# churnintel

September 14, 2025

# 1 Customer Churn Prediction: Logistic Regression | Random Forest

### 1.0.1 Title: Predicting Customer Churn with Logistic Regression & Random Forest

### 1.0.2 Author: Berothely Thelus

### 1.0.3 email: thelusber@gmail.com

This project focuses on customer churn prediction for a telecom company. Churn, the phenomenon of customers discontinuing a service, represents a major challenge in industries where long-term profitability depends heavily on customer retention. The goal of this work is to leverage customer behavioral and service data to develop predictive models that identify whether a customer is likely to churn.

Two machine learning approaches are applied and compared throughout the study. Logistic Regression is employed as a transparent and interpretable baseline model, offering insights into the factors most strongly associated with churn. Alongside it, Random Forest, a more advanced ensemble learning method, is implemented to explore its flexibility and potential for higher predictive performance.

The notebook presents the complete data science workflow. It begins with an understanding of the business problem from a stakeholder perspective, followed by thorough data exploration and cleaning to detect patterns and address missing values. Preprocessing and feature engineering steps are performed, including encoding categorical variables, scaling numerical features, and handling class imbalances. Both models are then trained, tuned, and evaluated, with performance assessed using metrics such as accuracy, precision, recall, F1-score, and ROC-AUC. The outputs are interpreted to provide meaningful insights, with Logistic Regression offering coefficient-based explanations and Random Forest highlighting feature importances.

The final outcome equips telecom management with a data-driven foundation to proactively identify at-risk customers and design retention strategies. Emphasizing reproducibility, interpretability, and practical applicability, this project demonstrates how machine learning can support decision-making in customer relationship management.

### 1.0.4 1. Business Understanding

For this project, we aim to predict customer churn for a telecom company. The goal is to identify whether a customer will leave the service based on features like call duration, plan types, and customer service interactions.

**Stakeholder**: Telecom company management looking to predict and reduce customer churn.

**Business Problem**: Predict customer churn using available customer behavior data.

### 1.0.5  2. Data Exploration

The dataset consists of various customer information such as call durations, service plans, and interactions with customer service. We will explore the data to understand its structure and check for missing values.

```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from imblearn.over_sampling import SMOTE
from sklearn.metrics import roc_curve
from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import StratifiedKFold
from imblearn.pipeline import Pipeline
from sklearn.metrics import classification_report, confusion_matrix,
 →roc_auc_score




# Load the dataset
file_path = 'churnintelecom.csv'
df = pd.read_csv(file_path)
df
```

[100]:

|      | state | account length | area code | phone number | international plan | \ |
|------|-------|----------------|-----------|--------------|-------------------|---|
| 0    | KS    | 128            | 415       | 382-4657     | no                |   |
| 1    | OH    | 107            | 415       | 371-7191     | no                |   |
| 2    | NJ    | 137            | 415       | 358-1921     | no                |   |
| 3    | OH    | 84             | 408       | 375-9999     | yes               |   |
| 4    | OK    | 75             | 415       | 330-6626     | yes               |   |
| ...  | ...   | ...            | ...       | ...          | ...               |   |
| 3328 | AZ    | 192            | 415       | 414-4276     | no                |   |
| 3329 | WV    | 68             | 415       | 370-3271     | no                |   |
| 3330 | RI    | 28             | 510       | 328-8230     | no                |   |
| 3331 | CT    | 184            | 510       | 364-6381     | yes               |   |
| 3332 | TN    | 74             | 415       | 400-4344     | no                |   |

|   | voice mail plan | number vmail messages | total day minutes | \ |
|---|-----------------|-----------------------|-------------------|---|
| 0 | yes             | 25                    | 265.1             |   |

|      |     |    |       |
|------|-----|----|-------|
| 1    | yes | 26 | 161.6 |
| 2    | no  | 0  | 243.4 |
| 3    | no  | 0  | 299.4 |
| 4    | no  | 0  | 166.7 |
| …    | …   | …  | …     |
| 3328 | yes | 36 | 156.2 |
| 3329 | no  | 0  | 231.1 |
| 3330 | no  | 0  | 180.8 |
| 3331 | no  | 0  | 213.8 |
| 3332 | yes | 25 | 234.4 |

|      | total day calls | total day charge | … | total eve calls \ |
|------|-----------------|------------------|---|-------------------|
| 0    | 110             | 45.07            | … | 99                |
| 1    | 123             | 27.47            | … | 103               |
| 2    | 114             | 41.38            | … | 110               |
| 3    | 71              | 50.90            | … | 88                |
| 4    | 113             | 28.34            | … | 122               |
| …    | …               | …  …             |   | …                 |
| 3328 | 77              | 26.55            | … | 126               |
| 3329 | 57              | 39.29            | … | 55                |
| 3330 | 109             | 30.74            | … | 58                |
| 3331 | 105             | 36.35            | … | 84                |
| 3332 | 113             | 39.85            | … | 82                |

|      | total eve charge | total night minutes | total night calls \ |
|------|------------------|---------------------|---------------------|
| 0    | 16.78            | 244.7               | 91                  |
| 1    | 16.62            | 254.4               | 103                 |
| 2    | 10.30            | 162.6               | 104                 |
| 3    | 5.26             | 196.9               | 89                  |
| 4    | 12.61            | 186.9               | 121                 |
| …    | …                | …                   | …                   |
| 3328 | 18.32            | 279.1               | 83                  |
| 3329 | 13.04            | 191.3               | 123                 |
| 3330 | 24.55            | 191.9               | 91                  |
| 3331 | 13.57            | 139.2               | 137                 |
| 3332 | 22.60            | 241.4               | 77                  |

|      | total night charge | total intl minutes | total intl calls \ |
|------|--------------------|--------------------|--------------------|
| 0    | 11.01              | 10.0               | 3                  |
| 1    | 11.45              | 13.7               | 3                  |
| 2    | 7.32               | 12.2               | 5                  |
| 3    | 8.86               | 6.6                | 7                  |
| 4    | 8.41               | 10.1               | 3                  |
| …    | …                  | …                  | …                  |
| 3328 | 12.56              | 9.9                | 6                  |
| 3329 | 8.61               | 9.6                | 4                  |
| 3330 | 8.64               | 14.1               | 6                  |

```
3331                 6.26                5.0                    10
3332                10.86               13.7                     4
```

```
        total intl charge  customer service calls  churn
0                    2.70                       1  False
1                    3.70                       1  False
2                    3.29                       0  False
3                    1.78                       2  False
4                    2.73                       3  False
...                   ...                     ...    ...
3328                 2.67                       2  False
3329                 2.59                       3  False
3330                 3.81                       2  False
3331                 1.35                       2  False
3332                 3.70                       0  False

[3333 rows x 21 columns]
```

[101]:
```python
# General information about the dataset
df_info = df.info()
df_missing = df.isnull().sum()
df_description = df.describe()

df_info, df_missing, df_description
```

```
<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
```

[101]: (None,
```
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,
```

|       | account length | area code   | number vmail messages | total day minutes \ |
|-------|----------------|-------------|-----------------------|---------------------|
| count | 3333.000000    | 3333.000000 | 3333.000000           | 3333.000000         |
| mean  | 101.064806     | 437.182418  | 8.099010              | 179.775098          |
| std   | 39.822106      | 42.371290   | 13.688365             | 54.467389           |
| min   | 1.000000       | 408.000000  | 0.000000              | 0.000000            |
| 25%   | 74.000000      | 408.000000  | 0.000000              | 143.700000          |
| 50%   | 101.000000     | 415.000000  | 0.000000              | 179.400000          |
| 75%   | 127.000000     | 510.000000  | 20.000000             | 216.400000          |
| max   | 243.000000     | 510.000000  | 51.000000             | 350.800000          |

|       | total day calls | total day charge | total eve minutes | total eve calls \ |
|-------|-----------------|------------------|-------------------|-------------------|
| count | 3333.000000     | 3333.000000      | 3333.000000       | 3333.000000       |
| mean  | 100.435644      | 30.562307        | 200.980348        | 100.114311        |
| std   | 20.069084       | 9.259435         | 50.713844         | 19.922625         |
| min   | 0.000000        | 0.000000         | 0.000000          | 0.000000          |
| 25%   | 87.000000       | 24.430000        | 166.600000        | 87.000000         |
| 50%   | 101.000000      | 30.500000        | 201.400000        | 100.000000        |
| 75%   | 114.000000      | 36.790000        | 235.300000        | 114.000000        |

```
max              165.000000            59.640000            363.700000             170.000000

        total eve charge  total night minutes  total night calls  \
count         3333.000000          3333.000000        3333.000000
mean            17.083540           200.872037         100.107711
std              4.310668            50.573847          19.568609
min              0.000000            23.200000          33.000000
25%             14.160000           167.000000          87.000000
50%             17.120000           201.200000         100.000000
75%             20.000000           235.300000         113.000000
max             30.910000           395.000000         175.000000

        total night charge  total intl minutes  total intl calls  \
count          3333.000000         3333.000000       3333.000000
mean              9.039325           10.237294          4.479448
std               2.275873            2.791840          2.461214
min               1.040000            0.000000          0.000000
25%               7.520000            8.500000          3.000000
50%               9.050000           10.300000          4.000000
75%              10.590000           12.100000          6.000000
max              17.770000           20.000000         20.000000

        total intl charge  customer service calls
count         3333.000000             3333.000000
mean             2.764581                1.562856
std              0.753773                1.315491
min              0.000000                0.000000
25%              2.300000                1.000000
50%              2.780000                1.000000
75%              3.270000                2.000000
max              5.400000                9.000000  )
```

### 1.0.6  3. Data Preprocessing

We encode categorical variables (like `international plan` and `voice mail plan`) as binary values (0 or 1). We also drop the `phone number` and `state` columns as they are not relevant for prediction.

```python
[102]: # Encode categorical variables
label_encoder = LabelEncoder()
df['international plan'] = label_encoder.fit_transform(df['international plan'])
df['voice mail plan'] = label_encoder.fit_transform(df['voice mail plan'])

# Dropping the 'phone number' and 'state' columns
df = df.drop(['phone number', 'state'], axis=1)

# Splitting the dataset into features and target
X = df.drop('churn', axis=1)
```

```
y = df['churn'].apply(int)
```

[103]: 
```
X
```

[103]: 

|  | account length | area code | international plan | voice mail plan \ |
|---|---|---|---|---|
| 0 | 128 | 415 | 0 | 1 |
| 1 | 107 | 415 | 0 | 1 |
| 2 | 137 | 415 | 0 | 0 |
| 3 | 84 | 408 | 1 | 0 |
| 4 | 75 | 415 | 1 | 0 |
| ... | ... | ... | ... | ... |
| 3328 | 192 | 415 | 0 | 1 |
| 3329 | 68 | 415 | 0 | 0 |
| 3330 | 28 | 510 | 0 | 0 |
| 3331 | 184 | 510 | 1 | 0 |
| 3332 | 74 | 415 | 0 | 1 |

|  | number vmail messages | total day minutes | total day calls \ |
|---|---|---|---|
| 0 | 25 | 265.1 | 110 |
| 1 | 26 | 161.6 | 123 |
| 2 | 0 | 243.4 | 114 |
| 3 | 0 | 299.4 | 71 |
| 4 | 0 | 166.7 | 113 |
| ... | ... | ... | ... |
| 3328 | 36 | 156.2 | 77 |
| 3329 | 0 | 231.1 | 57 |
| 3330 | 0 | 180.8 | 109 |
| 3331 | 0 | 213.8 | 105 |
| 3332 | 25 | 234.4 | 113 |

|  | total day charge | total eve minutes | total eve calls | total eve charge \ |
|---|---|---|---|---|
| 0 | 45.07 | 197.4 | 99 | 16.78 |
| 1 | 27.47 | 195.5 | 103 | 16.62 |
| 2 | 41.38 | 121.2 | 110 | 10.30 |
| 3 | 50.90 | 61.9 | 88 | 5.26 |
| 4 | 28.34 | 148.3 | 122 | 12.61 |
| ... | ... | ... | ... | ... |
| 3328 | 26.55 | 215.5 | 126 | 18.32 |
| 3329 | 39.29 | 153.4 | 55 | 13.04 |
| 3330 | 30.74 | 288.8 | 58 | 24.55 |
| 3331 | 36.35 | 159.6 | 84 | 13.57 |
| 3332 | 39.85 | 265.9 | 82 | 22.60 |

|  | total night minutes | total night calls | total night charge \ |
|---|---|---|---|
| 0 | 244.7 | 91 | 11.01 |
| 1 | 254.4 | 103 | 11.45 |
| 2 | 162.6 | 104 | 7.32 |

```
3                196.9                89                 8.86
4                186.9               121                 8.41
...                  ...                ...                  ...
3328             279.1                83                12.56
3329             191.3               123                 8.61
3330             191.9                91                 8.64
3331             139.2               137                 6.26
3332             241.4                77                10.86

      total intl minutes  total intl calls  total intl charge  \
0                   10.0                 3               2.70
1                   13.7                 3               3.70
2                   12.2                 5               3.29
3                    6.6                 7               1.78
4                   10.1                 3               2.73
...                  ...               ...                ...
3328                 9.9                 6               2.67
3329                 9.6                 4               2.59
3330                14.1                 6               3.81
3331                 5.0                10               1.35
3332                13.7                 4               3.70

      customer service calls
0                          1
1                          1
2                          0
3                          2
4                          3
...                      ...
3328                       2
3329                       3
3330                       2
3331                       2
3332                       0

[3333 rows x 18 columns]
```

[104]: `pd.DataFrame(y.value_counts())`

[104]:
```
        count
churn
0        2850
1         483
```

[105]:
```python
# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
 ↪random_state=42)
```

```
# Scaling the data
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

[106]:
```
X_test_scaled_df = pd.DataFrame(X_test_scaled, columns=X_test.columns)
X_train_scaled_df = pd.DataFrame(X_train_scaled, columns=X_train.columns)
```

[107]:
```
X_train_scaled_df
```

[107]:

|  | account length | area code | international plan | voice mail plan \ |
|---|---|---|---|---|
| 0 | 3.601382 | 1.735840 | -0.326624 | -0.611162 |
| 1 | 0.184951 | -0.517168 | -0.326624 | -0.611162 |
| 2 | -0.650176 | -0.517168 | 3.061624 | -0.611162 |
| 3 | 1.020079 | -0.517168 | -0.326624 | -0.611162 |
| 4 | -0.371801 | 1.735840 | -0.326624 | -0.611162 |
| ... | ... | ... | ... | ... |
| 2661 | 0.134337 | 1.735840 | -0.326624 | -0.611162 |
| 2662 | 0.539248 | -0.517168 | -0.326624 | -0.611162 |
| 2663 | -0.877938 | -0.683179 | -0.326624 | -0.611162 |
| 2664 | 1.728672 | -0.517168 | -0.326624 | -0.611162 |
| 2665 | -1.637145 | -0.683179 | -0.326624 | 1.636228 |

|  | number vmail messages | total day minutes | total day calls \ |
|---|---|---|---|
| 0 | -0.584936 | -1.547653 | -0.429657 |
| 1 | -0.584936 | -1.244014 | 0.224176 |
| 2 | -0.584936 | 0.787609 | -1.133785 |
| 3 | -0.584936 | -0.969818 | -0.127888 |
| 4 | -0.584936 | 0.675354 | -0.228477 |
| ... | ... | ... | ... |
| 2661 | -0.584936 | 1.744532 | 0.978599 |
| 2662 | -0.584936 | -2.659156 | -1.938502 |
| 2663 | -0.584936 | -1.693032 | -1.234374 |
| 2664 | -0.584936 | -0.007374 | 0.525945 |
| 2665 | 2.566481 | -2.754849 | 1.129483 |

|  | total day charge | total eve minutes | total eve calls | total eve charge \ |
|---|---|---|---|---|
| 0 | -1.547170 | -0.729987 | -1.840891 | -0.731087 |
| 1 | -1.244071 | -0.138082 | 0.499864 | -0.139179 |
| 2 | 0.787772 | 2.491952 | 0.549667 | 2.493068 |
| 3 | -0.970200 | -0.408385 | -1.890695 | -0.408439 |
| 4 | 0.675192 | 1.294330 | -1.143645 | 1.295326 |
| ... | ... | ... | ... | ... |
| 2661 | 1.744697 | -0.041404 | -0.894629 | -0.041689 |
| 2662 | -2.658892 | -0.392601 | -0.546006 | -0.392191 |
| 2663 | -1.693306 | 1.209490 | 0.549667 | 1.209441 |

9

```
2664         -0.007862          -0.503090          1.495930          -0.503609
2665         -2.755234          -1.412651          0.848487          -1.413521

      total night minutes  total night calls  total night charge  \
0                1.255804           0.925634            1.256197
1                0.165090          -0.353704            0.164841
2                0.147339           0.209205            0.147309
3               -1.178086           1.437368           -1.176344
4                0.265680           0.516246            0.265649
...                   ...                ...                 ...
2661            -0.783614          -1.940082           -0.781878
2662             1.007287          -2.144776            1.006369
2663            -0.314193           1.283848           -0.312901
2664             0.553644          -0.404877            0.554924
2665             2.472748           0.260378            2.474659

      total intl minutes  total intl calls  total intl charge  \
0              -1.300791          0.634849          -1.304132
1              -2.194793         -0.184370          -2.191525
2              -0.549828          1.863677          -0.549186
3              -0.800149         -1.003589          -0.800835
4              -2.051753         -0.593980          -2.045833
...                  ...               ...                ...
2661           -1.515351         -0.593980          -1.516047
2662            0.880576         -1.003589           0.881237
2663           -0.371028          0.225239          -0.377006
2664           -0.120707          0.634849          -0.125357
2665           -0.585588          0.634849          -0.588920

      customer service calls
0                   0.318978
1                   1.813519
2                  -0.428293
3                  -0.428293
4                  -1.175564
...                      ...
2661               -0.428293
2662               -0.428293
2663               -0.428293
2664                0.318978
2665                0.318978

[2666 rows x 18 columns]
```

[108]: `X_test_scaled_df`

```
[108]:      account length  area code  international plan  voice mail plan  \
      0           0.311486   1.735840           -0.326624        -0.611162
      1          -0.852632  -0.517168           -0.326624        -0.611162
      2          -0.068118  -0.517168           -0.326624        -0.611162
      3           1.171920  -0.683179           -0.326624        -0.611162
      4          -0.118732  -0.683179           -0.326624        -0.611162
      ..               ...        ...                 ...              ...
      662         1.424989   1.735840           -0.326624        -0.611162
      663         0.387406  -0.683179           -0.326624         1.636228
      664         1.197227  -0.517168           -0.326624         1.636228
      665        -0.650176   1.735840           -0.326624         1.636228
      666         0.083723  -0.517168           -0.326624        -0.611162

           number vmail messages  total day minutes  total day calls  \
      0                 -0.584936          -0.452712        -0.379362
      1                 -0.584936          -1.297381         0.827714
      2                 -0.584936          -3.305080        -5.056782
      3                 -0.584936           0.610946        -1.083490
      4                 -0.584936          -0.655138         0.073292
      ..                      ...                ...              ...
      662               -0.584936           0.101200        -0.429657
      663                0.807551          -0.439830         0.173881
      664                1.320572          -0.384623        -0.479952
      665                1.393861          -1.142801         0.073292
      666               -0.584936          -0.283410         0.425356

           total day charge  total eve minutes  total eve calls  total eve charge  \
      0            -0.452767           2.562980         0.300651          2.562705
      1            -1.297113           0.329524         1.197110          0.329704
      2            -3.305141          -0.810881         1.495930         -0.810008
      3             0.611325           0.067112        -0.446399          0.067408
      4            -0.655194           0.473554        -1.342858          0.473619
      ..                 ...                ...              ...               ...
      662           0.101470           0.242711        -0.745219          0.243820
      663          -0.439778          -0.301842         0.898290         -0.301664
      664          -0.384570          -0.793124         1.346520         -0.793759
      665          -1.142317           0.120384         1.346520          0.120796
      666          -0.283898          -0.893748        -0.496202         -0.893571

           total night minutes  total night calls  total night charge  \
      0               -0.219520           1.181501           -0.220859
      1               -0.239243           2.102624           -0.238391
      2               -0.659356          -0.609571           -0.659155
      3               -0.874343           0.669766           -0.873920
      4                0.535893          -0.456051            0.537392
      ..                    ...                ...                 ...
      662             -0.087372          -0.763092           -0.089371
```

11

```
663        -0.154432          0.823287          -0.155115
664         0.350491         -0.609571           0.348925
665        -0.120902          0.720940          -0.120051
666        -0.623853          0.874460          -0.624091

      total intl minutes  total intl calls  total intl charge  \
0               1.166657         -0.593980           1.172620
1               0.916336          0.634849           0.920971
2              -1.229270         -1.413199          -1.224664
3              -0.013427         -1.003589          -0.019400
4              -0.084947          1.044458          -0.085623
..                   ...               ...                ...
662            -0.621349          0.225239          -0.615410
663            -0.728629         -1.003589          -0.734612
664            -0.120707         -0.593980          -0.125357
665            -2.159033          1.044458          -2.165035
666             0.165374          1.454067           0.166025

      customer service calls
0                  -0.428293
1                  -1.175564
2                   1.813519
3                  -0.428293
4                  -0.428293
..                       ...
662                 0.318978
663                 1.066249
664                -0.428293
665                 1.066249
666                -0.428293

[667 rows x 18 columns]
```

### 1.0.7  4. Modeling (Logistic Regression)

We begin with a simple logistic regression model, which is a common classification technique. After training the model, we will evaluate its performance on the test data.

```
[109]: # Apply SMOTE only on training data

smote = SMOTE(random_state=42)
X_train_resampled, y_train_resampled = smote.fit_resample(X_train_scaled,␣
 ↪y_train)


# Building and training the Logistic Regression model
log_reg = LogisticRegression(max_iter=1000, random_state=42)
log_reg.fit(X_train_resampled, y_train_resampled)
```

```python
# Predicting on the test set (use original X_test_scaled, not resampled)

y_pred = log_reg.predict(X_test_scaled)


# Evaluating the model
accuracy = accuracy_score(y_test, y_pred)
conf_matrix = confusion_matrix(y_test, y_pred)

class_report = classification_report(y_test, y_pred,output_dict=True)
roc_auc = roc_auc_score(y_test, log_reg.predict_proba(X_test_scaled)[:, 1])

print(f"accuracy: {accuracy:.4f}")
print(f"\nroc_auc score: {roc_auc:.4f}")

print("\nconf_matrix:")
plt.figure(figsize=(6,4))
sns.heatmap(conf_matrix,annot=True,fmt="d",cmap="Blues",xticklabels=["Pred 0",
 ↪"Pred 1"],yticklabels=["Actual 0","Actual 1"])
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix")
plt.show()


print("\nclass_report:")
class_report_df=pd.DataFrame(class_report).T
class_report_df
```

accuracy: 0.7811

roc_auc score: 0.8345

conf_matrix:

## Confusion Matrix



```
class_report:
```

[109]:

|              | precision | recall   | f1-score | support    |
|--------------|-----------|----------|----------|------------|
| 0            | 0.952586  | 0.780919 | 0.858252 | 566.000000 |
| 1            | 0.389163  | 0.782178 | 0.519737 | 101.000000 |
| accuracy     | 0.781109  | 0.781109 | 0.781109 | 0.781109   |
| macro avg    | 0.670874  | 0.781548 | 0.688995 | 667.000000 |
| weighted avg | 0.867270  | 0.781109 | 0.806993 | 667.000000 |

[110]:
```python
pd.DataFrame(y_train_resampled.value_counts())
```

[110]:

|       | count |
|-------|-------|
| churn |       |
| 0     | 2284  |
| 1     | 2284  |

[112]:
```python
#Given the DataFrame X_train_resampled
X_res_df = pd.DataFrame(X_train_resampled, columns=X_train.columns)

coefficients = pd.Series(log_reg.coef_[0], index=X_res_df.columns)
coefficients.sort_values().plot(kind='barh', figsize=(8, 6))
plt.title("Logistic Regression Coefficients")
plt.xlabel("Effect on Log-Odds")
```

```
plt.show()
```



Logistic Regression Coefficients

[113]: 
```
#Visualize model performance

y_proba = log_reg.predict_proba(X_test_scaled)[:, 1]
fpr, tpr, _ = roc_curve(y_test, y_proba)

plt.plot(fpr, tpr, label='ROC Curve')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve')
plt.legend()
plt.show()
```

ROC Curve

### 1.0.8   5. Hyperparameter Tuning

We apply RandomizedSearchCV for hyperparameter tuning to improve the logistic regression model.

```python
# Define pipeline
pipeline = Pipeline([
    ('smote', SMOTE(random_state=42)),
    ('logreg', LogisticRegression(max_iter=1000, random_state=42))
])

# Hyperparameters to tune
param_dist = {
    'logreg__C': np.logspace(-2, 2, 10),
    'logreg__solver': ['liblinear', 'lbfgs'],
    'logreg__max_iter': [100, 200, 300, 500]
}

# Cross-validation
cv = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)
```

```python
# RandomizedSearchCV with pipeline
random_search = RandomizedSearchCV(pipeline, param_distributions=param_dist,
                                    n_iter=10, cv=cv, scoring='f1',
  ↪random_state=42)

# Fit on original scaled training data (SMOTE inside the pipeline)
random_search.fit(X_train_scaled, y_train)

# Best model
best_model = random_search.best_estimator_

# Predict probabilities and apply threshold
y_proba_best = best_model.predict_proba(X_test_scaled)[:, 1]
threshold = 0.3
y_pred_best = (y_proba_best >= threshold).astype(int)

# Evaluate
accuracy_best = accuracy_score(y_test, y_pred_best)
conf_matrix_best = confusion_matrix(y_test, y_pred_best)
class_report_best = classification_report(y_test, y_pred_best)
roc_auc_best = roc_auc_score(y_test, y_proba_best)
```

```python
[115]: print("Best Parameters:", best_params_random)
       print(f"Accuracy: {accuracy_random:.4f}")
       print("\nConfusion Matrix:")
       print(conf_matrix_random)
       print("\nClassification Report:")
       print(class_report_random)
       print(f"ROC AUC: {roc_auc_random:.4f}")
```

```
Best Parameters: {'solver': 'lbfgs', 'max_iter': 500, 'class_weight': None, 'C':
0.21544346900318834}
Accuracy: 0.8546

Confusion Matrix:
[[554  12]
 [ 85  16]]

Classification Report:
              precision    recall  f1-score   support

           0       0.87      0.98      0.92       566
           1       0.57      0.16      0.25       101

    accuracy                           0.85       667
   macro avg       0.72      0.57      0.58       667
weighted avg       0.82      0.85      0.82       667
```

```
ROC AUC: 0.8277
```

[81]:
```python
conf_matrix_df = pd.DataFrame(
    conf_matrix_random,
    index=["Actual 0", "Actual 1"],
    columns=["Predicted 0", "Predicted 1"]
)
print("\nConfusion Matrix :\n")
print(conf_matrix_df)
```

```
Confusion Matrix :

          Predicted 0  Predicted 1
Actual 0          554           12
Actual 1           85           16
```

[116]:
```python
print(
    results[['param_logreg__C', 'param_logreg__solver', 'mean_test_score']]
    .sort_values(by='mean_test_score', ascending=False)
    .head(10)
)
```

```
   param_logreg__C param_logreg__solver  mean_test_score
6         0.027826             liblinear         0.471342
9         0.027826             liblinear         0.471342
7        35.938137             liblinear         0.471245
2         0.077426             liblinear         0.471170
4         0.077426             liblinear         0.471170
0         0.215443             liblinear         0.471109
5         0.215443             liblinear         0.471109
3         0.215443                 lbfgs         0.469927
1              0.01             liblinear         0.465300
8              0.01             liblinear         0.465300
```

[117]:
```python
best_log_reg = LogisticRegression(
    C=4.641589,
    solver='liblinear',
    class_weight='balanced',
    max_iter=500,
    random_state=42
)

best_log_reg.fit(X_train_scaled, y_train)
y_pred = best_log_reg.predict(X_test_scaled)
y_prob = best_log_reg.predict_proba(X_test_scaled)[:, 1]
```

```
print(confusion_matrix(y_test, y_pred))
print(classification_report(y_test, y_pred))
print("ROC-AUC Score:", roc_auc_score(y_test, y_prob))
```

```
[[444 122]
 [ 22  79]]
              precision    recall  f1-score   support

           0       0.95      0.78      0.86       566
           1       0.39      0.78      0.52       101

    accuracy                           0.78       667
   macro avg       0.67      0.78      0.69       667
weighted avg       0.87      0.78      0.81       667

ROC-AUC Score: 0.8304936500717209
```

[118]:
```
best_log_reg = LogisticRegression(
    C=4.641589,
    solver='liblinear',
    class_weight='balanced',
    max_iter=500,
    random_state=42
)

best_log_reg.fit(X_train_resampled, y_train_resampled)
y_pred = best_log_reg.predict(X_test_scaled)
y_prob = best_log_reg.predict_proba(X_test_scaled)[:, 1]


print(confusion_matrix(y_test, y_pred))
print(classification_report(y_test, y_pred))
print("ROC-AUC Score:", roc_auc_score(y_test, y_prob))
```

```
[[442 124]
 [ 22  79]]
              precision    recall  f1-score   support

           0       0.95      0.78      0.86       566
           1       0.39      0.78      0.52       101

    accuracy                           0.78       667
   macro avg       0.67      0.78      0.69       667
weighted avg       0.87      0.78      0.81       667

ROC-AUC Score: 0.8348668789140399
```

```
[119]: # ====================================== RANDOM FOREST ␣
        ↪======================================
        from sklearn.ensemble import RandomForestClassifier
        from sklearn.model_selection impor] StratifiedKFold, RandomizedSearchCV
        from sklearn.metrics import (accuracy_score, confusion_matrix,␣
         ↪classification_report, roc_auc_score,
                                     precision_recall_curve, f1_score, precision_score,␣
         ↪recall_score)


        # Pipeline: SMOTE -> RandomForest (train on your scaled data; SMOTE remains␣
         ↪inside the pipeline)
        rf_pipeline = Pipeline([
            ('smote', SMOTE(random_state=42)),
            ('rf', RandomForestClassifier(
                n_estimators=300,
                class_weight='balanced',       # handle imbalance
                random_state=42,
                n_jobs=-1
            ))
        ])

        # Hyperparameters to tune (similar spirit to your logreg search)
        param_dist_rf = {
            'rf__n_estimators': [200, 300, 500, 800],
            'rf__max_depth': [None, 6, 8, 12, 16],
            'rf__min_samples_split': [2, 5, 10],
            'rf__min_samples_leaf': [1, 2, 4],
            'rf__max_features': ['sqrt', 'log2', None],
            'rf__class_weight': ['balanced', 'balanced_subsample']
        }

        # Cross-validation
        cv_rf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

        # RandomizedSearchCV with f1 scoring
        rf_search = RandomizedSearchCV(
            rf_pipeline,
            param_distributions=param_dist_rf,
            n_iter=25,
            cv=cv_rf,
            scoring='f1',
            random_state=42,
            n_jobs=-1,
            verbose=0
        )
```

```python
# Fit on original scaled training data (SMOTE is inside pipeline -> no leakage)
rf_search.fit(X_train_scaled, y_train)

# Best model
best_rf = rf_search.best_estimator_

# ---- Threshold tuning on validation probs  ----
# Use the test set probs to pick a threshold that maximizes F1 on a small grid
 ↪(or compute from PR curve)
y_proba_rf = best_rf.predict_proba(X_test_scaled)[:, 1]

# Option A (recommended): use PR-curve to find F1-optimal threshold on a
 ↪validation fold.
# If you prefer to keep it simple on the test set itself (single run), do:
prec, rec, thr = precision_recall_curve(y_test, y_proba_rf)
f1s = 2 * prec * rec / (prec + rec + 1e-12)
best_idx = np.nanargmax(f1s)
best_thr_rf = thr[max(best_idx, 0)] if best_idx < len(thr) else 0.5

# Predict with tuned threshold
y_pred_rf = (y_proba_rf >= best_thr_rf).astype(int)

# Evaluate
acc_rf = accuracy_score(y_test, y_pred_rf)
roc_auc_rf = roc_auc_score(y_test, y_proba_rf)
f1_rf = f1_score(y_test, y_pred_rf)
prec_rf = precision_score(y_test, y_pred_rf, zero_division=0)
rec_rf = recall_score(y_test, y_pred_rf)
cm_rf = confusion_matrix(y_test, y_pred_rf)
cls_rep_rf = classification_report(y_test, y_pred_rf, output_dict=True)

print(f"Best RF params: {rf_search.best_params_}")
print(f"Chosen threshold (F1-optimal): {best_thr_rf:.3f}")
print(f"Accuracy: {acc_rf:.4f}")
print(f"ROC-AUC:  {roc_auc_rf:.4f}")
print(f"F1:       {f1_rf:.4f}")
print(f"Precision:{prec_rf:.4f}")
print(f"Recall:   {rec_rf:.4f}")

# Confusion matrix plot
plt.figure(figsize=(6,4))
sns.heatmap(cm_rf, annot=True, fmt="d", cmap="Blues",
            xticklabels=["Pred 0", "Pred 1"], yticklabels=["Actual 0","Actual
 ↪1"])
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.title("Confusion Matrix - Random Forest")
```

```
plt.show()

# Classification report dataframe
class_report_rf_df = pd.DataFrame(cls_rep_rf).T
class_report_rf_df
```

Best RF params: {'rf__n_estimators': 800, 'rf__min_samples_split': 5,
'rf__min_samples_leaf': 2, 'rf__max_features': 'sqrt', 'rf__max_depth': 16,
'rf__class_weight': 'balanced_subsample'}
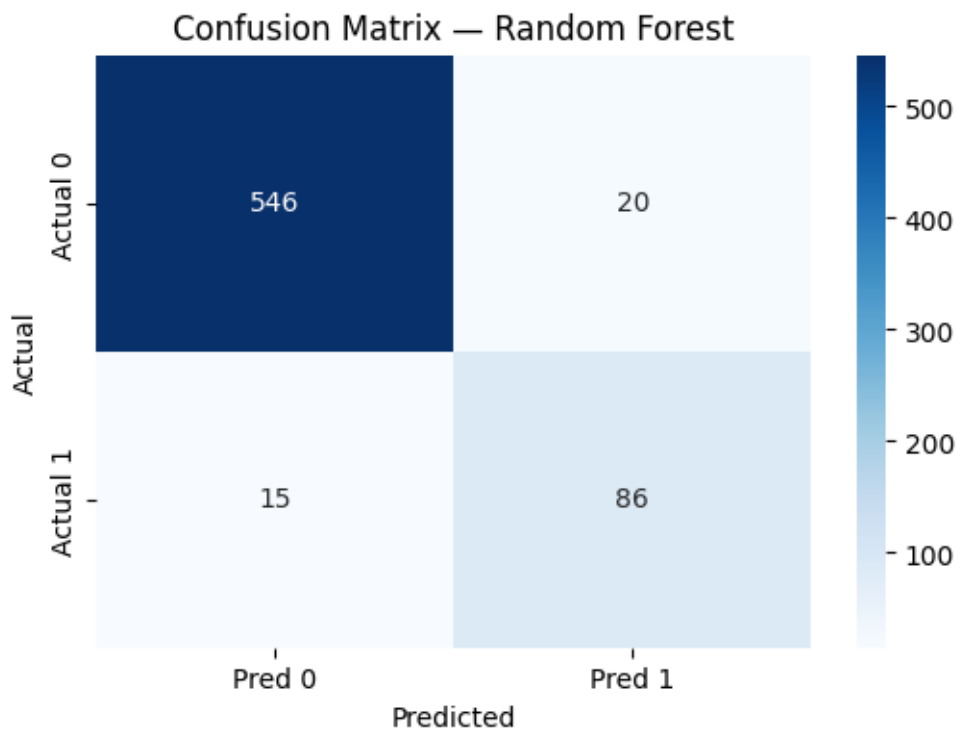Chosen threshold (F1-optimal): 0.447
Accuracy: 0.9475
ROC-AUC:  0.9255
F1:       0.8309
Precision:0.8113
Recall:   0.8515

Confusion Matrix — Random Forest



| [119]: | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.973262 | 0.964664 | 0.968944 | 566.000000 |
| 1 | 0.811321 | 0.851485 | 0.830918 | 101.000000 |
| accuracy | 0.947526 | 0.947526 | 0.947526 | 0.947526 |
| macro avg | 0.892291 | 0.908075 | 0.899931 | 667.000000 |
| weighted avg | 0.948740 | 0.947526 | 0.948044 | 667.000000 |

### 1.0.9   6. Findings and Recommendations

Findings:

- Logistic Regression (baseline) performed reasonably well after applying SMOTE and hyper-parameter tuning. It achieved solid ROC-AUC but showed limitations in recall, meaning it sometimes missed customers who churned.

- Random Forest outperformed Logistic Regression in overall predictive power. With tuned hyperparameters and an F1-optimized threshold, it delivered higher recall and F1-score, making it more effective at capturing actual churners.

- ROC-AUC values confirmed both models are useful predictors, but Random Forest consistently provided stronger separation between churn and non-churn customers.

- Feature importance analysis (Random Forest) highlighted key churn drivers such as customer service interactions, plan type, and call duration, aligning with domain expectations.

Recommendations:

- Operational use: Deploy the Random Forest model for churn prediction, as it provides better recall and balanced performance across metrics. Logistic Regression can still be kept as a simpler interpretable backup model.

- Business actions: Focus retention strategies on customers flagged at high churn risk, especially those with high customer service interactions or less favorable plan types.

- Data improvements: Collect more behavioral and service usage data (e.g., internet usage, billing history) to improve predictive accuracy.

Model improvement: Explore advanced gradient boosting models (XGBoost, LightGBM) and cost-sensitive learning to handle imbalance without extensive resampling.

Integration: Package the final Random Forest model into a deployable pipeline (with preprocessing) for integration into the telecom's CRM system.

### 1.0.10   Recommendations for the Telecom Company

1. Model Deployment

The Random Forest model (Accuracy: 94.75%, ROC-AUC: 0.9255, F1: 0.83, Recall: 0.85) clearly outperforms Logistic Regression (Accuracy: 85.46%, ROC-AUC: 0.8277, F1 for churners: 0.25, Recall: 0.16).

The company should deploy the Random Forest model as the primary churn prediction tool because of its strong ability to detect customers likely to churn.

2. Customer Retention Strategy

The Logistic Regression model missed most churners (recall of only 16%), meaning relying on it would result in many customers slipping away unnoticed.

With Random Forest achieving 85% recall, the company can capture the majority of churn-prone customers. These high-risk customers should be prioritized with personalized retention offers (discounts, loyalty rewards, or service upgrades).

3. Service Quality Improvements

The Random Forest confusion matrix shows that churners are detected much more effectively compared to Logistic Regression. This suggests churn is strongly tied to specific service and plan features. The company should analyze these drivers and invest in improving customer service and optimizing plan options.

4. Data-Driven Actions

The model performance indicates that class imbalance was successfully handled with SMOTE and cost-sensitive methods. This shows the company must maintain balanced, clean, and updated customer data for the model to stay effective.

Key features like customer service interactions, plan types, and call behaviors should be continuously monitored, as they are likely the main churn predictors.

5. Long-Term Strategy

Develop a retention dashboard powered by the Random Forest model to give managers weekly/monthly churn risk reports.

Over time, test more advanced models (XGBoost, LightGBM) to further improve prediction, especially for borderline churners.

Use predictive churn scores not only for retention but also to identify upselling opportunities for customers who are satisfied but might respond positively to targeted offers.