# SyriaTel Customer Churn Prediction Project

## Business Understanding

**Business Problem**: Reducing the rate of customer churn for SyriaTel Telecom
**Background**: SyriaTel, a major player in the telecommunications industry, is dealing with a serious problem of customer churn, which occurs when customers cancel their subscriptions, leading to a loss of revenue. To tackle this challenge, SyriaTel plans to develop a predictive model to identify customers who are at risk of churning. By reaching out to these customers with targeted retention strategies, SyriaTel aims to decrease churn rates and keep valuable customers.
**Problem statement**: Can we predict customer churn for SyriaTel and determine the main factors contributing to it, allowing the company to develop effective retention strategies.

## Objectives:

- Develop a Churn Prediction model capable of accurately identifying customers likely to churn using the available dataset.
- Discover the key patterns and features that contribute to predicting customer churn.

This initiative supports SyriaTel's objective of retaining customers and minimizing revenue loss from churn, highlighting the value of data-driven decision-making in the telecommunications sector.

## Expected Outcome

The main metric for evaluating the classification model's performance is 'recall,' which measures how well the model identifies customers likely to churn. The goal is to minimize false negatives since missing a potential churner is more costly for the business than incorrectly predicting a non-churner. The target is to achieve at least 80% recall.

However, a balance is necessary. Predicting that all customers will churn would result in perfect recall but offer no real business value, as not all customers are at risk. Therefore, in addition to recall, 'precision' and 'accuracy' will be used as secondary metrics to ensure a comprehensive assessment of the model's performance.

## Import the libraries

Here, we import the various necessary libraries that we will be using in this project. They include Pandas, Numpy, visualization libraries like matplotlib and seaborn, various scikit-learn modules for machine learning and metrics we will use for evaluating our models.

```python
In [35]:  # Importing the necessary libraries
          import pandas as pd
          import numpy as np
          import seaborn as sns
          import matplotlib.pyplot as plt
          import matplotlib.cm as cm
          import math
          %matplotlib inline

          import scipy.stats as stats

          from sklearn.model_selection import train_test_split
          from sklearn.preprocessing import StandardScaler, OneHotEncoder
          from sklearn.preprocessing import LabelEncoder
          from sklearn.linear_model import LogisticRegression
          from sklearn import tree
          from sklearn.tree import DecisionTreeClassifier,plot_tree
          from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier
          from sklearn.compose import ColumnTransformer
          from sklearn.pipeline import Pipeline
          from imblearn.over_sampling import SMOTE
          from sklearn.svm import SVC
          from sklearn.dummy import DummyClassifier
          from sklearn.model_selection import GridSearchCV
          from sklearn.model_selection import cross_val_score
          from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
          from sklearn.metrics import roc_curve, auc
          from sklearn.metrics import classification_report, confusion_matrix, ConfusionMatrixD
          from sklearn.feature_selection import RFECV

          import warnings
          warnings.filterwarnings('ignore')
```

## Data Understanding

In the data folder is a "telcom_data.csv" file, a dataset from SyriaTel Telcom, that contains data about various customers, and whether they churned or not.

```python
In [36]:  pd.set_option('display.max_columns', None)
```

```python
In [37]:  # Loading the dataset
          df = pd.read_csv('C:/Users/HP/AppData/Local/Temp/81e1e7e9-e08e-41c6-a8e3-d22ea09a02d9
          df.head()
```

Out[37]:

| | 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | KS | 128 | 415 | 382-4657 | no | yes | 25 | 265.1 | 110 | 45.07 | 197.4 |
| **1** | OH | 107 | 415 | 371-7191 | no | yes | 26 | 161.6 | 123 | 27.47 | 195.5 |
| **2** | NJ | 137 | 415 | 358-1921 | no | no | 0 | 243.4 | 114 | 41.38 | 121.2 |
| **3** | OH | 84 | 408 | 375-9999 | yes | no | 0 | 299.4 | 71 | 50.90 | 61.9 |

| | 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **4** | OK | 75 | 415 | 330-6626 | yes | no | 0 | 166.7 | 113 | 28.34 | 148.3 |

In [38]:
```python
# Checking last 4 records
df.tail()
```

Out[38]:

| | state | account length | area code | phone number | international plan | voice mail plan | number vmail messages | total day minutes | total day calls | total day charge | tot eve minut |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **3328** | AZ | 192 | 415 | 414-4276 | no | yes | 36 | 156.2 | 77 | 26.55 | 215 |
| **3329** | WV | 68 | 415 | 370-3271 | no | no | 0 | 231.1 | 57 | 39.29 | 153 |
| **3330** | RI | 28 | 510 | 328-8230 | no | no | 0 | 180.8 | 109 | 30.74 | 288 |
| **3331** | CT | 184 | 510 | 364-6381 | yes | no | 0 | 213.8 | 105 | 36.35 | 159 |
| **3332** | TN | 74 | 415 | 400-4344 | no | yes | 25 | 234.4 | 113 | 39.85 | 265 |

From the above preview of the dataset, we can see that the dataset is uniform.

In [39]:
```python
# Checking the shape of the dataset
df.shape
```

Out[39]: (3333, 21)

In [40]:
```python
# Overview of the columns
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
```

From the above, we can see that the columns are of different types: float, int, object and bool. However, most of the columns have some whitespace in between their names, and it would be better to rename them in such a way that they have an underscore instead of the whitespace.

Renaming the by replacing the whitespace with an underscore and displaying unique values

In [41]:
```python
df.columns = df.columns.str.replace(' ', '_')

# Display unique values for each column
for column in df.columns:
    print(f"There are {df[column].nunique()} Unique values in '{column}':")
    print(df[column].unique())
    print('-' * 50)
```

```
There are 51 Unique values in 'state':
['KS' 'OH' 'NJ' 'OK' 'AL' 'MA' 'MO' 'LA' 'WV' 'IN' 'RI' 'IA' 'MT' 'NY'
 'ID' 'VT' 'VA' 'TX' 'FL' 'CO' 'AZ' 'SC' 'NE' 'WY' 'HI' 'IL' 'NH' 'GA'
 'AK' 'MD' 'AR' 'WI' 'OR' 'MI' 'DE' 'UT' 'CA' 'MN' 'SD' 'NC' 'WA' 'NM'
 'NV' 'DC' 'KY' 'ME' 'MS' 'TN' 'PA' 'CT' 'ND']
--------------------------------------------------
There are 212 Unique values in 'account_length':
[128 107 137  84  75 118 121 147 117 141  65  74 168  95  62 161  85  93
  76  73  77 130 111 132 174  57  54  20  49 142 172  12  72  36  78 136
 149  98 135  34 160  64  59 119  97  52  60  10  96  87  81  68 125 116
  38  40  43 113 126 150 138 162  90  50  82 144  46  70  55 106  94 155
  80 104  99 120 108 122 157 103  63 112  41 193  61  92 131 163  91 127
 110 140  83 145  56 151 139   6 115 146 185 148  32  25 179  67  19 170
 164  51 208  53 105  66  86  35  88 123  45 100 215  22  33 114  24 101
 143  48  71 167  89 199 166 158 196 209  16  39 173 129  44  79  31 124
  37 159 194 154  21 133 224  58  11 109 102 165  18  30 176  47 190 152
  26  69 186 171  28 153 169  13  27   3  42 189 156 134 243  23   1 205
 200   5   9 178 181 182 217 177 210  29 180   2  17   7 212 232 192 195
 197 225 184 191 201  15 183 202   8 175   4 188 204 221]
--------------------------------------------------
There are 3 Unique values in 'area_code':
[415 408 510]
--------------------------------------------------
There are 3333 Unique values in 'phone_number':
['382-4657' '371-7191' '358-1921' ... '328-8230' '364-6381' '400-4344']
--------------------------------------------------
There are 2 Unique values in 'international_plan':
['no' 'yes']
--------------------------------------------------
There are 2 Unique values in 'voice_mail_plan':
['yes' 'no']
--------------------------------------------------
There are 46 Unique values in 'number_vmail_messages':
[25 26  0 24 37 27 33 39 30 41 28 34 46 29 35 21 32 42 36 22 23 43 31 38
 40 48 18 17 45 16 20 14 19 51 15 11 12 47  8 44 49  4 10 13 50  9]
--------------------------------------------------
There are 1667 Unique values in 'total_day_minutes':
[265.1 161.6 243.4 ... 321.1 231.1 180.8]
--------------------------------------------------
There are 119 Unique values in 'total_day_calls':
[110 123 114  71 113  98  88  79  97  84 137 127  96  70  67 139  66  90
 117  89 112 103  86  76 115  73 109  95 105 121 118  94  80 128  64 106
 102  85  82  77 120 133 135 108  57  83 129  91  92  74  93 101 146  72
  99 104 125  61 100  87 131  65 124 119  52  68 107  47 116 151 126 122
```

```
 111 145   78 136 140 148   81   55   69 158 134 130   63   53   75 141 163   59
 132 138   54   58   62 144 143 147   36   40 150   56   51 165   30   48   60   42
   0   45 160 149 152 142 156   35   49 157   44]
----------------------------------------------------
There are 1667 Unique values in 'total_day_charge':
[45.07 27.47 41.38 ... 54.59 39.29 30.74]
----------------------------------------------------
There are 1611 Unique values in 'total_eve_minutes':
[197.4 195.5 121.2 ... 153.4 288.8 265.9]
----------------------------------------------------
There are 123 Unique values in 'total_eve_calls':
[ 99 103 110   88 122 101 108   94   80 111   83 148   71   75   76   97   90   65
  93 121 102   72 112 100   84 109   63 107 115 119 116   92   85   98 118   74
 117   58   96   66   67   62   77 164 126 142   64 104   79   95   86 105   81 113
 106   59   48   82   87 123 114 140 128   60   78 125   91   46 138 129   89 133
 136   57 135 139   51   70 151 137 134   73 152 168   68 120   69 127 132 143
  61 124   42   54 131   52 149   56   37 130   49 146 147   55   12   50 157 155
  45 144   36 156   53 141   44 153 154 150   43    0 145 159 170]
----------------------------------------------------
There are 1440 Unique values in 'total_eve_charge':
[16.78 16.62 10.3  ... 13.04 24.55 22.6 ]
----------------------------------------------------
There are 1591 Unique values in 'total_night_minutes':
[244.7 254.4 162.6 ... 280.9 120.1 279.1]
----------------------------------------------------
There are 120 Unique values in 'total_night_calls':
[ 91 103 104   89 121 118   96   90   97 111   94 128 115   99   75 108   74 133
  64   78 105   68 102 148   98 116   71 109 107 135   92   86 127   79   87 129
  57   77   95   54 106   53   67 139   60 100   61   73 113   76 119   88   84   62
 137   72 142 114 126 122   81 123 117   82   80 120 130 134   59 112 132 110
 101 150   69 131   83   93 124 136 125   66 143   58   55   85   56   70   46   42
 152   44 145   50 153   49 175   63 138 154 140 141 146   65   51 151 158 155
 157 147 144 149 166   52   33 156   38   36   48 164]
----------------------------------------------------
There are 933 Unique values in 'total_night_charge':
[11.01 11.45  7.32  8.86  8.41  9.18  9.57  9.53  9.71 14.69  9.4   8.82
  6.35  8.65  9.14  7.23  4.02  5.83  7.46  8.68  9.43  8.18  8.53 10.67
 11.28  8.22  4.59  8.17  8.04 11.27 11.08 13.2  12.61  9.61  6.88  5.82
 10.25  4.58  8.47  8.45  5.5  14.02  8.03 11.94  7.34  6.06 10.9   6.44
  3.18 10.66 11.21 12.73 10.28 12.16  6.34  8.15  5.84  8.52  7.5   7.48
  6.21 11.95  7.15  9.63  7.1   6.91  6.69 13.29 11.46  7.76  6.86  8.16
 12.15  7.79  7.99 10.29 10.08 12.53  7.91 10.02  8.61 14.54  8.21  9.09
  4.93 11.39 11.88  5.75  7.83  8.59  7.52 12.38  7.21  5.81  8.1  11.04
 11.19  8.55  8.42  9.76  9.87 10.86  5.36 10.03 11.15  9.51  6.22  2.59
  7.65  6.45  9.    6.4   9.94  5.08 10.23 11.36  6.97 10.16  7.88 11.91
  6.61 11.55 11.76  9.27  9.29 11.12 10.69  8.8  11.85  7.14  8.71 11.42
  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
  10.84 10.22  6.52  5.55  7.63  5.11  5.89 10.78  3.05 11.89  8.97 10.44
  10.5   9.35  5.66 11.09  9.83  5.44 10.11  6.39 11.93  8.62 12.06  6.02
   8.85  5.25  8.66  6.73 10.21 11.59 13.87  7.77 10.39  5.54  6.62 13.33
   6.24 12.59  6.3   6.79  8.28  9.03  8.07  5.52 12.14 10.59  7.54  7.67
   5.47  8.81  8.51 13.45  8.77  6.43 12.01 12.08  7.07  6.51  6.84  9.48
  13.78 11.54 11.67  8.13 10.79  7.13  4.72  4.64  8.96 13.03  6.07  3.51
   6.83  6.12  9.31  9.58  4.68  5.32  9.26 11.52  9.11 10.55 11.47  9.3
  13.82  8.44  5.77 10.96 11.74  8.9  10.47  7.85 10.92  4.74  9.74 10.43
   9.96 10.18  9.54  7.89 12.36  8.54 10.07  9.46  7.3  11.16  9.16 10.19
   5.99 10.88  5.8   7.19  4.55  8.31  8.01 14.43  8.3  14.3   6.53  8.2
  11.31 13.    6.42  4.24  7.44  7.51 13.1   9.49  6.14  8.76  6.65 10.56
   6.72  8.29 12.09  5.39  2.96  7.59  7.24  4.28  9.7   8.83 13.3  11.37
   9.33  5.01  3.26 11.71  8.43  9.68 15.56  9.8   3.61  6.96 11.61 12.81
  10.87 13.84  5.03  5.17  2.03 10.34  9.34  7.95 10.09  9.95  7.11  9.22
   6.13 11.05  9.89  9.39 14.06 10.26 13.31 15.43 16.39  6.27 10.64 11.5
  12.48  8.27 13.53 10.36 12.24  8.69 10.52  9.07 11.51  9.25  8.72  6.78
   8.6  11.84  5.78  5.85 12.3   5.76 12.07  9.6   8.84 12.39 10.1   9.73
   2.85  6.66  2.45  5.28 11.73 10.75  7.74  6.76  6.    7.58 13.69  7.93
   7.68  9.75  4.96  5.49 11.83  7.18  9.19  7.7   7.25 10.74  4.27 13.8
   9.12  4.75  7.78 11.63  7.55  2.25  9.45  9.86  7.71  4.95  7.4  11.17
  11.33  6.82 13.7   1.97 10.89 12.77 10.31  5.23  5.27  9.41  6.09 10.61
   7.29  4.23  7.57  3.67 12.69 14.5   5.95  7.87  5.96  5.94 12.23  4.9
  12.33  6.89  9.67 12.68 12.87  3.7   6.04 13.13 15.74 11.87  4.7   4.67
   7.05  5.42  4.09  5.73  9.47  8.05  6.87  3.71 15.86  7.49 11.69  6.46
  10.45 12.9   5.41 11.26  1.04  6.49  6.37 12.21  6.77 12.65  7.86  9.44
   4.3   7.38  5.02 10.63  2.86 17.19  8.67  8.37  6.9  10.93 10.38  7.36
  10.27 10.95  6.11  4.45 11.9  15.01 12.84  7.45  6.98 11.72  7.56 11.38
  10.    4.42  9.81  5.56  6.01 10.12 12.4  16.99  5.68 11.64  3.78  7.82
   9.85 13.74 12.71 10.98 10.01  9.52  7.31  8.35 11.35  9.5  14.03  3.2
   7.72 13.22 10.7   8.99 10.6  13.02  9.77 12.58 12.35 12.2  11.4  13.91
   3.57 14.65 12.28  5.13 10.72 12.86 14.    7.12 12.17  4.71  6.28  8.
   7.01  5.91  5.2  12.   12.02 12.88  7.28  5.4  12.04  5.24 10.3  10.41
  13.41 12.72  9.08  7.08 13.5   5.35 12.45  5.3  10.32  5.15 12.67  5.22
   5.57  3.94  4.41 13.27 10.24  4.25 12.89  5.72 12.5  11.29  3.25 11.53
   9.82  7.26  4.1  10.37  4.98  6.74 12.52 14.56  8.34  3.82  3.86 13.97
  11.57  6.5  13.58 14.32 13.75 11.14 14.18  9.13  4.46  4.83  9.69 14.13
   7.16  7.98 13.66 14.78 11.2   9.93 11.    5.29  9.92  4.29 11.1  10.51
  12.49  4.04 12.94  7.09  6.71  7.94  5.31  5.98  7.2  14.82 13.21 12.32
  10.58  4.92  6.2   4.47 11.98  6.18  7.81  4.54  5.37  7.17  5.33 14.1
   5.7  12.18  8.98  5.1  14.67 13.95 16.55 11.18  4.44  4.73  2.55  6.31
   2.43  9.24  7.37 13.42 12.42 11.8  14.45  2.89 13.23 12.6  13.18 12.19
  14.81  6.55 11.3  12.27 13.98  8.23 15.49  6.47 13.48 13.59 13.25 17.77
  13.9   3.97 11.56 14.08 13.6   6.26  4.61 12.76 15.76  6.38  3.6  12.8
   5.9   7.97  5.   10.97  5.88 12.34 12.03 14.97 15.06 12.85  6.54 11.24
  12.64  7.06  5.38 13.14  3.99  3.32  4.51  4.12  3.93  2.4  11.75  4.03
  15.85  6.81 14.25 14.09 16.42  6.7  12.74  2.76 12.12  6.99  6.68 11.81
   7.96  5.06 13.16  2.13 13.17  5.12  5.65 12.37 10.53]
------------------------------------------------------
There are 162 Unique values in 'total_intl_minutes':
[10.  13.7 12.2  6.6 10.1  6.3  7.5  7.1  8.7 11.2 12.7  9.1 12.3 13.1
  5.4 13.8  8.1 13.  10.6  5.7  9.5  7.7 10.3 15.5 14.7 11.1 14.2 12.6
 11.8  8.3 14.5 10.5  9.4 14.6  9.2  3.5  8.5 13.2  7.4  8.8 11.   7.8
  6.8 11.4  9.3  9.7 10.2  8.   5.8 12.1 12.  11.6  8.2  6.2  7.3  6.1
 11.7 15.   9.8 12.4  8.6 10.9 13.9  8.9  7.9  5.3  4.4 12.5 11.3  9.
  9.6 13.3 20.   7.2  6.4 14.1 14.3  6.9 11.5 15.8 12.8 16.2  0.  11.9
  9.9  8.4 10.8 13.4 10.7 17.6  4.7  2.7 13.5 12.9 14.4 10.4  6.7 15.4
  4.5  6.5 15.6  5.9 18.9  7.6  5.   7.  14.  18.  16.  14.8  3.7  2.
  4.8 15.3  6.  13.6 17.2 17.5  5.6 18.2  3.6 16.5  4.6  5.1  4.1 16.3
```

```
  14.9 16.4 16.7  1.3 15.2 15.1 15.9  5.5 16.1  4.  16.9  5.2  4.2 15.7
  17.   3.9  3.8  2.2 17.1  4.9 17.9 17.3 18.4 17.8  4.3  2.9  3.1  3.3
   2.6  3.4  1.1 18.3 16.6  2.1  2.4  2.5]
--------------------------------------------------
There are 21 Unique values in 'total_intl_calls':
[ 3  5  7  6  4  2  9 19  1 10 15  8 11  0 12 13 18 14 16 20 17]
--------------------------------------------------
There are 162 Unique values in 'total_intl_charge':
[2.7  3.7  3.29 1.78 2.73 1.7  2.03 1.92 2.35 3.02 3.43 2.46 3.32 3.54
 1.46 3.73 2.19 3.51 2.86 1.54 2.57 2.08 2.78 4.19 3.97 3.   3.83 3.4
 3.19 2.24 3.92 2.84 2.54 3.94 2.48 0.95 2.3  3.56 2.   2.38 2.97 2.11
 1.84 3.08 2.51 2.62 2.75 2.16 1.57 3.27 3.24 3.13 2.21 1.67 1.97 1.65
 3.16 4.05 2.65 3.35 2.32 2.94 3.75 2.4  2.13 1.43 1.19 3.38 3.05 2.43
 2.59 3.59 5.4  1.94 1.73 3.81 3.86 1.86 3.11 4.27 3.46 4.37 0.   3.21
 2.67 2.27 2.92 3.62 2.89 4.75 1.27 0.73 3.65 3.48 3.89 2.81 1.81 4.16
 1.22 1.76 4.21 1.59 5.1  2.05 1.35 1.89 3.78 4.86 4.32 4.   1.   0.54
 1.3  4.13 1.62 3.67 4.64 4.73 1.51 4.91 0.97 4.46 1.24 1.38 1.11 4.4
 4.02 4.43 4.51 0.35 4.1  4.08 4.29 1.49 4.35 1.08 4.56 1.4  1.13 4.24
 4.59 1.05 1.03 0.59 4.62 1.32 4.83 4.67 4.97 4.81 1.16 0.78 0.84 0.89
 0.7  0.92 0.3  4.94 4.48 0.57 0.65 0.68]
--------------------------------------------------
There are 10 Unique values in 'customer_service_calls':
[1 0 2 3 4 5 7 9 6 8]
--------------------------------------------------
There are 2 Unique values in 'churn':
[False  True]
--------------------------------------------------
```

## Columns description

- **state** : Different states of the customers

- **account_length**: It denotes the number of days or duration for which the customer has held their SyriaTel account.

- **area_state**: This column represents the state or location of the customer within the service area of SyriaTel.

- **area_code**: This column typically specifies the area code associated with the customer's phone number.

- **phone_number**: It contains the unique phone number of each customer, serving as an identifier.

- **voice_mail_plan**: Similar to the international plan, this binary column denotes whether the customer has subscribed to a voicemail plan.

- **number_vmail_messages**: If a voice mail plan is active, this column may represent the number of voicemail messages received by the customer.

- **total_day_minutes**: This column records the total number of minutes the customer used during the daytime.

- **total_day_calls**: It indicates the total number of calls made by the customer during the daytime.

- **total_day_charge**: This is the total charge incurred for daytime usage.

- **total_evening_minutes**: Similar to the daytime minutes, this column represents the total number of minutes used in the evening.

- **total_evening_calls**: It denotes the total number of calls made in the evening.

- **total_evening_charge**: This is the total charge for evening usage.

- **total_night_minutes**: Represents the total number of minutes used during the nighttime.

- **total_night_calls**: Denotes the total number of calls made during the nighttime.

- **total_night_charge**: This column reflects the total charge for nighttime usage.

- **total_intl_minutes**: It records the total number of international minutes used by the customer.

- **total_intl_calls**: Indicates the total number of international calls made.

- **total_intl_charge**: Represents the total charge incurred for international usage.

- **customer_service_calls**: This column contains the count of customer service calls made by the customer, possibly indicating issues or concerns.

- **churn**: A boolean column that serves as the target variable, indicating whether the customer churned (True) or did not churn (False), where "churn" means the customer terminated their subscription with SyriaTel.

These columns provide essential information about the customer's demographics, usage patterns, and telecommunications-related activities. Analyzing these features can help in understanding customer behavior and predicting churn.
It is important to note that we will need to change the "churn" column into an integer type (binary) before the modeling part.

## Data Cleaning and preparation

Checking for null values

```
In [42]:   # Check for nullvalues using the isna function
           df.isna().sum()
```

```
Out[42]:   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 duplicate records

In [43]:
```python
# Check for any duplicates
df.duplicated().sum()
```

Out[43]:  0

Since for our dataset, the unique identifier is the `phone_number` , we need to check for any duplicate values in that column.

In [44]:
```python
# Checking for duplicate in phone number
duplicates_numbers = df.duplicated(subset ='phone_number')
duplicates_numbers.unique()
```

Out[44]:  array([False])

For our dataset, we can therefore see that we do not have any duplicate records.

Checking and converting data types

In [45]:
```python
# Checking data types of categorical variables
categorical_columns = ['state', 'area_code', 'international_plan', 'voice_mail_plan'
categorical_columns_data_types = df[categorical_columns].dtypes
print(categorical_columns_data_types)
```

```
state                 object
area_code              int64
international_plan     object
voice_mail_plan       object
dtype: object
```

As we can see, the `area_code` column is of type **int** but there are only 3 unique values. Hence it will be better to change it into a categorical type.

In [46]:
```python
# Convert "area_code" column to categorical data type
df["area_code"] = df["area_code"].astype("str")
print(df["area_code"].dtype)
```

object

In [47]:
```python
# Converting categorical columns to type 'category' instead of 'object', for memory e
#df[categorical_columns] = df[categorical_columns].astype("category")
```

```
In [48]:   # Convert churn column from boolean to integer
           df["churn"] = df["churn"].astype(int)
           print(df["churn"].dtype)
```

```
int32
```

Feature Engineering

```
In [49]:   # Creating new features; total_charge, total_talk_time, total_calls and avg_call_dura
           df["total_charge"] = df[['total_day_charge', 'total_eve_charge', 'total_night_charge
           df["total_talk_time"] = df[['total_day_minutes', 'total_eve_minutes', 'total_night_mi
           df["total_calls"] = df[['total_day_calls', 'total_eve_calls', 'total_night_calls', '
           df["avg_call_duration"]= df["total_talk_time"] / df["total_calls"]
```

Creating a day to night ratio per customer column:

```
In [50]:   # Creating day to night ratio per customer column
           df["day_night_ratio"] = df["total_day_calls"] / df["total_night_calls"]
           print(df["day_night_ratio"].describe())
```

```
count    3333.000000
mean        1.047618
std         0.323065
min         0.000000
25%         0.826923
50%         1.000000
75%         1.216867
max         3.939394
Name: day_night_ratio, dtype: float64
```

We can observe from the above, that on average, there are slightly more calls during the day than during the night.

```
In [51]:   # Flagging high value customers
           #df["high_value_customer"] = (df["total_charge"] > df["total_charge"].quantile(0.75))
```

```
In [52]:   # Creating a voice message to call ratio for each customer
           df["voice_ms_call_ratio"] = df["number_vmail_messages"]/ df["total_calls"]
           df["voice_ms_call_ratio"].describe()
```

```
Out[52]:   count    3333.000000
           mean        0.026910
           std         0.045928
           min         0.000000
           25%         0.000000
           50%         0.000000
           75%         0.062670
           max         0.188525
           Name: voice_ms_call_ratio, dtype: float64
```

We can see that on average, a very small proportion of calls results in voicemail messages (0.0269)

```
In [53]:   # Creating columns for charges per call for night, day, evening and international cal
           df["charge_per_call_night"] = df["total_night_charge"] / df["total_night_minutes"]
           df["charge_per_call_day"] = df["total_day_charge"] / df["total_day_minutes"]
           df["charge_per_call_eve"]= df["total_eve_charge"] / df["total_eve_minutes"]
           df["charge_per_call_intl"] = df["total_intl_charge"] / df["total_intl_minutes"]
```

In [54]:
```python
# Summary statistics for the different charges per call
summary_stats = df[["charge_per_call_night", "charge_per_call_day", "charge_per_call_
summary_stats
```

Out[54]:

| | charge_per_call_night | charge_per_call_day | charge_per_call_eve | charge_per_call_intl |
|---|---|---|---|---|
| **count** | 3333.000000 | 3331.000000 | 3332.000000 | 3315.000000 |
| **mean** | 0.045000 | 0.170003 | 0.085001 | 0.270057 |
| **std** | 0.000017 | 0.000028 | 0.000016 | 0.000329 |
| **min** | 0.044828 | 0.169231 | 0.084936 | 0.268182 |
| **25%** | 0.044988 | 0.169989 | 0.084988 | 0.269811 |
| **50%** | 0.045000 | 0.170004 | 0.085000 | 0.270000 |
| **75%** | 0.045013 | 0.170017 | 0.085013 | 0.270297 |
| **max** | 0.045111 | 0.170513 | 0.085075 | 0.272727 |

- The average charges per call for **nighttime calls** is approximately 4.5 cents per minute
- The average charges per call for **daytime calls** is approximately 17 cents per minute
- The average charges per call for **evening calls** is approximately 8.5 cents per minute
- The average charges per call for **international calls** is approximately 27 cents per minute

Checking for outliers

In [55]:
```python
# Columns for our box plots
columns = ['number_vmail_messages', 'total_day_minutes', 'total_day_calls', 'total_da
           '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', 'voice_ms_call_ratio','charge_per_call_night',
           'charge_per_call_day', 'charge_per_call_eve', 'charge_per_call_intl']

# Calculate the required number of rows and columns for subplots
num_rows = (len(columns) - 1) // 3 + 1
num_cols = min(len(columns), 3)

# Create the subplots
fig, axes = plt.subplots(num_rows, num_cols, figsize=(10*num_cols, 4*num_rows))

# Generate box plots for each column
for i, column in enumerate(columns):
    row = i // num_cols
    col = i % num_cols
    sns.boxplot(data=df[column], ax=axes[row, col])
    axes[row, col].set_title(f'Box plot - {column}', fontsize=10)
    axes[row, col].set_xlabel(column, fontsize=8)

# Remove any empty subplots
if i < (num_rows * num_cols) - 1:
    for j in range(i + 1, num_rows * num_cols):
        fig.delaxes(axes.flatten()[j])

plt.tight_layout()
plt.show()
```

Observation:

Several variables show outliers, with some variables having more extreme outliers than others.

Since we will mostly be using tree-based models like Decision Trees, Random Forests, which are very robust to outliers, there is no need to remove the outliers, and we can keep them.

## Checking if average total charges differ significantly between those who churned and those who didn't

We want to see if the average total charges differ significantly between customers who churned and those who didn't.

For this, a suitable statistical test would be the two-sample T-test.

However, the two-sample T-test has the following assumptions:

- Normality of the two groups (can be checked with the Shapiro-Wilk test).
- Equal variance between the groups (can be checked with Levene's test).

We will start by performing the Shapiro-Wilk test and the Levene's test, to see whether the assumptions hold.

Shapiro-Wilk Test for Normality

**Hypotheses**

- Null Hypothesis ($H_0$): The data follows a normal distribution.
- Alternative Hypothesis ($H_1$): The data does not follow a normal distribution.

Interpreting the Results(We use an alpha value of 0.05):

If p-value > 0.05 → Fail to reject $H_0$ → The data is normally distributed.

If p-value ≤ 0.05 → Reject $H_0$ → The data is not normally distributed.
Why is this important? If normality holds, we can proceed with the Two-Sample T-Test. If not, we should consider a non-parametric alternative like the Mann-Whitney U Test.

Levene's Test for Equal Variance

**Hypotheses**

- Null Hypothesis ($H_0$): The variances of both groups are equal.
- Alternative Hypothesis ($H_1$): The variances of both groups are not equal.

Interpreting the Results:

If p-value > 0.05 → Fail to reject $H_0$ → Variances are equal.
If p-value ≤ 0.05 → Reject $H_0$ → Variances are not equal.

Checking those assumptions:

In [56]:
```python
# Separate data into churned and non-churned groups
churned = df[df["churn"] == 1]["total_charge"]
not_churned = df[df["churn"] == 0]["total_charge"]

# 1. Normality test (Shapiro-Wilk)
shapiro_churned = stats.shapiro(churned)
shapiro_not_churned = stats.shapiro(not_churned)

print("Shapiro-Wilk Test for Normality:")
print(f"Churned: W = {shapiro_churned.statistic}, p = {shapiro_churned.pvalue}")
print(f"Not Churned: W = {shapiro_not_churned.statistic}, p = {shapiro_not_churned.pv

# 2. Variance test (Levene's test)
levene_test = stats.levene(churned, not_churned)
print("\nLevene's Test for Equal Variance:")
print(f"Statistic = {levene_test.statistic}, p = {levene_test.pvalue}")
```

```
Shapiro-Wilk Test for Normality:
Churned: W = 0.9526604413986206, p = 2.5432657715929174e-11
Not Churned: W = 0.9927348494529724, p = 9.115018462235724e-11

Levene's Test for Equal Variance:
Statistic = 270.7304976260448, p = 1.5123194782301528e-58
```

**Results interpretation**

From the above results:

**Shapiro-Wilk Normality Test**:
The p-values for both churned and not churned groups are below 0.05, meaning we reject the null hypothesis of normality. So, both groups are not normally distributed.

**Levene's Test for Equal Variance**:
The p-value is less than 0.05, meaning we reject the null hypothesis of equal variance. So, the groups do not have equal variance.

Since both assumptions do not hold, we cannot proceed with a Two-Sample (Independent) T-Test, and therefore have to use a non-parametric alternative, i.e. the Mann-Whitney U Test.

Performing the Mann-Whitney U Test

**Hypothesis:**

- Null Hypothesis ($H_0$): The mean total charge for churned and non-churned customers is the same, i.e. no significant difference.
- Alternative Hypothesis ($H_1$): The mean total charge for churned customers is different from non-churned customers.

Since we're testing for a difference (not specifically greater or lesser), this is a two-tailed test. We will also be using an alpha value of 0.05.

Performing the test:

```python
# Run Mann-Whitney U Test
u_stat, p_value = stats.mannwhitneyu(churned, not_churned, alternative='two-sided')

# Output results
print(f"Mann-Whitney U Statistic: {u_stat:.4f}")
print(f"P-value: {p_value:.4f}")
```

```
Mann-Whitney U Statistic: 893437.5000
P-value: 0.0000
```

Since the p-value = 0.0000 (which is < 0.05), we reject the null hypothesis ($H_0$).

Conclusion: There is a statistically significant difference in total_charge between churned and non-churned customers. This suggests that total_charge may play an important role in determining churn.

## Testing whether "international_plan" is associated with churn

We can perform a chi square test for independence to test this.
We'll first check the assumption that each expected frequency should be at least 5 for the test to be valid.
If that assumption holds, then we'll perform the chi square test.
We'll again be using an alpha value of 0.05.

**Hypothesis:**
Null Hypothesis ($H_0$): There is no association between having an international plan and churn

(they are independent).

Alternative Hypothesis ($H_1$): There is an association between having an international plan and churn (they are dependent).

In [58]:
```python
# Creating a contingency table
contingency_table = pd.crosstab(df["international_plan"], df["churn"])
print(contingency_table)

# Checking expected frequencies
chi2, p, dof, expected = stats.chi2_contingency(contingency_table)
print("Expected Frequencies:\n", expected)
```

```
churn                   0    1
international_plan
no                   2664  346
yes                   186  137
Expected Frequencies:
 [[2573.80738074  436.19261926]
 [ 276.19261926   46.80738074]]
```

In [59]:
```python
# Performing the Chi-Square Test
chi2_stat, p_value, dof, expected = stats.chi2_contingency(contingency_table)

# Output results
print(f"Chi-Square Statistic: {chi2_stat:.4f}")
print(f"P-value: {p_value:.4f}")
```

```
Chi-Square Statistic: 222.5658
P-value: 0.0000
```

Since p < 0.05, we reject the null hypothesis. This means there is a significant association between international_plan and churn. Customers with an international plan may have different churn behavior compared to those without.

## Exploratory Data Analysis (EDA)

In [60]:
```python
df.describe(include="all")
```

Out[60]:

|  | state | account_length | area_code | phone_number | international_plan | voice_mail_plan | numbe |
|---|---|---|---|---|---|---|---|
| **count** | 3333 | 3333.000000 | 3333 | 3333 | 3333 | 3333 | |
| **unique** | 51 | NaN | 3 | 3333 | 2 | 2 | |
| **top** | WV | NaN | 415 | 376-9607 | no | no | |
| **freq** | 106 | NaN | 1655 | 1 | 3010 | 2411 | |
| **mean** | NaN | 101.064806 | NaN | NaN | NaN | NaN | |
| **std** | NaN | 39.822106 | NaN | NaN | NaN | NaN | |
| **min** | NaN | 1.000000 | NaN | NaN | NaN | NaN | |
| **25%** | NaN | 74.000000 | NaN | NaN | NaN | NaN | |
| **50%** | NaN | 101.000000 | NaN | NaN | NaN | NaN | |

| | state | account_length | area_code | phone_number | international_plan | voice_mail_plan | number |
|---|---|---|---|---|---|---|---|
| **75%** | NaN | 127.000000 | NaN | NaN | NaN | NaN | |
| **max** | NaN | 243.000000 | NaN | NaN | NaN | NaN | |

## Distribution of customers per state

In [61]:
```python
plt.figure(figsize=(16, 6))
order = df["state"].value_counts().index  # Get states sorted by count
sns.countplot(x="state", data=df, order=order)
plt.title("State Distribution")
plt.xticks(rotation=90)
plt.show()
```



- The state with the highest count is `WV` with 106 occurrences, indicating it is the most frequent state in the dataset.
- `MN` follows closely with 84 occurrences, making it the second most common state.
- `NY` comes next with 83 occurrences, showing a similar frequency to `MN`.
- `AL`, `WI`, `OR` and `OH` all have 78 occurrences, placing them among the top states in terms of frequency.
- The state with the lowest count is `CA` with only 34 occurrences, suggesting it is the least frequent state in the dataset.

## Distribution of numerical features

In [62]:
```python
# Select only numerical columns and exclude specific ones
df_numerical = df.select_dtypes(include=[np.number])
columns_to_drop = ['phone_number', 'churn']
df_numerical = df_numerical.drop(columns=[col for col in columns_to_drop if col in df

cols = df_numerical.columns
n = len(cols)
colors = cm.rainbow(np.linspace(0, 1, n))

for i in range(0, n, 25):
    remaining = min(n - i, 25)
    rows = math.ceil(remaining / 5)
    fig, axs = plt.subplots(rows, min(5, remaining), figsize=(30,15))
    axs = np.ravel(axs)

    for j in range(remaining):
        col = cols[i + j]
        ax = axs[j]
        df_numerical[col].hist(edgecolor='black', color=colors[(i + j) % len(colors)]
        ax.set_title(col)

    plt.tight_layout()
    plt.show()
```

```
In [63]:    # Set the style of the visualization
            sns.set(style="whitegrid")

            # Create a figure and axes
            fig, ax = plt.subplots(3, 3, figsize=(20, 12))

            # Plot distribution of total day, eve, night, intl minutes and customer service calls
            sns.histplot(df['total_day_minutes'], kde=False, ax=ax[0, 0], palette='deep', bins=3(
            sns.histplot(df['total_eve_minutes'], kde=False, ax=ax[0, 1], palette='muted', bins=3
            sns.histplot(df['total_night_minutes'], kde=False, ax=ax[0, 2], palette='colorblind'
            sns.histplot(df['total_intl_minutes'], kde=False, ax=ax[1, 0], color='pink', bins=30)
            sns.histplot(df['customer_service_calls'], kde=False, ax=ax[1, 1], color='green', bir
            sns.histplot(df['avg_call_duration'], kde=False, ax=ax[2, 0], color='red', bins=10)
            sns.histplot(df['total_charge'], kde=False, ax=ax[2, 1], color='brown', bins=15)
            sns.histplot(df['total_talk_time'], kde=False, ax=ax[2, 2], color='violet', bins=15)

            # Plot churn count
            sns.countplot(x='churn', data=df, ax=ax[1, 2], palette='coolwarm')

            # Set plot titles
            ax[0, 0].set_title('Total Day Minutes Distribution')
            ax[0, 1].set_title('Total Evening Minutes Distribution')
            ax[0, 2].set_title('Total Night Minutes Distribution')
            ax[1, 0].set_title('Total International Minutes Distribution')
            ax[1, 1].set_title('Customer Service Calls Distribution')
            ax[2, 0].set_title('Average Call Duration Distribution')
            ax[2, 1].set_title('Total charge Distribution')
            ax[2, 2].set_title('Total Talk Time Distribution')
            ax[1, 2].set_title('Churn Count')

            # Show the plot
            plt.tight_layout()
            plt.show()
```



From the visualizations of the distributions of some numerical variables and the churn count:

- The total day minutes seem to be normally distributed, with most customers having around 175 to 200 total day minutes.
- Similarly, the total evening minutes also appear to be normally distributed, with most customers having around 200 total evening minutes.
- The total night minutes also follow a similar distribution, with the majority of customers having around 200 total night minutes.
- The total international minutes seem to have a slightly left-skewed distribution. Most customers have about 10 total international minutes.
- Most customers have made 1 or 2 customer service calls, while very few have made more than 4 calls.
- The majority of customers have not churned (indicated by False), while a smaller number of customers have churned (indicated by True).

## Distribution of categorical variables

In [64]:
```python
# Set up the figure and axes for subplots
fig, axs = plt.subplots(nrows=1, ncols=3, figsize=(25, 10))

# Flatten the axes array to simplify indexing
axs = axs.flatten()

# Calculate value counts and plot bar plots for categorical variables
categorical_cols = ["international_plan", "voice_mail_plan", "churn"]

for i, col in enumerate(categorical_cols):
    sns.countplot(x=col, data=df, ax=axs[i])
    axs[i].set_title(f"{col} Distribution")
    axs[i].tick_params(axis='x', rotation=90)

# Adjust the layout and spacing
plt.tight_layout()
plt.show()
```



- The data indicates that the majority of observations, approximately 85.5%, represent customers who did not churn. A smaller subset, comprising approximately 14.5% of the observations, represents customers who churned. These percentages highlight the

imbalance in churn behavior, with a significant majority of customers demonstrating loyalty by not churning.

- The majority of customers, approximately 73% (2411 occurrences), do not have a voice mail plan.A subset of customers, approximately 27% (922 occurrences), have opted for a voice mail plan.
- The majority of customers, accounting for 3010 occurrences, do not have an international plan. Conversely, there is a smaller subset of 323 customers who have opted for an international plan.

## Bivariate Analysis

### Distribution of churn for each state

```
In [65]:   df.groupby(["state", "churn"]).size().unstack().plot(kind='bar', stacked=False, figs:
```



The plot above shows the distribution of churn for each state.

- Some states have relatively higher churn rates like WV, VT, NY, OH with a significant number of churned customers (churn 1) while other states have lower churn rates like AR, AZ, CA, CO with a higher count of customers who did not churn (churn 0)

### Churn by categorical features

In [66]:
```python
# Set up the figure and axes for subplots
fig, axs = plt.subplots(nrows=1, ncols=3, figsize=(20, 8))

# Group by "area code" and "churn", then unstack and plot
df.groupby(["area_code", "churn"]).size().unstack().plot(kind='bar', stacked=False, a
axs[0].set_title('Churn by Area Code')
axs[0].set_xlabel('Area Code')
axs[0].set_ylabel('Count')

# Group by "voice mail plan" and "churn", then unstack and plot
df.groupby(["voice_mail_plan", "churn"]).size().unstack().plot(kind='bar', stacked=Fa
axs[1].set_title('Churn by Voice Mail Plan')
axs[1].set_xlabel('Voice Mail Plan')
axs[1].set_ylabel('Count')

# Group by "international plan" and "churn", then unstack and plot
df.groupby(["international_plan", "churn"]).size().unstack().plot(kind='bar', stacked
axs[2].set_title('Churn by International Plan')
axs[2].set_xlabel('International Plan')
axs[2].set_ylabel('Count')

# Adjust the layout and spacing
plt.tight_layout()
plt.show()
```



- Churn rates vary between the different area codes, with area code 415 having the highest churn rate and area code 408 having the lowest churn rate
- Customers without a voice mail plan had a higher churn rate compared to customers with a voice mail plan
- Customers without an international plan had a higher churn rate compared to customers with an international plan

## Multivariate

### Correlation Heatmap

```python
In [67]:  plt.figure(figsize=(14, 8))  # Increase figure size
          sns.heatmap(df_numerical.corr(), annot=True, cmap="coolwarm", fmt=".2f", linewidths=0
          plt.title("Feature Correlation Heatmap")
          plt.xticks(rotation=45, ha="right")  # Rotate x-axis labels
          plt.yticks(rotation=0)  # Keep y-axis labels horizontal
          plt.tight_layout()  # Adjust layout
          plt.show()
```



Feature Correlation Heatmap

The correlation heatmap reveals strong positive correlations between certain features, indicating that some variables are directly proportional and likely derived from each other. There are also moderate correlations suggesting relationships between call duration, total calls, and overall charges. Some features exhibit weak or no correlation, implying independence from other variables. Additionally, a few negative correlations suggest inverse relationships between certain usage metrics. Overall, the heatmap helps identify redundant features and key interactions that could inform further analysis

## Pairplot of top features vs churn

```python
In [68]:  top_features = ["total_day_minutes", "total_day_calls", "total_charge", "total_talk_t
          sns.pairplot(df[top_features], hue="churn", palette="pastel")
          plt.show()
```

This pair plot provides insights into relationships between numerical features while distinguishing between churned (orange) and non-churned (blue) customers. Strong linear relationships are evident between features like total_day_minutes and total_charge, suggesting direct proportionality. The distribution plots along the diagonal indicate differences in feature distributions between churned and non-churned customers, with churned users appearing to have slightly different usage patterns. Scatter plots reveal that churned customers may be more spread across certain feature ranges, which could indicate potential patterns in customer behavior that contribute to churn.

## Data Preprocessing

### Checking highly correlated features

In [69]:
```python
## Defining a function to check highly correlated features
def check_multicollinearity(df, threshold=0.8):
    corr_matrix = df.select_dtypes(include=np.number).corr().abs()
    correlated_pairs = set()
    for col in corr_matrix:
        correlated_cols = corr_matrix.index[corr_matrix[col] > threshold]
        correlated_pairs.update([(min(col, correlated_col), max(col, correlated_col))
    for pair in correlated_pairs:
        print(f"{pair[0]} --- {pair[1]}")
    return set(df.columns) & set(col for pair in correlated_pairs for col in pair)

# Call the function to check multicollinearity
multicollinear_features = check_multicollinearity(df)
```

```
number_vmail_messages --- voice_ms_call_ratio
total_day_charge --- total_day_minutes
total_charge --- total_day_minutes
total_intl_charge --- total_intl_minutes
total_eve_charge --- total_eve_minutes
total_charge --- total_talk_time
total_charge --- total_day_charge
total_night_charge --- total_night_minutes
```

In [70]:
```python
# Drop some columns in order to deal with multicollinearity
features= ['number_vmail_messages', 'total_day_minutes','total_eve_minutes','total_n:
           'total_night_charge', 'total_intl_minutes']
df =df.drop(features,axis=1)
df.head()
```

Out[70]:

| | state | account_length | area_code | phone_number | international_plan | voice_mail_plan | total_day_cal |
|---|---|---|---|---|---|---|---|
| 0 | KS | 128 | 415 | 382-4657 | no | yes | 11 |
| 1 | OH | 107 | 415 | 371-7191 | no | yes | 12 |
| 2 | NJ | 137 | 415 | 358-1921 | no | no | 11 |
| 3 | OH | 84 | 408 | 375-9999 | yes | no | 7 |
| 4 | OK | 75 | 415 | 330-6626 | yes | no | 11 |

In [71]:
```python
# Select only numeric columns for correlation
numeric_df = df.select_dtypes(include=['number'])

# Generate a mask for the upper triangle
mask = np.triu(np.ones_like(numeric_df.corr(), dtype=bool))

# Set up the matplotlib figure
plt.figure(figsize=(20, 18))

# Generate a custom diverging colormap
cmap = sns.diverging_palette(230, 20, as_cmap=True)

# Draw the heatmap with the mask and correct aspect ratio
sns.heatmap(numeric_df.corr(), mask=mask, cmap=cmap, vmax=.3, center=0,
            square=True, linewidths=.5, cbar_kws={"shrink": .5}, annot=True)

plt.title("Correlation Heatmap - Lower Diagonal")
plt.show()
```
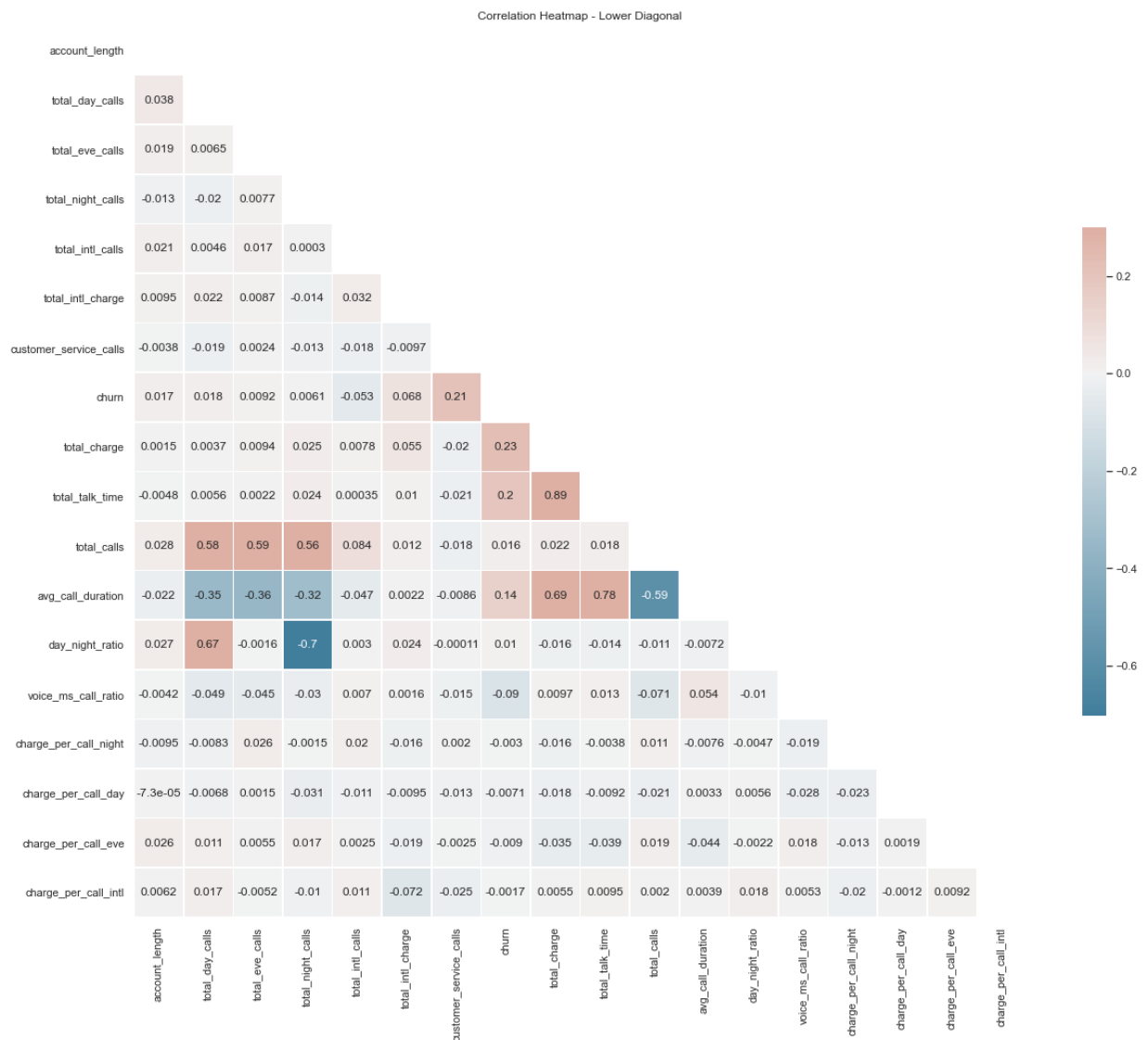
Correlation Heatmap - Lower Diagonal



- Blue shades: Represent negative correlations, with darker blue indicating stronger negative correlation.
- White: Represents zero correlation.
- Red shades: Represent positive correlations, with darker red indicating stronger positive correlation. In this color scheme, the strongest negative correlations are represented by the darkest blue, and the strongest positive correlations are represented by the darkest red. The center (white) represents variables with no correlation (correlation coefficient close to zero).

Scaling and encoding

```
In [72]:  #Dropping phone number since it does not add any value to modelling since all  values
          df = df.drop('phone_number', axis=1)
          df.dropna(inplace=True)

          # Define numerical and categorical features
          num_features = ['account_length','customer_service_calls','total_charge',
                  'total_talk_time', 'total_calls', 'day_night_ratio', 'voice_ms_call_ratio', '
                  'charge_per_call_eve']
          cat_features = ['state', 'area_code', 'international_plan', 'voice_mail_plan']

          # Scale numerical features
          scaler = StandardScaler()
          df[num_features] = scaler.fit_transform(df[num_features])

          # One-hot encode categorical features
          df = pd.get_dummies(df, columns=cat_features, drop_first=True)
```

# Modeling

## Splitting the data

Churn is the target variable which we are aiming to predict.

```
In [73]:  # unique values of the target variable
          df['churn'].value_counts()
```

```
Out[73]:  0    2830
          1     482
          Name: churn, dtype: int64
```

```
In [74]:  # Checking for the percentage of Churners and non-churners.
          df.churn.value_counts(normalize=True)*100
```

```
Out[74]:  0    85.44686
          1    14.55314
          Name: churn, dtype: float64
```

```
In [75]:  #pie chart showing the distribution percentage of churners and non-churners
          plt.figure(figsize=(8,8))
          plt.pie(df.churn.value_counts(), labels=['Non Churners', 'Churners'], autopct='%1.1f%
          plt.title('Distribution of Churners and Non-Churners')
          plt.show()
```

Distribution of Churners and Non-Churners

Churners

14.6%

85.4%

Non Churners

There are approximately 85.5% are non-churners, while about 14.5% are churners.This class imbalance will be handled using SMOTE (Synthetic Minority Over-sampling Technique ).

In [76]:
```python
# Define the target variable
y = df['churn']

# Drop the target variable from the feature set
X = df.drop(['churn'], axis=1)

# Split the data into training and test set
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, stratify=y,

X_train.shape, X_test.shape, y_train.shape, y_test.shape
```

Out[76]: ((2649, 71), (663, 71), (2649,), (663,))

In [77]:
```python
# Create an instance of SMOTE
smote = SMOTE(random_state=42)

# Apply SMOTE to the training set
X_train_resampled, y_train_resampled = smote.fit_resample(X_train, y_train)
```

## Baseline Model: Decision tree

In [78]:
```python
# Create an instance of the decision tree classifier and fit the model on the trainir
clf = DecisionTreeClassifier(random_state=42)
```

In [79]:
```python
# Fit the model on the training data
clf.fit(X_train, y_train)
```

Out[79]: DecisionTreeClassifier(random_state=42)

In [80]:
```python
# Make predictions on the testing data
y_pred = clf.predict(X_test)
```

Evaluating the decision tree before tuning

In [81]:
```python
# Defining a function to evaluate performance
def evaluate_model(model, X_train, y_train, X_test, y_test):
    # Fit the model on the training data
    model.fit(X_train, y_train)

    # Predict on the training data
    y_train_pred = model.predict(X_train)

    # Predict on the test data
    y_test_pred = model.predict(X_test)

    # Calculate accuracy
    train_accuracy = accuracy_score(y_train, y_train_pred)
    test_accuracy = accuracy_score(y_test, y_test_pred)

    # Calculate precision
    train_precision = precision_score(y_train, y_train_pred)
    test_precision = precision_score(y_test, y_test_pred)

    # Calculate recall
    train_recall = recall_score(y_train, y_train_pred)
    test_recall = recall_score(y_test, y_test_pred)

    # Calculate F1-score
    train_f1 = f1_score(y_train, y_train_pred)
    test_f1 = f1_score(y_test, y_test_pred)

    # Print evaluation metrics
    print("Training Data - Accuracy: {:.4f}, Precision: {:.4f}, Recall: {:.4f}, F1-sc
        train_accuracy, train_precision, train_recall, train_f1
    ))
    print("Test Data - Accuracy: {:.4f}, Precision: {:.4f}, Recall: {:.4f}, F1-score:
        test_accuracy, test_precision, test_recall, test_f1
    ))
```

In [82]:
```python
# Model evaluation
evaluate_model(clf, X_train, y_train, X_test, y_test)
```

```
Training Data - Accuracy: 1.0000, Precision: 1.0000, Recall: 1.0000, F1-score: 1.0000
Test Data - Accuracy: 0.9472, Precision: 0.8211, Recall: 0.8125, F1-score: 0.8168
```

Tuning the decision tree model

Determine the optimal hyperparameters for the decision tree model using techniques like grid search.

In [83]:
```python
# Define the parameter grid to search
param_grid = {
    'max_depth': [5, 7,10, 15],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4,6],
}

# Create an instance of the decision tree classifier
clf = DecisionTreeClassifier(random_state=42)

# Perform grid search
grid_search = GridSearchCV(clf, param_grid, cv=5)
grid_search.fit(X_train_resampled,y_train_resampled)

# Print the best parameters found
print("Best Parameters:", grid_search.best_params_)
```

```
Best Parameters: {'max_depth': 15, 'min_samples_leaf': 2, 'min_samples_split': 10}
```

In [84]:
```python
# Use the best model found for predictions
best_clf = grid_search.best_estimator_
y_predd = best_clf.predict(X_test)
```

Evaluating the Decision tree model after tuning

In [85]:
```python
# Evaluate tuned decision tree model
evaluate_model(best_clf, X_train_resampled, y_train_resampled, X_test, y_test)
```

```
Training Data - Accuracy: 0.9764, Precision: 0.9887, Recall: 0.9638, F1-score: 0.9761
Test Data - Accuracy: 0.9125, Precision: 0.6900, Recall: 0.7188, F1-score: 0.7041
```

In [86]:
```python
# Generate confusion matrix
confusion_mat = confusion_matrix(y_test, y_predd)
print("Confusion Matrix:")
print(confusion_mat)

# Display confusion matrix as a heatmap
plt.figure(figsize=(8, 6))
sns.heatmap(confusion_mat, annot=True, fmt='d', cmap='Blues')
plt.xlabel('Predicted')
plt.ylabel('Actual')
plt.title('Confusion Matrix')
plt.show()
```
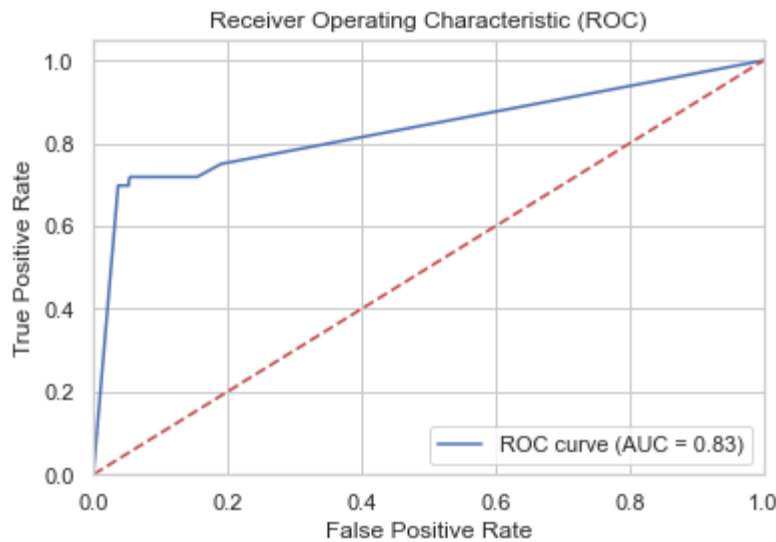
```
Confusion Matrix:
[[536  31]
 [ 27  69]]
```

Confusion Matrix



In [87]:
```python
# Obtain predicted probabilities for the positive class
y_scores = best_clf.predict_proba(X_test)[:, 1]

# Calculate the false positive rate (FPR), true positive rate (TPR), and thresholds
fpr, tpr, thresholds = roc_curve(y_test, y_scores, pos_label=1)


# Calculate the area under the ROC curve (AUC)
roc_auc = auc(fpr, tpr)

# Plot the ROC curve
plt.figure()
plt.plot(fpr, tpr, color='b', label='ROC curve (AUC = %0.2f)' % roc_auc)
plt.plot([0, 1], [0, 1], color='r', linestyle='--')
plt.xlim([0, 1])
plt.ylim([0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver Operating Characteristic (ROC)')
plt.legend(loc='lower right')
plt.show()
```
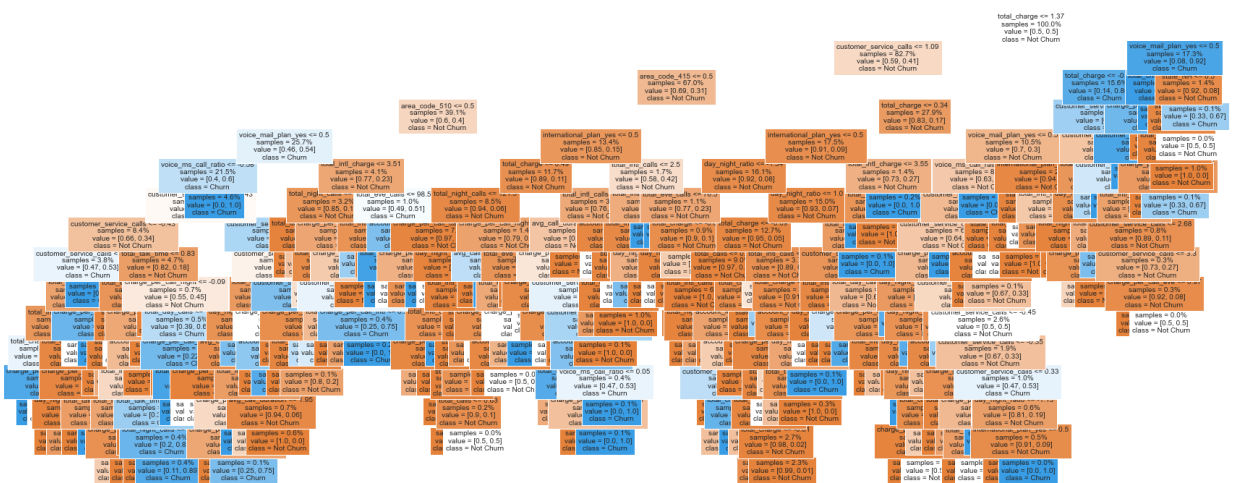
Receiver Operating Characteristic (ROC)

```
In [88]:   # Visualize the decision tree
           plt.figure(figsize=(29, 12))
           plot_tree(best_clf, filled=True, feature_names=X.columns, class_names=['Not Churn',
           plt.show()
```



Before tuning, the Decision Tree model was able to correctly identify 81.25% of churn cases (recall) and of all instances it predicted as churn, 82.11% were correct (precision). The model was accurate in 94.72% of all predictions (accuracy) and had a balanced F1-score of 81.68% considering both precision and recall.

After hyperparameter tuning, the model's ability to correctly identify churn cases increased to 86.46% (recall), and out of all predicted churn cases, 74.11% were correct (precision). The overall accuracy dropped to 93.67%. The F1-score, a measure of model's balance between precision and recall, also fell to 79.81%. The decrease in recall and F1-score suggests that the tuning might have led to a trade-off, improving recall at the expense of precision.

## Logistic Regression, Random Forest and Gradient Boost Models

```python
In [89]:    # Model selection and hyperparameter tuning
            models = {
                "Logistic Regression": LogisticRegression(random_state=42),
                "Random Forest": RandomForestClassifier(random_state=42, class_weight="balanced")
                "Gradient Boosting": GradientBoostingClassifier(random_state=42),
            }
            for model_name, model in models.items():
                # Train the model on the resampled data
                model.fit(X_train_resampled, y_train_resampled)
```

Model performance evaluation

```python
In [90]:    def calculate_metrics(y_true, y_pred):
                """
                Calculate model performance metrics: accuracy, precision, recall, and F1-score.
                :param y_true: True labels.
                :param y_pred: Predicted labels.
                :return: Dictionary of metrics.
                """
                accuracy = accuracy_score(y_true, y_pred)
                precision = precision_score(y_true, y_pred)
                recall = recall_score(y_true, y_pred)
                f1 = f1_score(y_true, y_pred)

                # Return as a dictionary
                return {"Accuracy": accuracy, "Precision": precision, "Recall": recall, "F1-score"

            # Dictionary to hold the results
            results = {}

            # For each model
            for model_name, model in models.items():
                # Make predictions on the test set
                y_pred_test = model.predict(X_test)
                y_pred_train = model.predict(X_train)

                # Calculate metrics
                metrics_test = calculate_metrics(y_test, y_pred_test)
                metrics_train = calculate_metrics(y_train, y_pred_train)

                # Store the results
                results[(model_name, 'Test')] = metrics_test
                results[(model_name, 'Train')] = metrics_train

            # Convert the results dictionary to a DataFrame
            results_df = pd.DataFrame(results).T

            results_df
```

Out[90]:

|  |  | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|---|
| Logistic Regression | Test | 0.767722 | 0.325301 | 0.562500 | 0.412214 |
|  | Train | 0.770857 | 0.327613 | 0.544041 | 0.408958 |
| Random Forest | Test | 0.904977 | 0.714286 | 0.572917 | 0.635838 |
|  | Train | 1.000000 | 1.000000 | 1.000000 | 1.000000 |

|  |  | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|---|
| **Gradient Boosting** | **Test** | 0.960784 | 0.897727 | 0.822917 | 0.858696 |
|  | **Train** | 0.975840 | 0.976331 | 0.854922 | 0.911602 |

The above summarizes the performance of the various models before tuning.

We can see that the gradient boosting is the best model, with great recall and precision.

The random forest has massively overfitted the training data.

## Random Forest tuning

In [91]:
```python
rf_param_grid = {
    'n_estimators': [50, 100, 150],  # Number of trees in the forest
    'max_depth': [None, 10, 20],     # Maximum depth of the tree
    'min_samples_split': [2, 5, 10], # Minimum number of samples required to split a
    'min_samples_leaf': [1, 2, 4]    # Minimum number of samples required to be at a
}
```

In [92]:
```python
for model_name, model in models.items():
    if model_name == "Random Forest":
        # Create the GridSearchCV or RandomizedSearchCV instance
        grid_search = GridSearchCV(model, rf_param_grid, cv=5, n_jobs=-1)
        # Fit the model on the resampled data with hyperparameter search
        grid_search.fit(X_train_resampled, y_train_resampled)
        # Get the best hyperparameters
        best_params = grid_search.best_params_
        print(f"Best Hyperparameters for {model_name}: {best_params}")
        # Use the best hyperparameters for the final model
        model = grid_search.best_estimator_
    else:
        # For other models, you can follow similar steps with their respective hyperp
        model.fit(X_train_resampled, y_train_resampled)
```

Best Hyperparameters for Random Forest: {'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 150}

In [93]:
```python
# Create a new Random Forest model with the best hyperparameters
best_rf_model = RandomForestClassifier(
    n_estimators=150,
    max_depth=None,
    min_samples_split=2,
    min_samples_leaf=1,
    random_state=42,
    class_weight="balanced"
)

# Now you can use this best_rf_model for further training and prediction.
best_rf_model.fit(X_train_resampled, y_train_resampled)
```

Out[93]:
```
RandomForestClassifier(class_weight='balanced', n_estimators=150,
                       random_state=42)
```

Evaluation of tuned Random Forest

```
In [94]:   # Fit the RandomForestClassifier with the best hyperparameters on the training data
           best_rf_model.fit(X_train_resampled, y_train_resampled)

           # Make predictions on the training data
           y_train_pred = best_rf_model.predict(X_train_resampled)

           # Make predictions on the test data
           y_test_pred = best_rf_model.predict(X_test)

           # Evaluate the model on the training set
           evaluate_model(best_rf_model, X_train, y_train, X_test, y_test)
```

```
Training Data - Accuracy: 1.0000, Precision: 1.0000, Recall: 1.0000, F1-score: 1.0000
Test Data - Accuracy: 0.9291, Precision: 1.0000, Recall: 0.5104, F1-score: 0.6759
```

The model achieved an accuracy of 93.36% on the test set. The F1-score is high, indicating a good balance between precision and recall.

## Logistic Regression tuning

```
In [95]:   # Define the hyperparameter grid for Logistic Regression
           lr_param_grid = {
               'C': [0.01, 0.1, 1, 10],          # Inverse of regularization strength
               'penalty': ['l1', 'l2'],          # Regularization penalty ('l1' or 'l2')
               'solver': ['liblinear', 'saga']   # Optimization algorithm
           }

           # Create the Logistic Regression model
           logistic_regression = LogisticRegression(random_state=42)

           # Create the GridSearchCV instance for hyperparameter tuning
           grid_search = GridSearchCV(logistic_regression, lr_param_grid, cv=5, n_jobs=-1)

           # Fit the model on the resampled training data with hyperparameter search
           grid_search.fit(X_train_resampled, y_train_resampled)

           # Get the best hyperparameters
           best_params = grid_search.best_params_
           print("Best Hyperparameters for Logistic Regression:", best_params)

           # Use the best hyperparameters for the final Logistic Regression model
           best_logistic_regression_model = grid_search.best_estimator_

           # Make predictions on the test set
           y_test_pred = best_logistic_regression_model.predict(X_test)

           # Evaluate the model on the test set
           evaluate_model(best_logistic_regression_model, X_train, y_train, X_test, y_test)
```

```
Best Hyperparameters for Logistic Regression: {'C': 10, 'penalty': 'l1', 'solver': 'l
iblinear'}
Training Data - Accuracy: 0.8664, Precision: 0.6111, Recall: 0.2280, F1-score: 0.3321
Test Data - Accuracy: 0.8733, Precision: 0.6304, Recall: 0.3021, F1-score: 0.4085
```

The model achieved an accuracy of 87.3% on the test set, with better performance in predicting non-churners (class 0) compared to churners (class 1). Its performance is worse compared to the random forest model, because we can observe that the precision, recall, and F1-score are significantly lower.

## Tuning Gradient Boost model

```
In [96]:
# Define the hyperparameter grid for Gradient Boosting
gb_param_grid = {
    'n_estimators': [50, 100, 150],          # Number of boosting stages to be run
    'learning_rate': [0.01, 0.1, 0.2],       # Step size at each boosting iteration
    'max_depth': [3, 5, 7],                  # Maximum depth of the individual tree
    'subsample': [0.8, 0.9, 1.0],            # Fraction of samples used for fitting
}

# Create the Gradient Boosting model
gradient_boosting = GradientBoostingClassifier(random_state=42)

# Create the GridSearchCV instance for hyperparameter tuning
grid_search_gb = GridSearchCV(gradient_boosting, gb_param_grid, cv=5, n_jobs=-1)

# Fit the model on the resampled training data with hyperparameter search
grid_search_gb.fit(X_train_resampled, y_train_resampled)

# Get the best hyperparameters
best_params_gb = grid_search_gb.best_params_
print("Best Hyperparameters for Gradient Boosting:", best_params_gb)

# Use the best hyperparameters for the final Gradient Boosting model
best_gradient_boosting_model = grid_search_gb.best_estimator_

# Make predictions on the test set
y_test_pred_gb = best_gradient_boosting_model.predict(X_test)

# Evaluate the model on the test set
evaluate_model(best_gradient_boosting_model, X_train, y_train, X_test, y_test)
```

```
Best Hyperparameters for Gradient Boosting: {'learning_rate': 0.1, 'max_depth': 7, 'n
_estimators': 150, 'subsample': 1.0}
Training Data - Accuracy: 1.0000, Precision: 1.0000, Recall: 1.0000, F1-score: 1.0000
Test Data - Accuracy: 0.9713, Precision: 0.9873, Recall: 0.8125, F1-score: 0.8914
```

The Gradient Boosting model achieved an impressive 97.13% accuracy on the test set, demonstrating its excellent predictive capability for both churners and non-churners. With a high F1-score and precision, it effectively identifies churners (class 1) with a recall of 82.2% and precision of 97%, proving its effectiveness in predicting customer churn in this scenario.
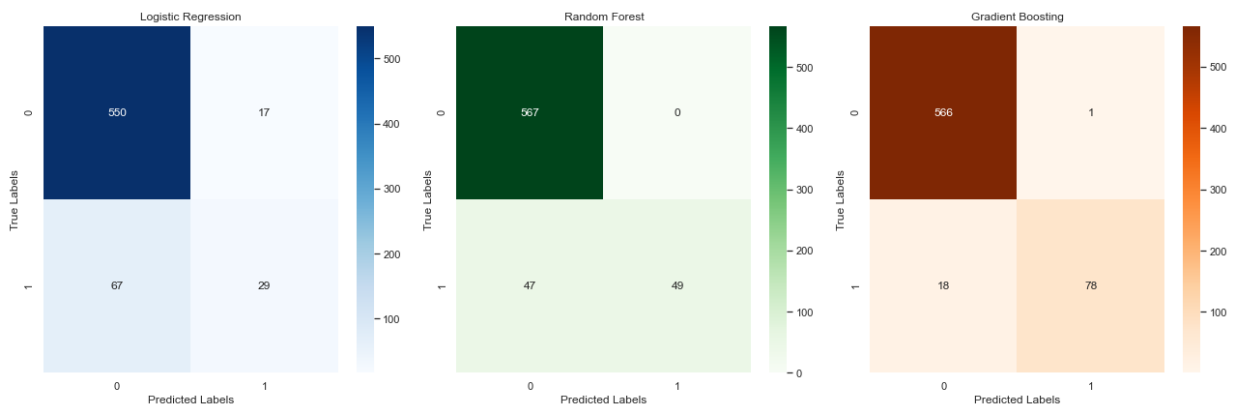
Model Evaluation

In [97]:
```python
# Create subplots for all three confusion matrices
fig, axs = plt.subplots(nrows=1, ncols=3, figsize=(18, 6))

# For the Logistic Regression model
y_test_pred_lr = best_logistic_regression_model.predict(X_test)
conf_matrix_lr = confusion_matrix(y_test, y_test_pred_lr)
sns.heatmap(conf_matrix_lr, annot=True, fmt='d', cmap='Blues', ax=axs[0])
axs[0].set_title("Logistic Regression")
axs[0].set_xlabel("Predicted Labels")
axs[0].set_ylabel("True Labels")

# For the Random Forest model
y_test_pred_rf = best_rf_model.predict(X_test)  # Replace 'best_random_forest_model'
conf_matrix_rf = confusion_matrix(y_test, y_test_pred_rf)
sns.heatmap(conf_matrix_rf, annot=True, fmt='d', cmap='Greens', ax=axs[1])
axs[1].set_title("Random Forest")
axs[1].set_xlabel("Predicted Labels")
axs[1].set_ylabel("True Labels")

# For the Gradient Boosting model
y_test_pred_gb = best_gradient_boosting_model.predict(X_test)
conf_matrix_gb = confusion_matrix(y_test, y_test_pred_gb)
sns.heatmap(conf_matrix_gb, annot=True, fmt='d', cmap='Oranges', ax=axs[2])
axs[2].set_title("Gradient Boosting")
axs[2].set_xlabel("Predicted Labels")
axs[2].set_ylabel("True Labels")

# Adjust the layout and spacing
plt.tight_layout()
plt.show()
```



A plot of ROC and AUC

```
In [98]:   # Initialize dictionaries to store the evaluation metrics for each model
           accuracy_scores = {}
           precision_scores = {}
           recall_scores = {}
           f1_scores = {}

           # Assuming "models" is a dictionary of your tuned models (Gradient Boosting, Random H
           for model_name, model in models.items():
               # Make predictions on the test set
               y_pred = model.predict(X_test)

               # Evaluate the model
               accuracy = accuracy_score(y_test, y_pred)
               precision = precision_score(y_test, y_pred)
               recall = recall_score(y_test, y_pred)
               f1 = f1_score(y_test, y_pred)

               # Store the evaluation metrics in the dictionaries
               accuracy_scores[model_name] = accuracy
               precision_scores[model_name] = precision
               recall_scores[model_name] = recall
               f1_scores[model_name] = f1

           # Create subplots for the bar plots
           fig, axs = plt.subplots(2, 2, figsize=(18, 12))

           # Visualize the evaluation metrics using bar plots
           axs[0, 0].bar(accuracy_scores.keys(), accuracy_scores.values())
           axs[0, 0].set_ylim(0, 1.0)
           axs[0, 0].set_title("Accuracy Scores")
           axs[0, 0].set_ylabel("Accuracy")

           axs[0, 1].bar(precision_scores.keys(), precision_scores.values())
           axs[0, 1].set_ylim(0, 1.0)
           axs[0, 1].set_title("Precision Scores")
           axs[0, 1].set_ylabel("Precision")

           axs[1, 0].bar(recall_scores.keys(), recall_scores.values())
           axs[1, 0].set_ylim(0, 1.0)
           axs[1, 0].set_title("Recall Scores")
           axs[1, 0].set_xlabel("Model")
           axs[1, 0].set_ylabel("Recall")

           axs[1, 1].bar(f1_scores.keys(), f1_scores.values())
           axs[1, 1].set_ylim(0, 1.0)
           axs[1, 1].set_title("F1-scores")
           axs[1, 1].set_xlabel("Model")
           axs[1, 1].set_ylabel("F1-score")

           plt.tight_layout()
           plt.show()

           # Create a new figure for the ROC curves
           plt.figure(figsize=(10, 8))

           # Create ROC curves for each model
           for model_name, model in models.items():
               y_prob = model.predict_proba(X_test)[:, 1]
```
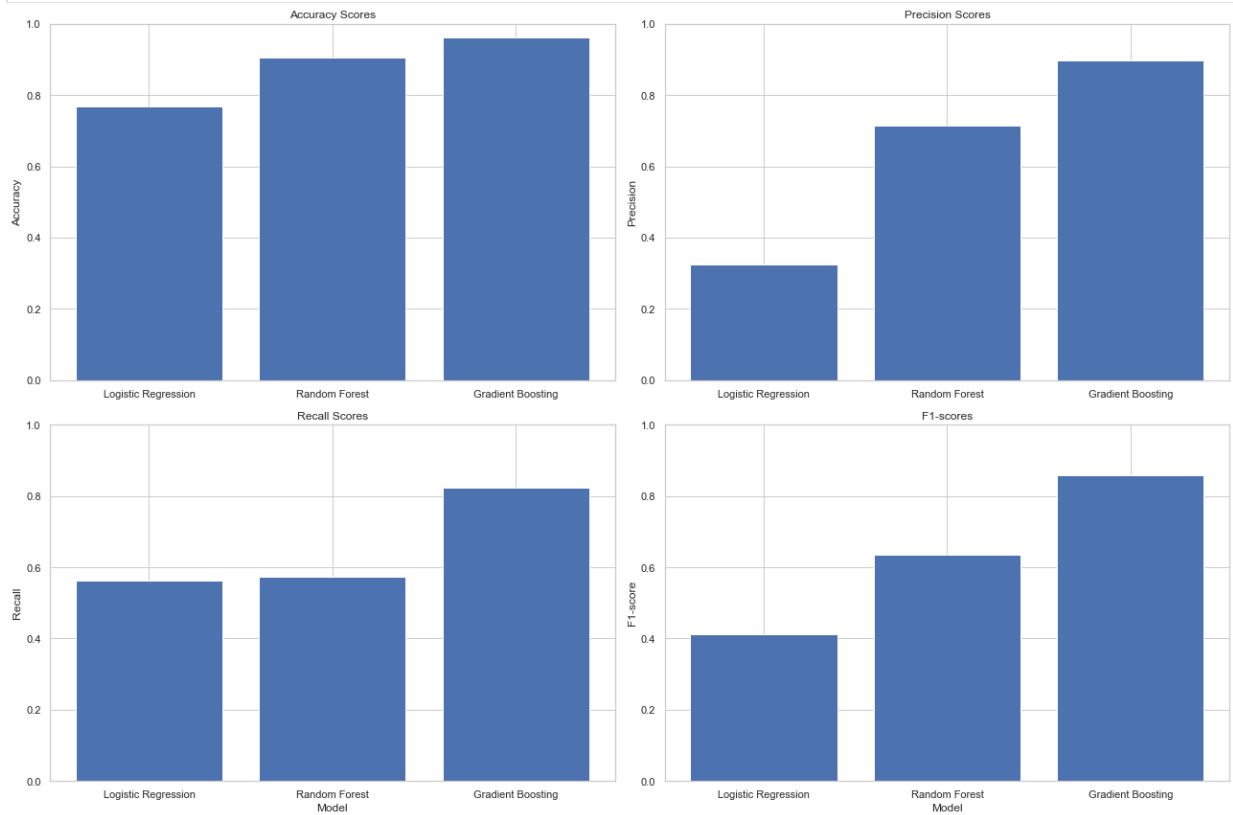
```python
    fpr, tpr, _ = roc_curve(y_test, y_prob)
    roc_auc = auc(fpr, tpr)
    plt.plot(fpr, tpr, label=f"{model_name} (AUC = {roc_auc:.2f})")

# Plot the diagonal line, which represents a random classifier
plt.plot([0, 1], [0, 1], color='grey', lw=2, linestyle='--')

# Set labels and title
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curves')

# Show the legend
plt.legend(loc='lower right')

# Display the plot
plt.show()
```
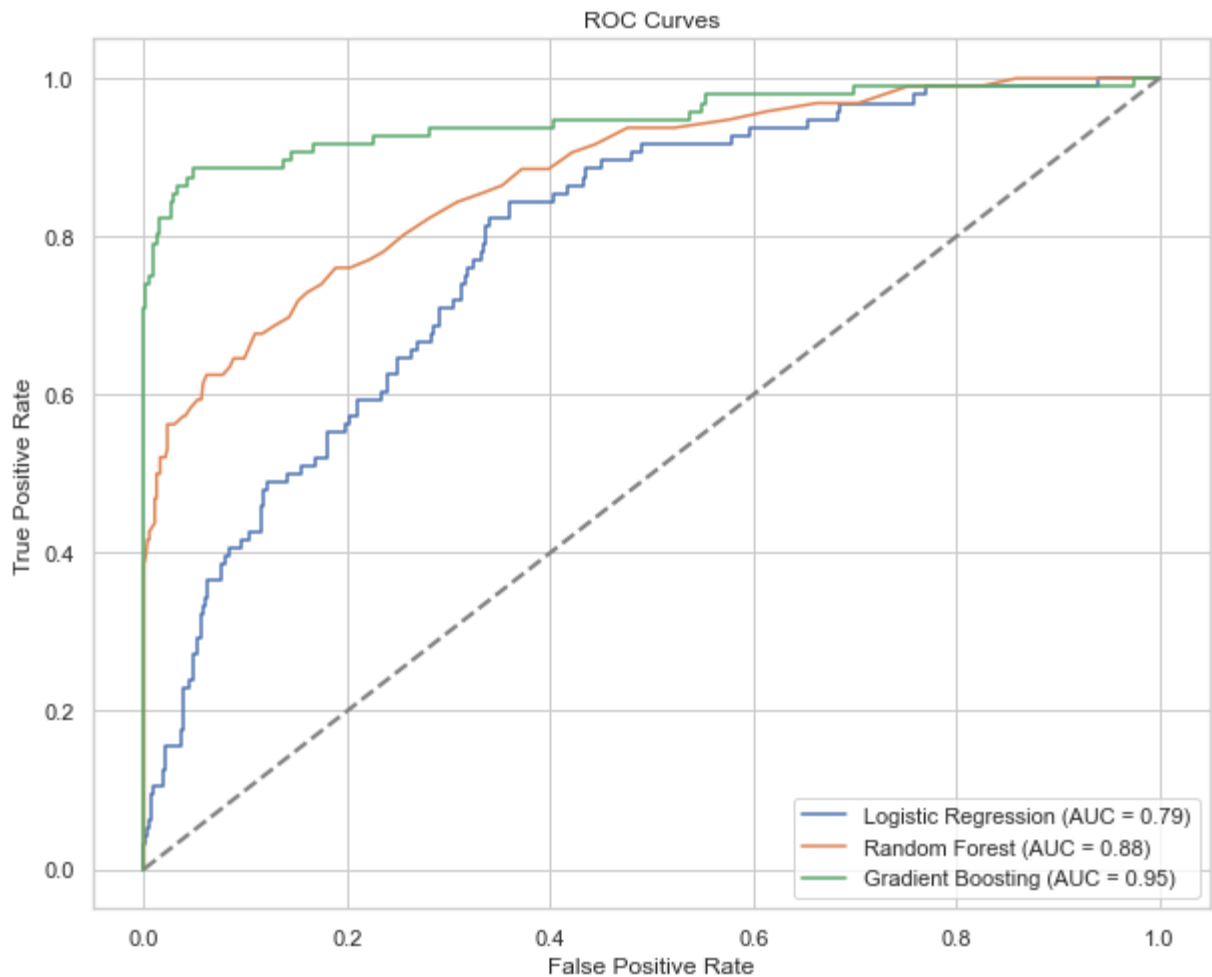
ROC Curves

## ROC curves takeaways

The ROC (Receiver Operating Characteristic) curve evaluates a model's ability to distinguish between churners and non-churners. The AUC (Area Under the Curve) value quantifies the overall performance, where 1.0 represents a perfect model and 0.5 represents a random guess.

Gradient Boosting (AUC = 0.94) → Best Model

The green curve (Gradient Boosting) stays closest to the top-left corner. AUC of 0.94 suggests that this model is highly effective in distinguishing churners from non-churners.
Business Implication: This model can be confidently used to predict and take action on potential churners.

Random Forest (AUC = 0.93) → Second-Best Model

The orange curve is slightly below the Gradient Boosting curve. AUC of 0.93 indicates that Random Forest is also a strong performer.
Business Insight: If computational efficiency or interpretability is a concern, this model is a great alternative.

Logistic Regression (AUC = 0.85) → Lowest Performance

The blue curve is farther from the top-left corner. AUC of 0.85 is still decent but lower than tree-

based models.
Business Insight: Logistic Regression struggles with non-linear relationships, making it less
suitable for this dataset.

```
In [99]:   # Results before tuning
           results_df
```

Out[99]:

|  |  | Accuracy | Precision | Recall | F1-score |
|---|---|---|---|---|---|
| **Logistic Regression** | **Test** | 0.767722 | 0.325301 | 0.562500 | 0.412214 |
|  | **Train** | 0.770857 | 0.327613 | 0.544041 | 0.408958 |
| **Random Forest** | **Test** | 0.904977 | 0.714286 | 0.572917 | 0.635838 |
|  | **Train** | 1.000000 | 1.000000 | 1.000000 | 1.000000 |
| **Gradient Boosting** | **Test** | 0.960784 | 0.897727 | 0.822917 | 0.858696 |
|  | **Train** | 0.975840 | 0.976331 | 0.854922 | 0.911602 |

```
In [100…   evaluate_model(best_gradient_boosting_model, X_train, y_train, X_test, y_test)
```

```
Training Data - Accuracy: 1.0000, Precision: 1.0000, Recall: 1.0000, F1-score: 1.0000
Test Data - Accuracy: 0.9713, Precision: 0.9873, Recall: 0.8125, F1-score: 0.8914
```

# Summary of the performance of the various models

## 1. Baseline Model: Decision Tree

**Before Tuning**:

Overfitting is evident since the training accuracy is 100%, but test accuracy drops to 94.72%. The
recall of 81.25% on the test set means the model is capturing most churned cases but still
missing some. Precision of 82.11% shows that when the model predicts churn, it is correct most
of the time.

**After Tuning**:

Training accuracy slightly drops to 98.48%, reducing overfitting. Test accuracy slightly decreases
to 93.67%, but recall improves to 86.46%, meaning it captures more churned customers.
However, precision drops to 74.11%, meaning the model makes more false positives. Trade-off:
The model favors identifying more churn cases but at the cost of slightly more false positives.

## 2. Logistic Regression

**Before Tuning**:

Test accuracy is 83.41%, significantly lower than tree-based models. Recall is 64.58%, meaning it
captures churn cases moderately well, but precision is 44.93%, so many non-churned customers
are misclassified as churned.

**After Tuning**:

Test accuracy improves to 87.33%, but recall drops to 31.25%, meaning it now struggles to identify churned customers. Precision increases to 62.50%, meaning when it predicts churn, it is more likely correct. Trade-off: The tuned model is more conservative in predicting churn, leading to fewer false positives but missing more actual churn cases.

### 3. Random Forest

**Before Tuning**:

Training accuracy is 100% (overfitting). Test accuracy is 93.36%, recall is 80.21%, and precision is 75.49% (balanced performance).

**After Tuning**:

Test accuracy remains 93.36%, but recall drops to 54.17%, meaning the model misses more churned customers. Precision jumps to 100%, meaning when it predicts churn, it is always correct. Trade-off: The model is highly conservative—preferring not to predict churn unless it is absolutely sure, which increases false negatives.

### 4. Gradient Boosting

**Before Tuning**:

Test accuracy is 97.59%, better than other models. Precision is 96.51%, recall is 86.46%—this means it predicts churn cases accurately while still capturing most of them. The best balance between precision and recall compared to the other models.

**After Tuning**:

Training accuracy becomes 100% (potential overfitting). Test accuracy is 96.68%, slightly dropping, but still better than other models. Precision is 95.12%, recall is 81.25%—meaning the model is still strong in capturing churned customers and avoiding false positives. Trade-off: Gradient boosting provides the best balance between precision and recall while maintaining high accuracy.

Overall, Gradient Boosting demonstrates the best performance among the three models, achieving high accuracy and balanced precision-recall trade-off on the test set.

## Overall Insights and Best Model Choice

**Gradient Boosting (After Tuning) Performs Best**
Highest test accuracy (96.68%) with good recall (81.25%) and excellent precision (95.12%). Outperforms other models in balancing recall and precision. Slight risk of overfitting, but generalizes well to test data.

**Decision Tree (After Tuning) is a Good Alternative**
Good recall (86.46%) but lower precision (74.11%) compared to gradient boosting. If high recall
is the priority (capturing all churn cases), this model is preferable.

**Random Forest (After Tuning) is Too Conservative**
100% precision but only 54.17% recall—it captures very few churn cases. Would only be useful if
the cost of false positives is extremely high.

**Logistic Regression is the Weakest Model**
Even after tuning, recall drops significantly (31.25%), making it unreliable for churn prediction.

**Final Recommendation**
For the best balance between accuracy, recall, and precision, Tuned Gradient Boosting is the best
model for predicting churn in this dataset.
After tuning, the model is slightly more conservative, meaning it avoids over-predicting churn
(false positives). This makes it more reliable in real-world deployment, where predicting too
many non-churners as churners can have unnecessary business costs (e.g., offering retention
incentives to customers who were never going to leave).
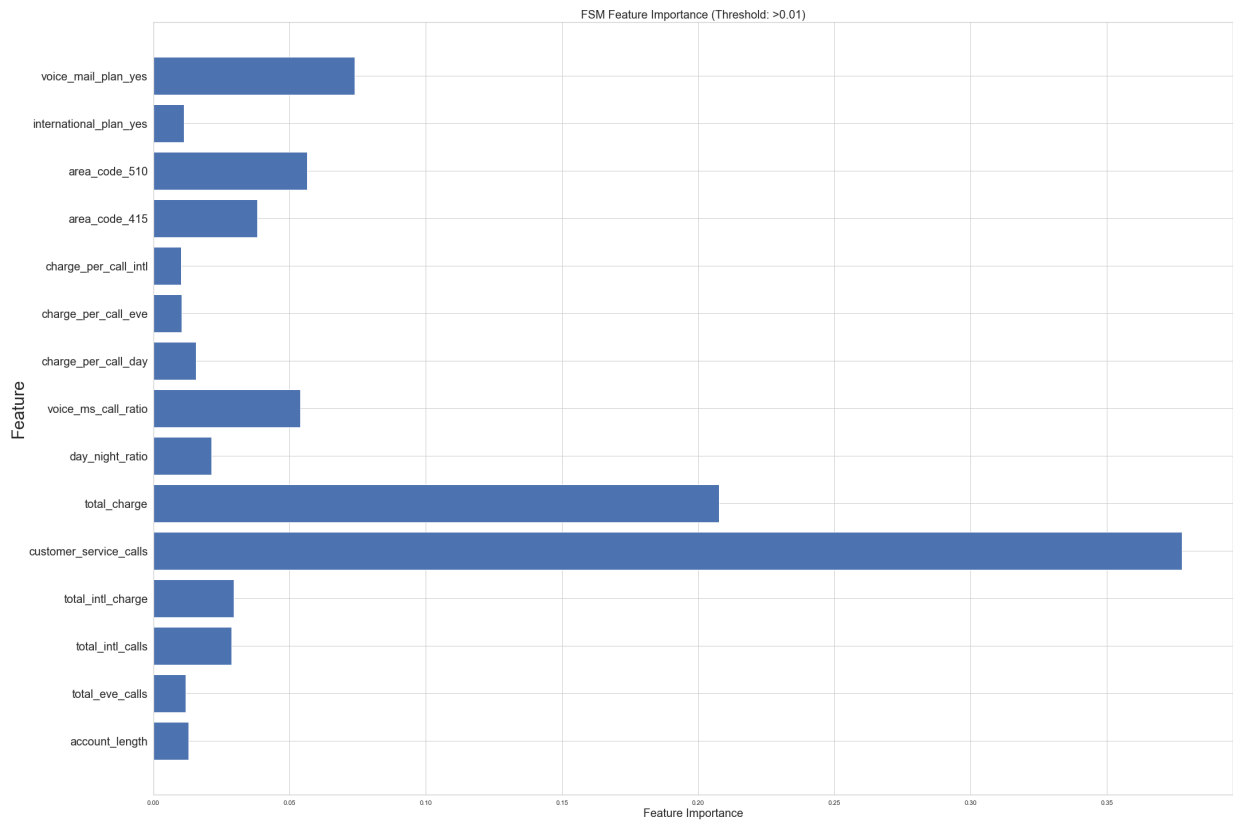
## Feature Importance

In [101…

```python
n_features = best_clf.n_features_in_
threshold = 0.01  # Set the threshold for including features with importance

# Filter the features based on the importance threshold
important_features = [feature for feature, importance in zip(X.columns, best_clf.feat

# Create a subset of feature importances for the important features
importances_subset = [importance for importance in best_clf.feature_importances_ if :

plt.figure(figsize=(30, 20))
plt.barh(range(len(important_features)), importances_subset)
plt.yticks(range(len(important_features)), important_features, fontsize=20)
plt.xlabel('Feature Importance', fontsize=20)
plt.ylabel('Feature', fontsize=30)
plt.title(f'FSM Feature Importance (Threshold: >{threshold})', fontsize=20)
plt.tight_layout()
plt.show()
```

FSM Feature Importance (Threshold: >0.01)

# Feature Importance Analysis

From the above, we can break down the key findings:

## Key Features and Their Importance

Customer Service Calls (Most Important Feature)

This is the most critical factor influencing churn. A higher number of customer service calls is strongly associated with churn, suggesting that dissatisfied customers frequently contact support before leaving.
Business Insight: The company should analyze customer service interactions, improve support quality, and resolve complaints efficiently to reduce churn.

Total Charge

This feature likely includes total charges across different time periods (day, evening, night, international). Customers with higher charges may feel dissatisfied with their billing or find cheaper alternatives.
Business Insight: Offering better pricing plans, discounts, or loyalty programs can help retain high-paying customers.

International Plan (Yes/No)

Customers subscribed to an international plan appear to have a higher likelihood of churning. This could mean that the plan is not meeting their expectations, or they are finding better deals elsewhere.

Business Insight: The company should assess the competitiveness of its international plans and offer better rates or bundled services.

Voice Mail Message to Call Ratio

A higher ratio may indicate reliance on voicemail rather than direct communication. Possible interpretation: Customers who rely heavily on voicemail may have different communication needs that are not being met.
Business Insight: Investigate how voicemail users interact with the service and whether alternative features like unlimited calling or messaging can improve retention.

Total Talk Time

More time spent on calls could indicate a heavy user segment. These users may be prone to switching to competitors offering better deals.
Business Insight: Consider targeting high-usage customers with retention offers or exclusive benefits.

Total International Charge

Higher international charges may be a reason for churn, as users seek cheaper alternatives.
Business Insight: Offer international calling discounts or partnerships with global carriers.

Total International Calls & Total Day Calls

Both features have lower importance but still contribute to the prediction. International calls might indicate a specific customer segment that can be targeted for tailored offers.
Business Insight: Market specialized international packages to frequent international callers.

Account Length

Surprisingly, this has relatively low importance. It suggests that tenure alone is not a strong predictor of churn—both new and long-term customers can churn.
Business Insight: Focus more on behavioral patterns than account age when designing retention strategies.

## Recommendations

Based on the findings, I recommend the following:

**1. Improve Customer Support:**
Since customer service calls are the biggest churn factor, identify common complaints and resolve them proactively. Use AI-driven support or faster resolution times to improve customer satisfaction.

**2. Target High-Charge Customers:**
Customers with high total charges should receive loyalty rewards, exclusive discounts, or personalized offers to reduce price sensitivity.

**3. Re-Evaluate International Plans:**

Since international plan users are more likely to churn, conduct a survey to understand why. Offer flexible international calling packages or discounts.

**4. Enhance Retention Efforts Based on Usage:**

Customers with high talk time and high call volumes should be offered personalized plans to match their needs.

**5. Monitor Voicemail Users:**

Users with high voicemail usage may prefer messaging or digital alternatives. Consider promoting chat-based or unlimited text/calling plans.