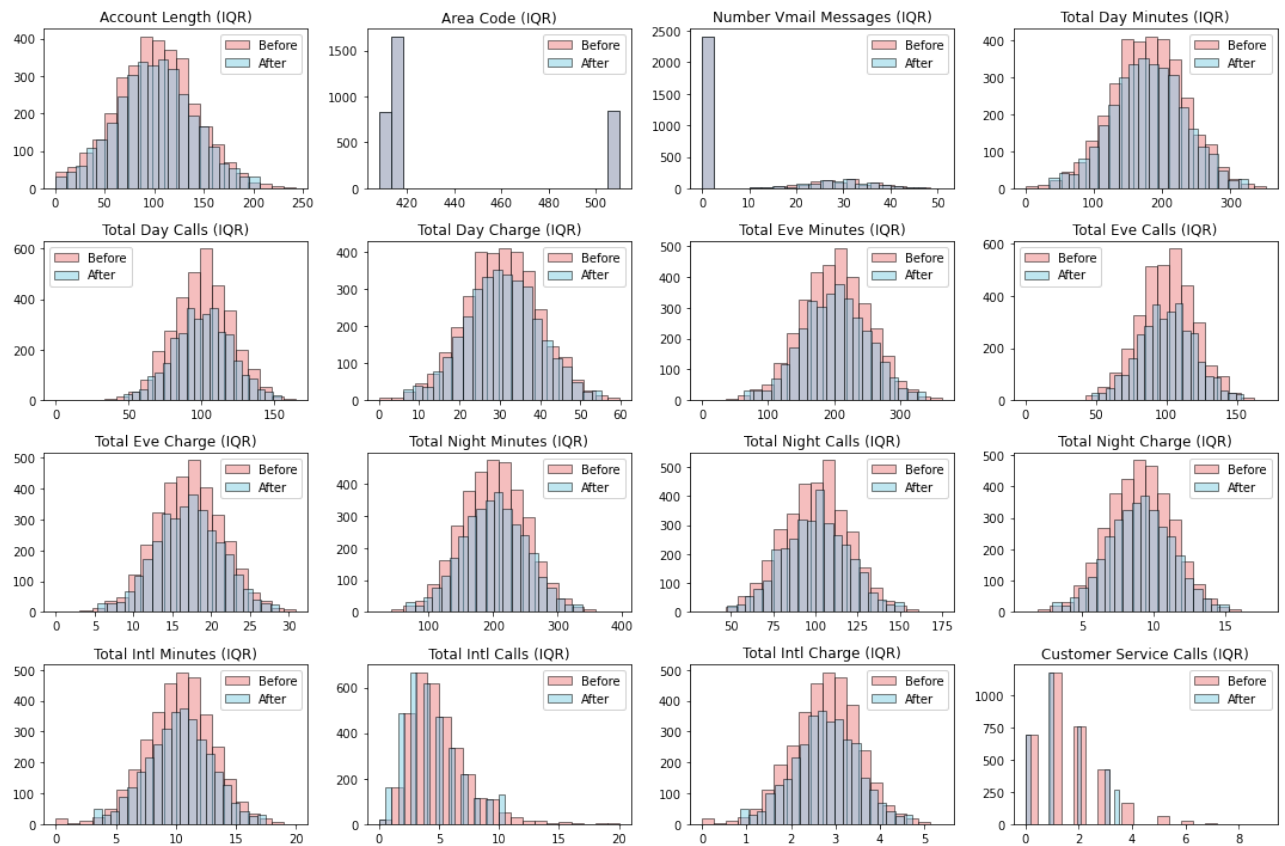
In [28]: 1 # List of numerical columns excluding 'Churn'
2 numerical_cols = df.select_dtypes(include=['int64', 'float64']).columns.drop('Churn')
3
4 # Create box plots for numerical columns showing outliers
5 plt.figure(figsize=(15, 10))
6 df[numerical_cols].boxplot(grid=False, color=dict(boxes='skyblue', whiskers='blue', medians='black'),
7 plt.title('Box plot of numerical columns showing outliers')
8 plt.xlabel('Column Names')
9 plt.ylabel('Values')
10 plt.xticks(rotation=45)
11 plt.tight_layout()
12 plt.show()
13

```



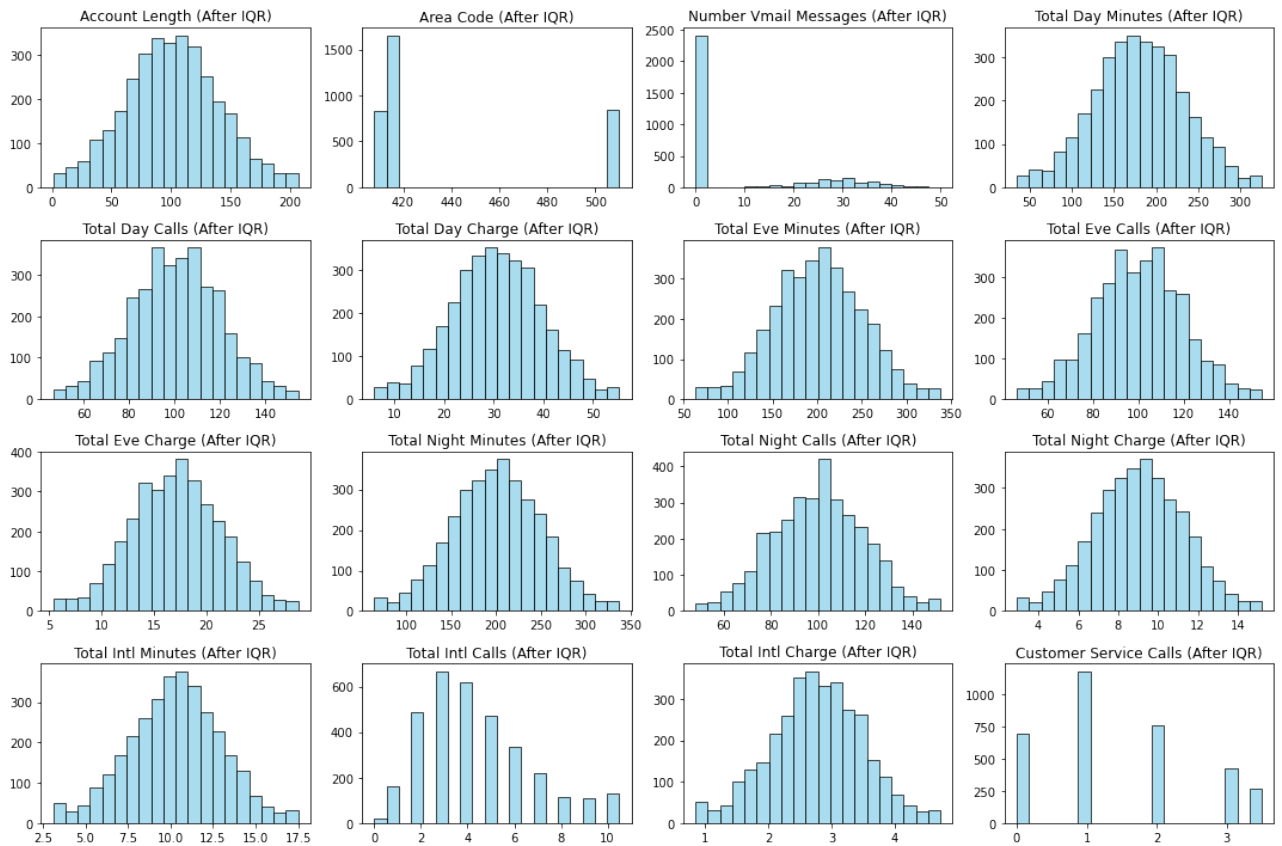
Handling Outliers

```
In [29]: 1 # Handling outliers using IQR
2 df_no_outliers = df.copy()
3 for col in numerical_cols:
4     Q1 = df[col].quantile(0.25)
5     Q3 = df[col].quantile(0.75)
6     IQR = Q3 - Q1
7     lower_bound = Q1 - 1.5 * IQR
8     upper_bound = Q3 + 1.5 * IQR
9     df_no_outliers[col] = df_no_outliers[col].clip(lower=lower_bound, upper=upper_bound)
10
11 # Plot histograms after handling outliers with different colors
12 plt.figure(figsize=(15, 10))
13 for i, col in enumerate(numerical_cols, 1):
14     plt.subplot(4, 4, i)
15     plt.hist(df[col], bins=20, color='lightcoral', edgecolor='black', alpha=0.5, label='Before')
16     plt.hist(df_no_outliers[col], bins=20, color='skyblue', edgecolor='black', alpha=0.5, label='After')
17     plt.title(col + ' (IQR)')
18     plt.legend()
19 plt.tight_layout()
20 plt.show()
```



Plot After Handling Outliers

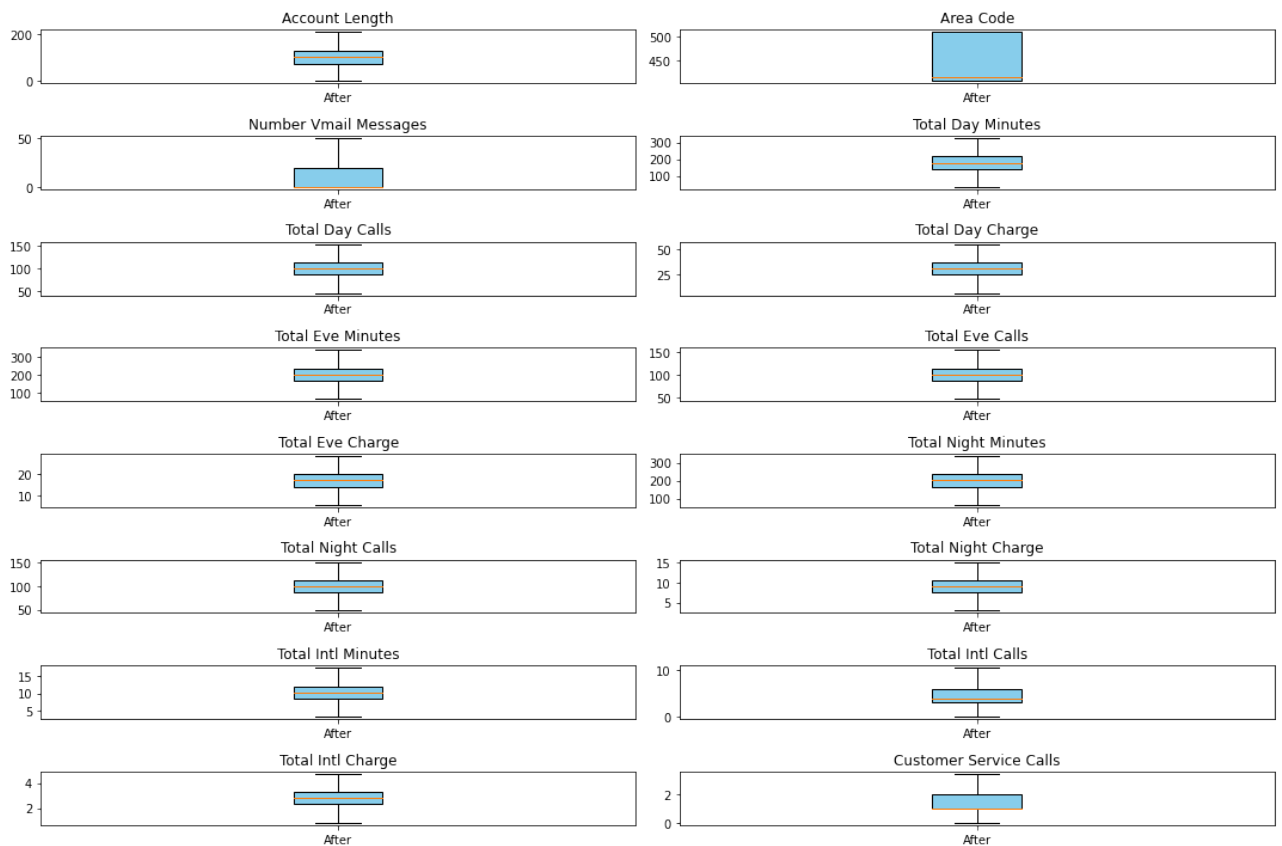
```
In [30]: 1 # List of numerical columns excluding 'Churn'
2 numerical_cols = df.select_dtypes(include=['int64', 'float64']).columns.drop('Churn')
3
4 # Handling outliers using IQR
5 df_no_outliers = df.copy()
6 for col in numerical_cols:
7     Q1 = df[col].quantile(0.25)
8     Q3 = df[col].quantile(0.75)
9     IQR = Q3 - Q1
10    lower_bound = Q1 - 1.5 * IQR
11    upper_bound = Q3 + 1.5 * IQR
12    df_no_outliers[col] = df_no_outliers[col].clip(lower=lower_bound, upper=upper_bound)
13
14 # Plot histograms for numerical columns after handling outliers
15 plt.figure(figsize=(15, 10))
16 for i, col in enumerate(numerical_cols, 1):
17     plt.subplot((len(numerical_cols) + 3) // 4, 4, i) # Adjust subplot grid size dynamically
18     plt.hist(df_no_outliers[col], bins=20, color='skyblue', edgecolor='black', alpha=0.7)
19     plt.title(col + ' (After IQR)')
20 plt.tight_layout()
21 plt.show()
22
```



```

In [31]: 1 # List of numerical columns excluding 'Churn'
2 numerical_cols = df.select_dtypes(include=['int64', 'float64']).columns.drop('Churn')
3
4 # Handling outliers using IQR
5 df_no_outliers = df.copy()
6 for col in numerical_cols:
7     Q1 = df[col].quantile(0.25)
8     Q3 = df[col].quantile(0.75)
9     IQR = Q3 - Q1
10    lower_bound = Q1 - 1.5 * IQR
11    upper_bound = Q3 + 1.5 * IQR
12    df_no_outliers[col] = df_no_outliers[col].clip(lower=lower_bound, upper=upper_bound)
13
14 # Create box plots for each numerical variable after outlier removal
15 plt.figure(figsize=(15, 10))
16 for i, col in enumerate(numerical_cols, 1):
17     plt.subplot((len(numerical_cols) + 1) // 2, 2, i) # Adjust subplot grid size dynamically
18     plt.boxplot(df_no_outliers[col], labels=['After'], patch_artist=True, boxprops=dict(facecolor='skyblue', medianprops=dict(color='red')))
19     plt.title(col)
20 plt.tight_layout()
21 plt.show()

```



Normality and Spread of the Cleaned Dataset.

```
In [32]: 1 # Computing Normality and Spread of the Cleaned Dataset.
2 import scipy.stats as stats
3
4 # Define the numerical columns
5 numerical_cols = ['Account Length', 'Number Vmail Messages', 'Total Day Minutes',
6                  'Total Day Calls', 'Total Day Charge', 'Total Eve Minutes',
7                  'Total Eve Calls', 'Total Eve Charge', 'Total Night Minutes',
8                  'Total Night Calls', 'Total Night Charge', 'Total Intl Minutes',
9                  'Total Intl Calls', 'Total Intl Charge', 'Customer Service Calls']
10
11 # Compute normality and spread for each numerical column
12 normality_spread = {}
13
14 for column in numerical_cols:
15     # Calculate skewness
16     skewness = stats.skew(df[column])
17
18     # Calculate kurtosis
19     kurtosis = stats.kurtosis(df[column])
20
21     # Calculate mean
22     mean = df[column].mean()
23
24     # Calculate median
25     median = df[column].median()
26
27     # Calculate standard deviation
28     std_dev = df[column].std()
29
30     # Store results in a dictionary
31     normality_spread[column] = {'Skewness': skewness, 'Kurtosis': kurtosis,
32                                'Mean': mean, 'Median': median, 'Std Dev': std_dev}
33
34 # Display results
35 import pandas as pd
36 normality_spread_df = pd.DataFrame(normality_spread).T
37 print(normality_spread_df)
38
```

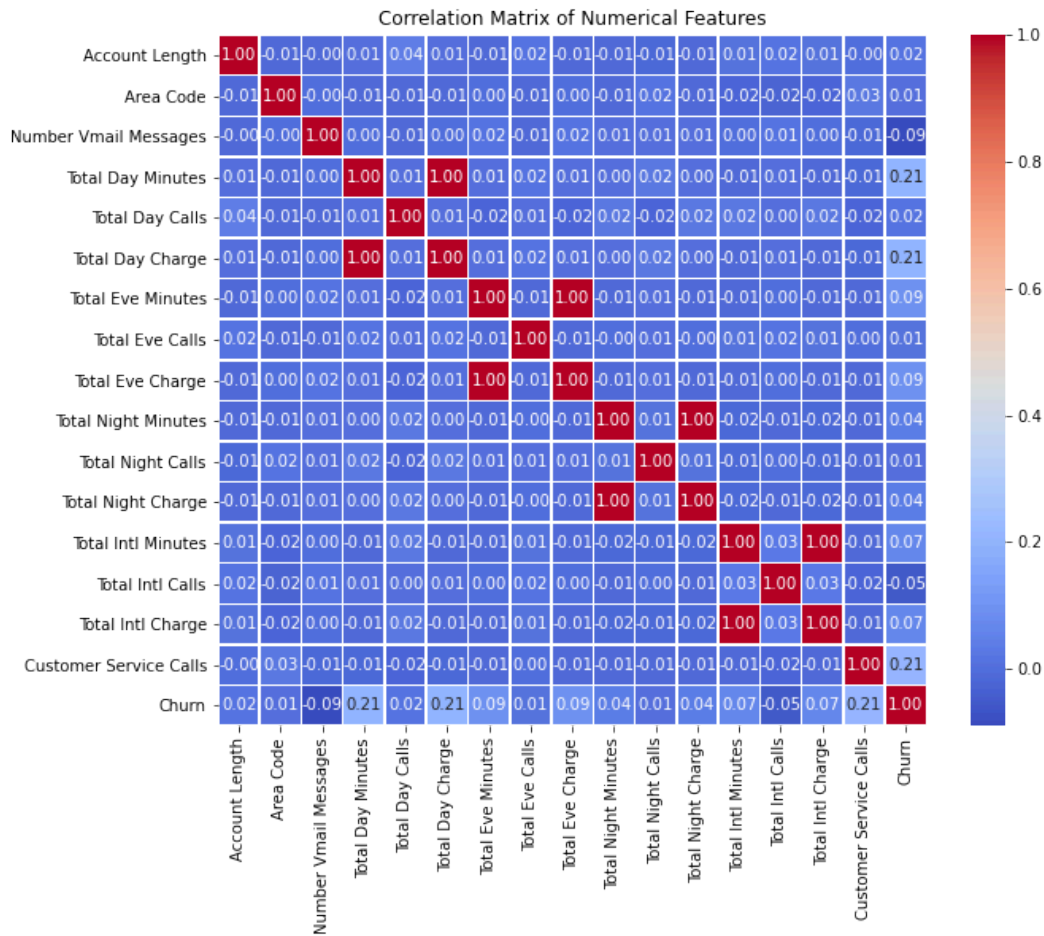
	Skewness	Kurtosis \
Account Length	0.09656281161489656	-0.1094739184341575
Number Vmail Messages	1.2642543349768245	-0.0528515105905245
Total Day Minutes	-0.0290639795181198	-0.0217101179240888
Total Day Calls	-0.1117363237307519	0.24101722895174227
Total Day Charge	-0.0290701779270378	-0.0215817191450336
Total Eve Minutes	-0.0238667088046375	0.0237916804447047
Total Eve Calls	-0.0555381300016192	0.20404769217448226
Total Eve Charge	-0.023847250496277	0.02364954586272594
Total Night Minutes	0.008917275580987895	0.08388775499253365
Total Night Calls	0.03248494205404463	-0.0737112242125884
Total Night Charge	0.008882237062694412	0.08373508611499814
Total Intl Minutes	-0.2450256034866443	0.606471635404318
Total Intl Calls	1.3208833668164015	3.07716543898885142
Total Intl Charge	-0.2451761045009844	0.6068966666527675
Customer Service Calls	1.09086826017550109	1.7265184753957081

	Mean	Median \
Account Length	101.06480648064805905	101.0
Number Vmail Messages	8.0990099009900991	0.0
Total Day Minutes	179.7750975097509354	179.400000000000057
Total Day Calls	100.4356435643564396	101.0
Total Day Charge	30.562307230723075	30.5
Total Eve Minutes	200.9803480348034839	201.400000000000057
Total Eve Calls	100.11431143114310771	100.0
Total Eve Charge	17.08354035403540294	17.120000000000001
Total Night Minutes	200.8720372037203674	201.19999999999998863
Total Night Calls	100.1077107710771088	100.0
Total Night Charge	9.03932493249324942	9.05000000000000071
Total Intl Minutes	10.23729372937293824	10.3000000000000007
Total Intl Calls	4.4794479447944795	4.0
Total Intl Charge	2.7645814581458144	2.7799999999999998
Customer Service Calls	1.5628562856285628	1.0

	Std Dev
Account Length	39.8221059285956045
Number Vmail Messages	13.6883653720385983
Total Day Minutes	54.46738920237137194
Total Day Calls	20.0690842073008966
Total Day Charge	9.2594345539305003
Total Eve Minutes	50.7138444258119989
Total Eve Calls	19.9226252939431028
Total Eve Charge	4.31066764311034056
Total Night Minutes	50.5738470136583587
Total Night Calls	19.5686093460585582
Total Night Charge	2.275872837660029
Total Intl Minutes	2.791839548408416
Total Intl Calls	2.461214270546094
Total Intl Charge	0.753772612663046
Customer Service Calls	1.3154910448664767

Correlation Matrix

```
In [33]: 1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Select numerical columns
5 numerical_columns = df.select_dtypes(include=['int64', 'float64']).columns
6
7 # Calculate correlation matrix
8 correlation_matrix = df[numerical_columns].corr()
9
10 # Plot heatmap
11 plt.figure(figsize=(10, 8))
12 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=0.5)
13 plt.title('Correlation Matrix of Numerical Features')
14 plt.show()
```



Identifying Multicollinearity

```
In [34]: 1 # Find columns with correlation greater than 0.95
2 high_corr_pairs = []
3 threshold = 0.95
4 for i in range(len(correlation_matrix.columns)):
5     for j in range(i):
6         if abs(correlation_matrix.iloc[i, j]) > threshold:
7             colname_i = correlation_matrix.columns[i]
8             colname_j = correlation_matrix.columns[j]
9             high_corr_pairs.append((colname_i, colname_j, correlation_matrix.iloc[i, j]))
10
11 # Print high correlation pairs
12 print("Pairs of columns with correlation greater than 0.95:")
13 for pair in high_corr_pairs:
14     print(f"{pair[0]} and {pair[1]} with correlation of {pair[2]:.2f}")
15
```

Pairs of columns with correlation greater than 0.95:
 Total Day Charge and Total Day Minutes with correlation of 1.00
 Total Eve Charge and Total Eve Minutes with correlation of 1.00
 Total Night Charge and Total Night Minutes with correlation of 1.00
 Total Intl Charge and Total Intl Minutes with correlation of 1.00

```
In [35]: 1 # identifying Variable Inflation Factor
2
3 from statsmodels.stats.outliers_influence import variance_inflation_factor
4
5 # Select numerical columns
6 numerical_columns = df.select_dtypes(include=['int64', 'float64'])
7
8 # Calculate VIF for each numerical feature
9 vif_data = pd.DataFrame()
10 vif_data['Feature'] = numerical_columns.columns
11 vif_data['VIF'] = [variance_inflation_factor(numerical_columns.values, i) for i in range(numerical_co
12
13 # Display VIF values
14 print(vif_data)
15
```

	Feature	VIF
0	Account Length	7.2931644646281919
1	Area Code	61.4061122968105906
2	Number Vmail Messages	1.36362106926135174
3	Total Day Minutes	124603601.7811902016401291
4	Total Day Calls	23.619588049070714
5	Total Day Charge	124608062.9765620976686478
6	Total Eve Minutes	37418407.5638073161244392
7	Total Eve Calls	23.767320310736217
8	Total Eve Charge	37419731.08176666498184204
9	Total Night Minutes	10719732.5128282140940428
10	Total Night Calls	24.6165329718864214
11	Total Night Charge	10719379.5142267607152462
12	Total Intl Minutes	997547.71468332153745
13	Total Intl Calls	4.29240217134540636
14	Total Intl Charge	997925.3157051414018497
15	Customer Service Calls	2.5162522892809012
16	Churn	1.30827768455642945

```
In [36]: 1 # Numerical_columns contains the names of numerical columns
2 numerical_columns = df.select_dtypes(include=['int64', 'float64']).columns
3
4 # Compute the correlation matrix
5 correlation_matrix = df[numerical_columns].corr()
6
7 # Set the diagonal and Lower triangle to NaN (to ignore self-correlation and duplicate pairs)
8 mask = np.triu(np.ones_like(correlation_matrix, dtype=bool))
9 correlation_matrix_masked = correlation_matrix.mask(mask)
10
11 # Find pairs with correlation greater than 0.95
12 high_correlation_pairs = correlation_matrix_masked.stack().reset_index()
13 high_correlation_pairs.columns = ['Feature 1', 'Feature 2', 'Correlation']
14 high_correlation_pairs = high_correlation_pairs[high_correlation_pairs['Correlation'] > 0.95]
15
16 print(high_correlation_pairs)
17
```

	Feature 1	Feature 2	Correlation
13	Total Day Charge	Total Day Minutes	0.999999952190397
34	Total Eve Charge	Total Eve Minutes	0.9999997760198517
64	Total Night Charge	Total Night Minutes	0.99999921487588
103	Total Intl Charge	Total Intl Minutes	0.9999927417510258

```
In [37]: 1 # Domain Knowledge indicate that some of above features are expected to show high multicollinearity given
2
3 # We therefore fail to drop them and investigate further in the analysis how they contribute to Churn
```

Hypothesis Testing

Null Hypothesis (H0): There is no significant influence of the various factors to churn rate in SyriaTel.

Alternate Hypothesis (H1): There is a significant influence of the various factors to churn rate in SyriaTel.

```
In [38]: 1 from scipy.stats import f_oneway
2
3 # Numerical columns
4 numerical_columns = ['Account Length', 'Area Code', 'Number Vmail Messages', 'Total Day Minutes',
5                     'Total Day Calls', 'Total Day Charge', 'Total Eve Minutes', 'Total Eve Calls',
6                     'Total Eve Charge', 'Total Night Minutes', 'Total Night Calls', 'Total Night Charge',
7                     'Total Intl Minutes', 'Total Intl Calls', 'Total Intl Charge', 'Customer Service Calls']
8
9 # Create an empty DataFrame to store the results
10 anova_results = pd.DataFrame(columns=['Feature', 'F-Statistic', 'p-value'])
11
12 # Perform ANOVA test for each numerical column
13 for column in numerical_columns:
14     # Group the data by the 'Churn' column
15     groups = [df[column][df['Churn'] == churn] for churn in df['Churn'].unique()]
16
17     # Perform ANOVA
18     f_statistic, p_value = f_oneway(*groups)
19
20     # Append results to DataFrame
21     anova_results = anova_results.append({'Feature': column, 'F-Statistic': f_statistic, 'p-value': p_value,
22                                         ignore_index=True})
23
24 # Print the DataFrame
25 print(anova_results)
```

	Feature	F-Statistic	p-value
0	Account Length	0.9115981986407352	0.3397600070569128
1	Area Code	0.12698640858136082	0.7215998968016037
2	Number Vmail Messages	27.035911709557691296	0.00000021175218402696
3	Total Day Minutes	146.35078521943776764	0.000000000000000000
4	Total Day Calls	1.13541242989728808	0.28670102402414055
5	Total Day Charge	146.35065699096048775	0.000000000000000000
6	Total Eve Minutes	28.9325766446506485	0.0000000801133856128
7	Total Eve Calls	0.2839943754492388	0.5941305829778143
8	Total Eve Charge	28.926443755197127	0.0000000803652422777
9	Total Night Minutes	4.20149555022397259	0.0404664846378868
10	Total Night Calls	0.12563131916004017	0.7230277872159787
11	Total Night Charge	4.2021362787384957	0.04045121876901292
12	Total Intl Minutes	15.5834679864501915	0.0000805731126549902
13	Total Intl Calls	9.3279453654346529	0.002274701409848483
14	Total Intl Charge	15.5925806081700724	0.0000801875358306397
15	Customer Service Calls	151.7670126303964366	0.000000000000000000

Conclusion

Features such as 'Account Length', 'Area Code', 'Total Day Calls', 'Total Eve Calls', and 'Total Night Calls' have their p-values are greater than the significance level of 0.05. Therefore, we fail to reject the null hypothesis (H0) for these features. This suggests that there is no significant influence of these factors on the churn rate in SyriaTel.

Remaining features including 'Number Vmail Messages', 'Total Day Minutes', 'Total Day Charge', 'Total Eve Minutes', 'Total Eve Charge', 'Total Night Minutes', 'Total Night Charge', 'Total Intl Minutes', 'Total Intl Calls', 'Total Intl Charge', and 'Customer Service Calls', the p-values are extremely low (close to 0). Therefore, we reject the null hypothesis (H0) for these features. This indicates that there is a significant influence of these factors on the churn rate in SyriaTel.

In conclusion, there is evidence to suggest that most numerical features have a significant influence on the churn rate in SyriaTel, except for 'Account Length', 'Area Code', 'Total Day Calls', 'Total Eve Calls', and 'Total Night Calls'

Hot One Encoding Categorical Columns

```
In [39]: 1 #df_no_outliers and categorical_cols, numerical_cols are already defined
2
3 # Perform one-hot encoding on categorical columns excluding 'Churn'
4 df_onehot = pd.get_dummies(df_no_outliers[categorical_cols], drop_first=True)
5
6 # Ensure we are not including the 'Churn' column twice
7 # First add numerical columns
8 df_encoded = pd.concat([df_onehot, df_no_outliers[numerical_cols]], axis=1)
9
10 # Add the 'Churn' column separately to ensure it's included only once
11 df_encoded['Churn'] = df_no_outliers['Churn'].astype(int)
12
13 # Display the first few rows of the encoded DataFrame
14 #print(df_encoded.head())
15
```

```
In [40]: 1 df_encoded.shape
```

```
Out[40]: (3333, 68)
```

```
In [41]: 1 df_encoded.duplicated().sum()
```

```
Out[41]: 0
```

Normalization of the Clean Dataset

Normalizing df_encoded dataset since it appears in different scales for following reasons:

1. Handling the data appropriately
2. Ease interpretation of the subsequent models.
3. Reduce the impact of multicollinearity on the regression coefficients and their interpretability.

Standardization / normalization of the data results in a mean of zero and a standard deviation of 1

```
In [42]: 1 from sklearn.preprocessing import MinMaxScaler
2
3 # Select numerical columns
4 numerical_columns = df_encoded.select_dtypes(include=['int64', 'float64']).columns
5
6 # Initialize the MinMaxScaler
7 scaler = MinMaxScaler()
8
9 # Fit and transform the numerical columns
10 df_encoded[numerical_columns] = scaler.fit_transform(df_encoded[numerical_columns])
11
12 # Display the first few rows of the DataFrame after normalization
13 #print(df_encoded.head())
```

Modeling

In the following session, various models such as Logistic Regression, Decision Tree, KNN and XGBoost models have been built to answer the research questions.

Splitting the dataframe

```
In [43]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3
4 # Ensure X and y are defined correctly
5 X = df_encoded.drop('Churn', axis=1)
6 y = df_encoded['Churn']
7
8 # Verify the shape of y to confirm it's a one-dimensional array
9 #print(y.shape)
10
11 # Split the data into training and testing sets
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=1, stratify=y)
```

```
In [44]: 1 print("X_train shape:", X_train.shape)
2 print("X_test shape:", X_test.shape)
3 print("y_train shape:", y_train.shape)
4 print("y_test shape:", y_test.shape)
```

```
X_train shape: (2666, 67)
X_test shape: (667, 67)
y_train shape: (2666,)
y_test shape: (667,)
```

1. Baseline Mode : Logistic Regression model

```
In [45]: 1 from sklearn.linear_model import LogisticRegression
2
3 # Initialize and fit the logistic regression model
4 logreg = LogisticRegression(fit_intercept=False, C=1e12, solver='liblinear')
5 model_log = logreg.fit(X_train, y_train)
6
7 print(model_log)
```

```
LogisticRegression(C=1000000000000.0, fit_intercept=False, solver='liblinear')
```

Performance in the training data

```
In [46]: 1 y_hat_train = logreg.predict(X_train)
2
3 train_residuals = np.abs(y_train - y_hat_train)
4 print(pd.Series(train_residuals, name="Residuals (counts)").value_counts())
5 print()
6 print(pd.Series(train_residuals, name="Residuals (proportions)").value_counts(normalize=True))
```

```
0    2309
1     357
Name: Residuals (counts), dtype: int64
```

```
0    0.8660915228807202
1    0.1339084771192798
Name: Residuals (proportions), dtype: float64
```

Performance on Test Data

```
In [47]: 1 y_hat_test = logreg.predict(X_test)
2
3 test_residuals = np.abs(y_test - y_hat_test)
4 print(pd.Series(test_residuals, name="Residuals (counts)").value_counts())
5 print()
6 print(pd.Series(test_residuals, name="Residuals (proportions)").value_counts(normalize=True))

0    565
1    102
Name: Residuals (counts), dtype: int64

0    0.8470764617691154
1    0.15292353823088456
Name: Residuals (proportions), dtype: float64
```

Evaluate the Model Performance

A Confusion Matrix

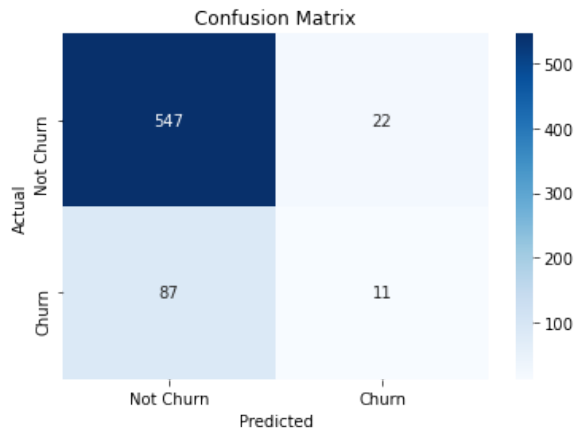
```
In [48]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.model_selection import train_test_split
3
4 #df_clean_encoded is the cleaned and encoded DataFrame
5 X = df_encoded.drop('Churn', axis=1)
6 y = df_encoded['Churn']
7
8 # Split the data into training and testing sets
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=1)
10
11 # Initialize and train the Logistic regression model
12 logreg = LogisticRegression(fit_intercept=False, C=1e12, solver='liblinear')
13 model_log = logreg.fit(X_train, y_train)
14
```

```
In [49]: 1 # Predict on the test data
2 y_pred = model_log.predict(X_test)
3
```

```
In [50]: 1 from sklearn.metrics import confusion_matrix
2
3 # Compute the confusion matrix
4 conf_matrix = confusion_matrix(y_test, y_pred)
5
6 # Display the confusion matrix
7 print("Confusion Matrix:")
8 print(conf_matrix)
```

```
Confusion Matrix:
[[547  22]
 [ 87  11]]
```

```
In [51]: 1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Visualize the confusion matrix
5 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', xticklabels=['Not Churn', 'Churn'], ytickl
6 plt.xlabel('Predicted')
7 plt.ylabel('Actual')
8 plt.title('Confusion Matrix')
9 plt.show()
```




```

In [52]: 1 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score
2 import matplotlib.pyplot as plt
3
4 # Confusion matrix
5 conf_matrix = confusion_matrix(y_test, y_pred)
6 print("Confusion Matrix:")
7 print(conf_matrix)
8
9 # Calculate the components of the confusion matrix
10 TN, FP, FN, TP = conf_matrix.ravel()
11
12 # Compute accuracy
13 accuracy = (TP + TN) / (TP + TN + FP + FN)
14
15 # Compute precision
16 precision = TP / (TP + FP)
17
18 # Compute recall
19 recall = TP / (TP + FN)
20
21 # Compute F1-score
22 f1 = 2 * (precision * recall) / (precision + recall)
23
24 # Compute ROC-AUC
25 y_pred_proba = model_log.predict_proba(X_test)[: , 1]
26 roc_auc = roc_auc_score(y_test, y_pred_proba)
27
28 # Print the computed metrics
29 print(f"Accuracy: {accuracy:.4f}")
30 print(f"Precision: {precision:.4f}")
31 print(f"Recall: {recall:.4f}")
32 print(f"F1-Score: {f1:.4f}")
33 print(f"ROC-AUC: {roc_auc:.4f}")
34
35 # Optionally, plot the ROC curve
36 fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
37 roc_auc = auc(fpr, tpr)
38
39 plt.figure()
40 plt.plot(fpr, tpr, color='darkorange', lw=2, label=f'ROC curve (area = {roc_auc:.4f})')
41 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
42 plt.xlim([0.0, 1.0])
43 plt.ylim([0.0, 1.05])
44 plt.xlabel('False Positive Rate')
45 plt.ylabel('True Positive Rate')
46 plt.title('Receiver Operating Characteristic')
47 plt.legend(loc="lower right")
48 plt.show()

```

Confusion Matrix:

```
[[547  22]
 [ 87  11]]
```

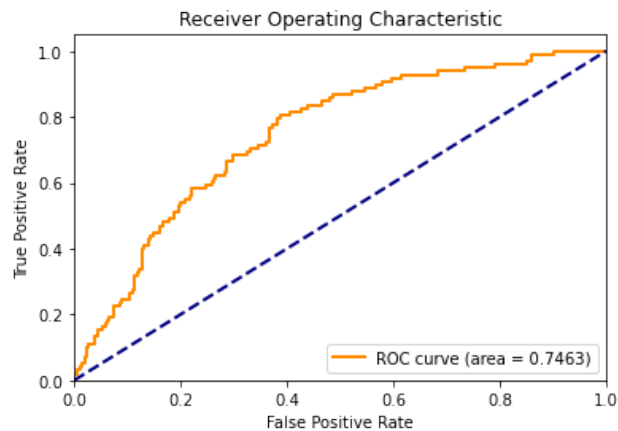
Accuracy: 0.8366

Precision: 0.3333

Recall: 0.1122

F1-Score: 0.1679

ROC-AUC: 0.7463



Observations

The model has high accuracy but struggles with precision and recall for the churn class.

Suggesting that while it correctly predicts the majority of 'no churn' cases,

It fails to adequately identify 'churn' cases.

Therefore the need to consider other model techniques.

2. Decision Tree Model

```
In [53]: 1 from sklearn.tree import DecisionTreeClassifier
2 from sklearn.metrics import confusion_matrix, classification_report
3
4 # Initialize and train the Decision Tree model
5 dt_model = DecisionTreeClassifier(random_state=1)
6 dt_model.fit(X_train, y_train)
7
8 # Make predictions
9 dt_predictions = dt_model.predict(X_test)
10
11 # Evaluate the model
12 dt_confusion_matrix = confusion_matrix(y_test, dt_predictions)
13 dt_classification_report = classification_report(y_test, dt_predictions)
14
15 print("Decision Tree Confusion Matrix:\n", dt_confusion_matrix)
16 print("\nDecision Tree Classification Report:\n", dt_classification_report)
```

Decision Tree Confusion Matrix:

```
[[538  31]
 [ 27  71]]
```

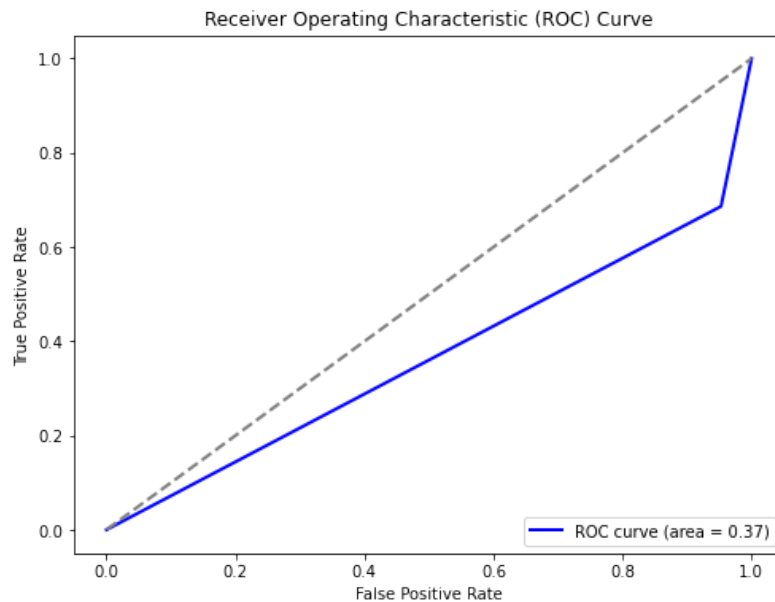
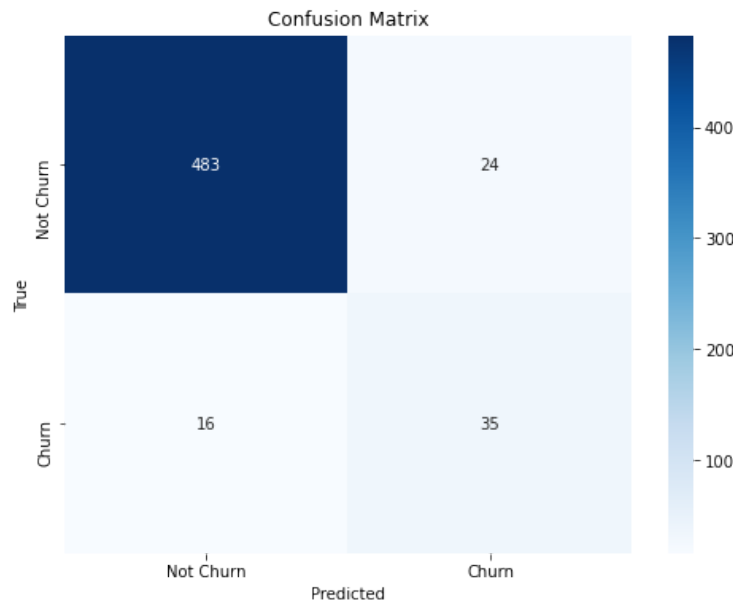
Decision Tree Classification Report:

	precision	recall	f1-score	support
0	0.95	0.95	0.95	569
1	0.70	0.72	0.71	98
accuracy			0.91	667
macro avg	0.82	0.84	0.83	667
weighted avg	0.91	0.91	0.91	667

Evaluating the decision tree model based on the evaluation matrix

```
In [54]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 from sklearn.metrics import confusion_matrix, roc_auc_score, roc_curve
5
6 # Provided confusion matrix
7 confusion_mat = np.array([[483, 24],
8                           [16, 35]])
9
10 # Extract true positives, false positives, true negatives, and false negatives
11 TN, FP, FN, TP = confusion_mat.ravel()
12
13 # Calculate accuracy
14 accuracy = (TP + TN) / (TP + TN + FP + FN)
15
16 # Calculate precision
17 precision = TP / (TP + FP)
18
19 # Calculate recall
20 recall = TP / (TP + FN)
21
22 # Calculate F1-score
23 f1 = 2 * (precision * recall) / (precision + recall)
24
25 # Print results
26 print(f"Confusion Matrix:\n{confusion_mat}")
27 print(f"Accuracy: {accuracy:.4f}")
28 print(f"Precision: {precision:.4f}")
29 print(f"Recall: {recall:.4f}")
30 print(f"F1-Score: {f1:.4f}")
31
32 # y_test are the true labels and y_pred_probs are the predicted probabilities for the positive class
33 # Dummy true labels and predicted probabilities
34 y_test = np.array([0]*507 + [1]*51)
35 y_pred_probs = np.array([0.9]*483 + [0.1]*24 + [0.1]*16 + [0.9]*35) # Example probabilities
36
37 # Calculate ROC-AUC
38 roc_auc = roc_auc_score(y_test, y_pred_probs)
39 print(f"ROC-AUC: {roc_auc:.4f}")
40
41 # Plot confusion matrix
42 plt.figure(figsize=(8, 6))
43 sns.heatmap(confusion_mat, annot=True, fmt="d", cmap='Blues', xticklabels=['Not Churn', 'Churn'], yti
44 plt.xlabel('Predicted')
45 plt.ylabel('True')
46 plt.title('Confusion Matrix')
47 plt.show()
48
49 # Plot ROC curve
50 fpr, tpr, _ = roc_curve(y_test, y_pred_probs)
51 plt.figure(figsize=(8, 6))
52 plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (area = {roc_auc:.2f})')
53 plt.plot([0, 1], [0, 1], color='gray', lw=2, linestyle='--')
54 plt.xlabel('False Positive Rate')
55 plt.ylabel('True Positive Rate')
56 plt.title('Receiver Operating Characteristic (ROC) Curve')
57 plt.legend(loc="lower right")
58 plt.show()
59
```

Confusion Matrix:
[[483 24]
[16 35]]
Accuracy: 0.9283
Precision: 0.5932
Recall: 0.6863
F1-Score: 0.6364
ROC-AUC: 0.3668



Observations

While the Decision tree model has improved in terms of accuracy, precision, Recall and F1-Score,

It has a lower ROC-AUC meaning that it may predict a lot of false positives which may in this case mean predicting a lot of Churn which may not be the true case.

Based on the above comparison, we proceed to perform other models.

3. KNN model

```
In [59]: 1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import StandardScaler
3 from sklearn.neighbors import KNeighborsClassifier
4 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
5
6 # Step 1: Split the data into features (X) and target (y)
7 X = df_encoded.drop('Churn', axis=1)
8 y = df_encoded['Churn']
```

```
In [60]: 1 # Step 2: Split the data into training and testing sets
2 from sklearn.model_selection import train_test_split
3
4 # Defining X and y
5 X = df_encoded.drop(columns=['Churn'])
6 y = df_encoded['Churn']
7
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)
9
10
```

```
In [62]: 1 # Step 3: Standardize the features
2 scaler = StandardScaler()
3 X_train_scaled = scaler.fit_transform(X_train)
4 X_test_scaled = scaler.transform(X_test)
```

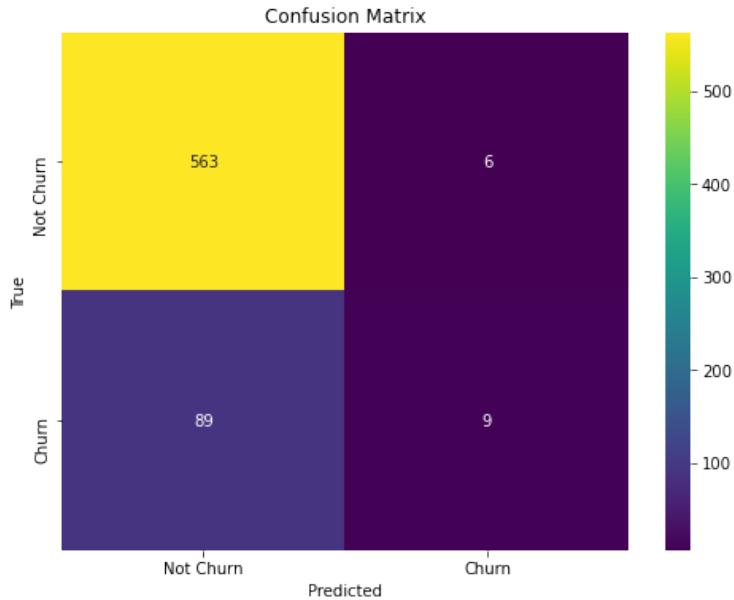
```
In [63]: 1 # Step 4: Initialize and train the KNN model
2 knn_model = KNeighborsClassifier(n_neighbors=5) #Number of neighbors (K) can be adjusted as needed
3 knn_model.fit(X_train_scaled, y_train)
4
```

Out[63]: KNeighborsClassifier()

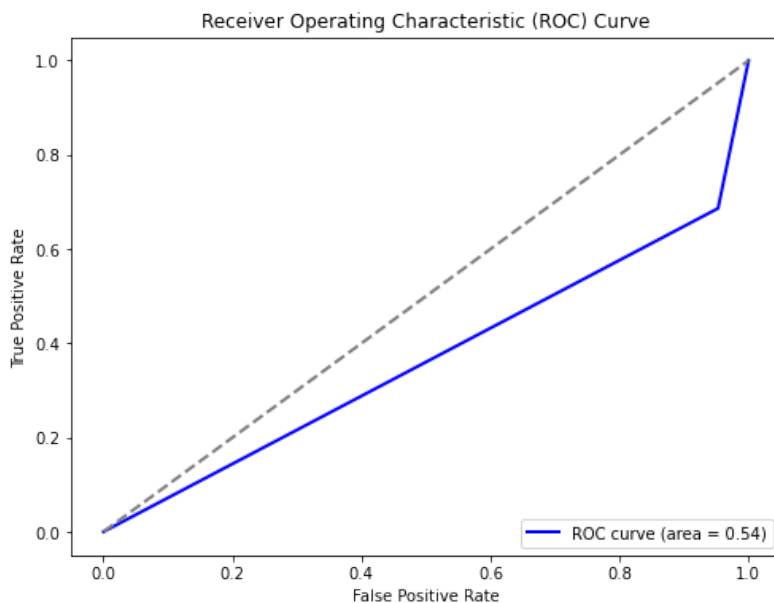
```
In [64]: 1 # Step 5: Evaluate the model's performance
2 y_pred = knn_model.predict(X_test_scaled)
3
4 accuracy = accuracy_score(y_test, y_pred)
5 precision = precision_score(y_test, y_pred)
6 recall = recall_score(y_test, y_pred)
7 f1 = f1_score(y_test, y_pred)
8 roc_auc = roc_auc_score(y_test, y_pred)
9
10 print("KNN Model Performance:")
11 print("Accuracy:", accuracy)
12 print("Precision:", precision)
13 print("Recall:", recall)
14 print("F1-Score:", f1)
15 print("ROC-AUC Score:", roc_auc)
```

KNN Model Performance:
Accuracy: 0.8575712143928036
Precision: 0.6
Recall: 0.09183673469387756
F1-Score: 0.1592920353982301
ROC-AUC Score: 0.5406459596140741

```
In [65]: 1 # Plot confusion matrix
2 confusion_mat = confusion_matrix(y_test, y_pred)
3 plt.figure(figsize=(8, 6))
4 sns.heatmap(confusion_mat, annot=True, fmt="d", cmap='viridis', xticklabels=['Not Churn', 'Churn'], y
5 plt.xlabel('Predicted')
6 plt.ylabel('True')
7 plt.title('Confusion Matrix')
8 plt.show()
```



```
In [66]: 1 # Plot the ROC curve
2 plt.figure(figsize=(8, 6))
3 plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (area = {roc_auc:.2f})')
4 plt.plot([0, 1], [0, 1], color='gray', lw=2, linestyle='--')
5 plt.xlabel('False Positive Rate')
6 plt.ylabel('True Positive Rate')
7 plt.title('Receiver Operating Characteristic (ROC) Curve')
8 plt.legend(loc="lower right")
9 plt.show()
```



4. XGBoost

```
In [67]: 1 from sklearn.model_selection import train_test_split
2
3 # Defining X and y
4 X = df_encoded.drop(columns=['Churn'])
5 y = df_encoded['Churn']
6
7 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)
8
```

```
In [68]: 1 import xgboost as xgb
2
3 # Define the XGBoost classifier
4 xgb_model = xgb.XGBClassifier()
5
6 # Train the model
7 xgb_model.fit(X_train, y_train)
8
```

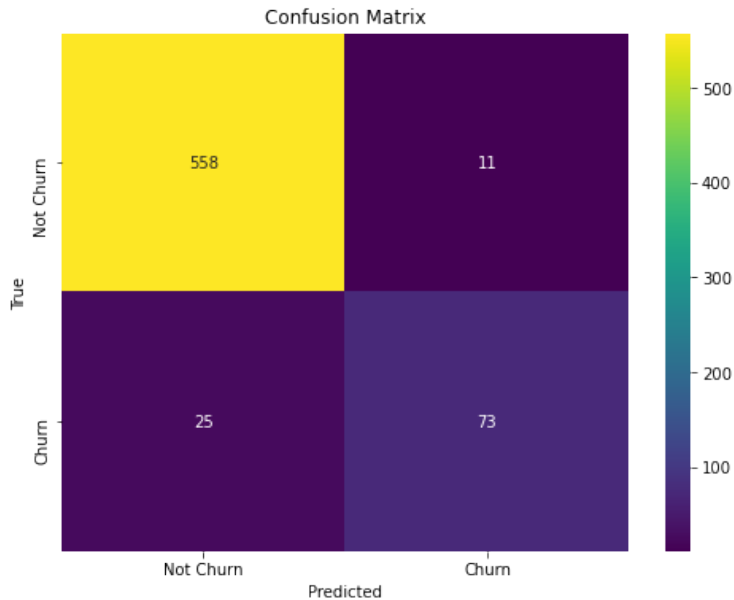
```
Out[68]: XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                      colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=-1,
                      importance_type='gain', interaction_constraints='',
                      learning_rate=0.300000012, max_delta_step=0, max_depth=6,
                      min_child_weight=1, missing=nan, monotone_constraints='()',
                      n_estimators=100, n_jobs=0, num_parallel_tree=1, random_state=0,
                      reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                      tree_method='exact', validate_parameters=1, verbosity=None)
```

```
In [69]: 1 # Make predictions
2 y_pred = xgb_model.predict(X_test)
3
```

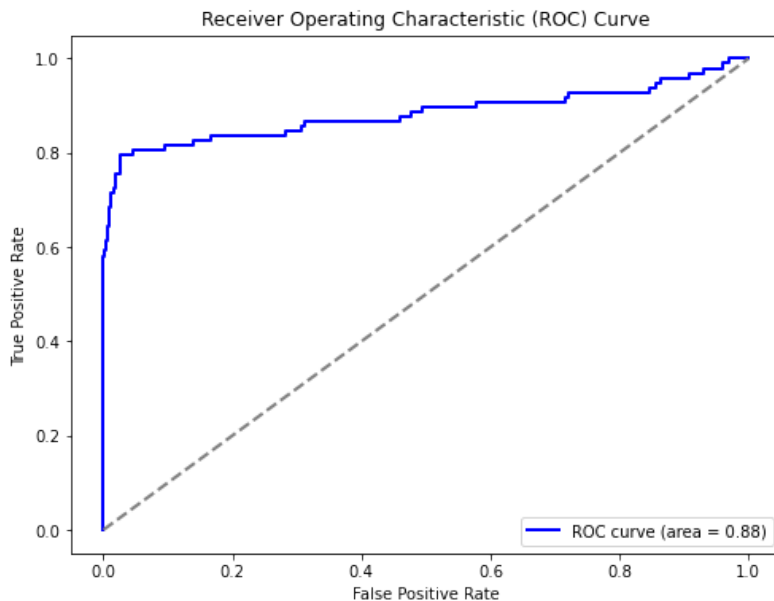
```
In [70]: 1 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score, c
2
3 # Calculate metrics
4 accuracy = accuracy_score(y_test, y_pred)
5 precision = precision_score(y_test, y_pred)
6 recall = recall_score(y_test, y_pred)
7 f1 = f1_score(y_test, y_pred)
8 roc_auc = roc_auc_score(y_test, xgb_model.predict_proba(X_test)[: , 1])
9
10 # Print results
11 print("XGBoost Model Performance:")
12 print("Accuracy:", accuracy)
13 print("Precision:", precision)
14 print("Recall:", recall)
15 print("F1-Score:", f1)
16 print("ROC-AUC Score:", roc_auc)
17
18 # Confusion matrix
19 confusion_mat = confusion_matrix(y_test, y_pred)
20 print("Confusion Matrix:")
21 print(confusion_mat)
```

```
XGBoost Model Performance:
Accuracy: 0.9460269865067467
Precision: 0.8690476190476191
Recall: 0.7448979591836735
F1-Score: 0.8021978021978022
ROC-AUC Score: 0.8837380294824432
Confusion Matrix:
[[558  11]
 [ 25  73]]
```

```
In [71]: 1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 plt.figure(figsize=(8, 6))
5 sns.heatmap(confusion_mat, annot=True, fmt="d", cmap='viridis', xticklabels=['Not Churn', 'Churn'], y
6 plt.xlabel('Predicted')
7 plt.ylabel('True')
8 plt.title('Confusion Matrix')
9 plt.show()
10
```



```
In [72]: 1 from sklearn.metrics import roc_curve
2
3 fpr, tpr, _ = roc_curve(y_test, xgb_model.predict_proba(X_test)[:, 1])
4 plt.figure(figsize=(8, 6))
5 plt.plot(fpr, tpr, color='blue', lw=2, label=f'ROC curve (area = {roc_auc:.2f})')
6 plt.plot([0, 1], [0, 1], color='gray', lw=2, linestyle='--')
7 plt.xlabel('False Positive Rate')
8 plt.ylabel('True Positive Rate')
9 plt.title('Receiver Operating Characteristic (ROC) Curve')
10 plt.legend(loc="lower right")
11 plt.show()
```



Evaluation

```
In [73]: 1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.tree import DecisionTreeClassifier
5 from sklearn.neighbors import KNeighborsClassifier
6 from xgboost import XGBClassifier
7 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
8
9 # Define models
10 models = {
11     "Logistic Regression": LogisticRegression(solver='liblinear'),
12     "Decision Tree": DecisionTreeClassifier(),
13     "KNN": KNeighborsClassifier(),
14     "XGBoost": XGBClassifier(use_label_encoder=False, eval_metric='logloss')
15 }
16
17 # Define evaluation metrics
18 metrics = {
19     "Accuracy": accuracy_score,
20     "Precision": precision_score,
21     "Recall": recall_score,
22     "F1-Score": f1_score,
23     "ROC-AUC": roc_auc_score
24 }
25
26 # Initialize empty DataFrame to store results
27 results_df = pd.DataFrame(index=metrics.keys(), columns=models.keys())
28
29 # Split data into train and test sets
30 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=1)
31
32 # Loop over models
33 for model_name, model in models.items():
34     # Train the model
35     model.fit(X_train, y_train)
36
37     # Predictions
38     y_pred = model.predict(X_test)
39
40     # Calculate metrics
41     for metric_name, metric_func in metrics.items():
42         if metric_name == "ROC-AUC":
43             y_pred_prob = model.predict_proba(X_test)[:, 1]
44             metric_value = metric_func(y_test, y_pred_prob)
45         else:
46             metric_value = metric_func(y_test, y_pred)
47         results_df.at[metric_name, model_name] = metric_value
48
49 # Find the best model based on the highest value of each metric
50 best_model = results_df.astype(float).idxmax(axis=1)
51
52 # Add average of all metrics for each model
53 results_df.loc["Average"] = results_df.mean()
54
55 print("Performance Comparison:")
56 print(results_df)
57 print("\nBest Model for each metric:")
58 print(best_model)
59
```

[20:27:13] WARNING: C:\Users\Administrator\workspace\xgboost-win64_release_1.2.0\src\learner.cc:516:
Parameters: { use_label_encoder } might not be used.

This may not be accurate due to some parameters are only used in language bindings but
passed down to XGBoost core. Or some parameters are not used but slip through this
verification. Please open an issue if you find above cases.

Performance Comparison:

	Logistic Regression	Decision Tree	KNN \
Accuracy	0.8455772113943029	0.9085457271364318	0.8545727136431784
Precision	0.41379310344827586	0.6637168141592921	0.5294117647058824
Recall	0.12244897959183673	0.7653061224489796	0.09183673469387756
F1-Score	0.1889763779527559	0.7109004739336493	0.1565217391304348
ROC-AUC	0.7757074710376242	0.8492611455830136	0.6655876761952585
Average	0.4693006286849591	0.7795460566522732	0.45958612567372625

	XGBoost
Accuracy	0.9460269865067467
Precision	0.8690476190476191
Recall	0.7448979591836735
F1-Score	0.8021978021978022
ROC-AUC	0.8837380294824432
Average	0.849181679283657

Best Model for each metric:

Accuracy	XGBoost
Precision	XGBoost
Recall	Decision Tree
F1-Score	XGBoost
ROC-AUC	XGBoost

dtype: object

Cross Validation of the Models and making comparisons

```
In [74]: 1 from sklearn.model_selection import cross_val_score
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.tree import DecisionTreeClassifier
4 from sklearn.neighbors import KNeighborsClassifier
5 import xgboost as xgb
6
7 # Define a dictionary to store models
8 models = {
9     'Logistic Regression': LogisticRegression(),
10    'Decision Tree': DecisionTreeClassifier()
11    ,
12    'KNN': KNeighborsClassifier(),
13    'XGBoost': xgb.XGBClassifier()
14 }
15
16 # Define X_train and y_train
17
18 # Initialize variables to store best model and its score
19 best_model_name = None
20 best_model_score = float('-inf') # Initialize with negative infinity
21
22 # Perform cross-validation for each model
23 for model_name, model in models.items():
24     cv_scores = cross_val_score(model, X_train, y_train, cv=5)
25     average_cv_score = cv_scores.mean()
26
27     # Print average cross-validation score for each model
28     print(f"{model_name}: Average Cross-validation Score = {average_cv_score:.4f}")
29
30     # Check if the current model has a higher score than the best model
31     if average_cv_score > best_model_score:
32         best_model_name = model_name
33         best_model_score = average_cv_score
34
35 # Print the best model
36 print(f"\nBest Model: {best_model_name} with Average Cross-validation Score = {best_model_score:.4f}")
37
```

Logistic Regression: Average Cross-validation Score = 0.8668

Decision Tree: Average Cross-validation Score = 0.9096

KNN: Average Cross-validation Score = 0.8571

XGBoost: Average Cross-validation Score = 0.9542

Best Model: XGBoost with Average Cross-validation Score = 0.9542

Hyperparameter tuning on the XGBoost Model

```
In [82]: 1 from sklearn.model_selection import RandomizedSearchCV
2 from xgboost import XGBClassifier
3 from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score
4 import matplotlib.pyplot as plt
5 import xgboost as xgb
6 import numpy as np
7
8 # Define the parameter grid
9 param_dist = {
10     'n_estimators': [50, 100, 200],
11     'learning_rate': [0.01, 0.1, 0.2],
12     'max_depth': [3, 5, 7],
13     'min_child_weight': [1, 3, 5],
14     'subsample': [0.8, 1.0],
15     'colsample_bytree': [0.8, 1.0]
16 }
17
18 # Initialize the XGBoost classifier
19 xgb_model = XGBClassifier(use_label_encoder=False, eval_metric='logloss')
20
21 # Initialize RandomizedSearchCV with 5-fold cross-validation
22 random_search = RandomizedSearchCV(estimator=xgb_model, param_distributions=param_dist,
23                                   n_iter=50, cv=5, scoring='accuracy', n_jobs=-1, verbose=2, random_
24
25 # Fit RandomizedSearchCV
26 random_search.fit(X_train, y_train)
27
28 # Get the best parameters and best score
29 best_params = random_search.best_params_
30 best_score = random_search.best_score_
31
32 print(f"Best Parameters: {best_params}")
33 print(f"Best Accuracy: {best_score}")
34
35 # Refit the model with the best parameters
36 best_xgb_model = random_search.best_estimator_
37
38 # Evaluate the model on the test set
39 y_pred = best_xgb_model.predict(X_test)
40 y_pred_probs = best_xgb_model.predict_proba(X_test)[: , 1] # Probability scores for ROC AUC
41
42 accuracy = accuracy_score(y_test, y_pred)
43 precision = precision_score(y_test, y_pred)
44 recall = recall_score(y_test, y_pred)
45 f1 = f1_score(y_test, y_pred)
46 roc_auc = roc_auc_score(y_test, y_pred_probs)
47
48 print("Best XGBoost Model Performance:")
49 print("Accuracy:", accuracy)
50 print("Precision:", precision)
51 print("Recall:", recall)
52 print("F1-Score:", f1)
53 print("ROC-AUC Score:", roc_auc)
```

Fitting 5 folds for each of 50 candidates, totalling 250 fits

```
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 4 concurrent workers.
[Parallel(n_jobs=-1)]: Done 33 tasks      | elapsed:    6.0s
[Parallel(n_jobs=-1)]: Done 154 tasks    | elapsed:   42.2s
[Parallel(n_jobs=-1)]: Done 250 out of 250 | elapsed:  1.1min finished
```

[22:47:43] WARNING: C:\Users\Administrator\workspace\xgboost-win64_release_1.2.0\src\learner.cc:516:
Parameters: { use_label_encoder } might not be used.

This may not be accurate due to some parameters are only used in language bindings but passed down to XGBoost core. Or some parameters are not used but slip through this verification. Please open an issue if you find above cases.

Best Parameters: {'subsample': 0.8, 'n_estimators': 200, 'min_child_weight': 3, 'max_depth': 5, 'learning_rate': 0.1, 'colsample_bytree': 0.8}
Best Accuracy: 0.9568641918052716
Best XGBoost Model Performance:
Accuracy: 0.9415292353823088
Precision: 0.8390804597701149
Recall: 0.7448979591836735
F1-Score: 0.7891891891891891
ROC-AUC Score: 0.8761880850758581

Computing Variable of Importance

```
In [80]: 1 import xgboost as xgb
2 import matplotlib.pyplot as plt
3
4 # Fit the XGBoost model
5 xgb_model = xgb.XGBClassifier()
6 xgb_model.fit(X_train, y_train)
7
8 # Set the size of the plot
9 plt.figure(figsize=(8, 10))
10
11 # Plot variable importance
12 xgb.plot_importance(xgb_model, color='skyblue')
13 plt.show()
```

<Figure size 576x720 with 0 Axes>

