

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.preprocessing import LabelEncoder
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import (
8 classification_report, confusion_matrix, accuracy_score, roc_curve, auc)
9 from sklearn.linear_model import LogisticRegression
10 from imblearn.over_sampling import SMOTE

```

```


1 # Load the dataset
2 df = pd.read_csv('Customer Churn Data.csv')

```

```

1 # Show first five rows
2 df.head()


```



	customerID	gender	SeniorCitizen	Partner	Dependents	tenure	PhoneService	MultipleLines	InternetService	OnlineSecurity	...
0	7590-VHVEG	Female	0	Yes	No	1	No	No phone service	DSL	No	...
1	5575-GNVDE	Male	0	No	No	34	Yes	No	DSL	Yes	...
2	3668-QPYBK	Male	0	No	No	2	Yes	No	DSL	Yes	...
3	7795-CFOCW	Male	0	No	No	45	No	No phone service	DSL	Yes	...
4	9237-HQITU	Female	0	No	No	2	Yes	No	Fiber optic	No	...

5 rows × 21 columns

```
1 df.info()
```

 <class 'pandas.core.frame.DataFrame'>
RangeIndex: 7043 entries, 0 to 7042
Data columns (total 21 columns):

#	Column	Non-Null Count	Dtype
0	customerID	7043 non-null	object
1	gender	7043 non-null	object
2	SeniorCitizen	7043 non-null	int64
3	Partner	7043 non-null	object
4	Dependents	7043 non-null	object
5	tenure	7043 non-null	int64
6	PhoneService	7043 non-null	object
7	MultipleLines	7043 non-null	object
8	InternetService	7043 non-null	object
9	OnlineSecurity	7043 non-null	object
10	OnlineBackup	7043 non-null	object
11	DeviceProtection	7043 non-null	object
12	TechSupport	7043 non-null	object
13	StreamingTV	7043 non-null	object
14	StreamingMovies	7043 non-null	object
15	Contract	7043 non-null	object
16	PaperlessBilling	7043 non-null	object
17	PaymentMethod	7043 non-null	object
18	MonthlyCharges	7043 non-null	float64
19	TotalCharges	7043 non-null	object
20	Churn	7043 non-null	object

dtypes: float64(1), int64(2), object(18)
memory usage: 1.1+ MB

```

1 #Breif Description of Data
2 df.describe()

```

	SeniorCitizen	tenure	MonthlyCharges
count	7043.000000	7043.000000	7043.000000
mean	0.162147	32.371149	64.761692
std	0.368612	24.559481	30.090047
min	0.000000	0.000000	18.250000
25%	0.000000	9.000000	35.500000
50%	0.000000	29.000000	70.350000
75%	0.000000	55.000000	89.850000
max	1.000000	72.000000	118.750000

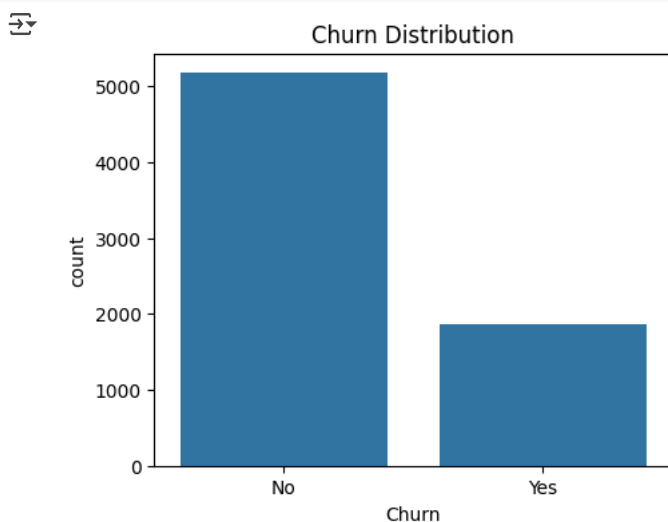
```
1 #all coulms names
2 df.columns
```

```
Index(['customerID', 'gender', 'SeniorCitizen', 'Partner', 'Dependents',
      'tenure', 'PhoneService', 'MultipleLines', 'InternetService',
      'OnlineSecurity', 'OnlineBackup', 'DeviceProtection', 'TechSupport',
      'StreamingTV', 'StreamingMovies', 'Contract', 'PaperlessBilling',
      'PaymentMethod', 'MonthlyCharges', 'TotalCharges', 'Churn'],
      dtype='object')
```

```
1 print(df.isnull().sum())
```

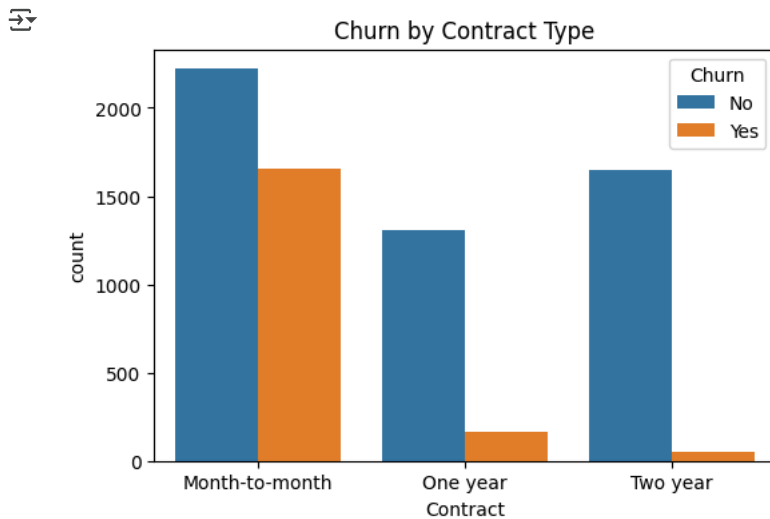
```
customerID      0
gender          0
SeniorCitizen   0
Partner         0
Dependents      0
tenure          0
PhoneService    0
MultipleLines   0
InternetService 0
OnlineSecurity  0
OnlineBackup    0
DeviceProtection 0
TechSupport     0
StreamingTV     0
StreamingMovies 0
Contract        0
PaperlessBilling 0
PaymentMethod   0
MonthlyCharges  0
TotalCharges    0
Churn           0
dtype: int64
```

```
1 #first we see the distribution of churn
2 plt.figure(figsize=(5,4))
3 sns.countplot(x='Churn', data=df)
4 plt.title('Churn Distribution')
5 plt.show()
```



```
1 #here we can see the Churn By Contract Type
2 plt.figure(figsize=(6,4))
3 sns.countplot(x='Contract', hue='Churn', data=df)
```

```
4 plt.title('Churn by Contract Type')
5 plt.show()
```



```
1 #handling the missing values
2 df['TotalCharges'] = pd.to_numeric(df['TotalCharges'], errors='coerce')
3 df['TotalCharges'] = df['TotalCharges'].fillna(df['TotalCharges'].median())
```

```
1 #dropping Unnecessary coulumn
2 df.drop('customerID', axis=1, inplace=True)
```

```
1 #Converts all string,object columns into numbers.
2 from sklearn.preprocessing import LabelEncoder
3
4 df['Churn'] = df['Churn'].map({'Yes': 1, 'No': 0})
5
6 le = LabelEncoder()
7 for col in df.select_dtypes(include=['object']).columns:
8     df[col] = le.fit_transform(df[col])
```

```
1 #splitting the data
2 X = df.drop('Churn', axis=1)
3 y = df['Churn']
4 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=60, stratify=y)
```

```
1 #handling imbalanced data with SMOTE
2 print("Before SMOTE:", y_train.value_counts())
3 smote = SMOTE(random_state=42)
4 X_train, y_train = smote.fit_resample(X_train, y_train)
5 print("After SMOTE:", pd.Series(y_train).value_counts())
```

```
Before SMOTE: Churn
0    3622
1    1308
Name: count, dtype: int64
After SMOTE: Churn
0    3622
1    3622
Name: count, dtype: int64
```

```
1 #traning the model
2 model = LogisticRegression(max_iter=1000, random_state=50)
3 model.fit(X_train, y_train)
```

```
/usr/local/lib/python3.11/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to converge (status=
STOP: TOTAL NO. OF ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html>

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression

```
n_iter_i = _check_optimize_result(
```

```
    LogisticRegression
```

```
LogisticRegression(max_iter=1000, random_state=50)
```

```
1 y_pred = model.predict(X_test)
2 print("Accuracy:", accuracy_score(y_test, y_pred))
```

```
3 print("\nConfusion Matrix:\n", confusion_matrix(y_test, y_pred))
4 print("\nClassification Report:\n", classification_report(y_test, y_pred))
```

↗ Accuracy: 0.7647893989588264

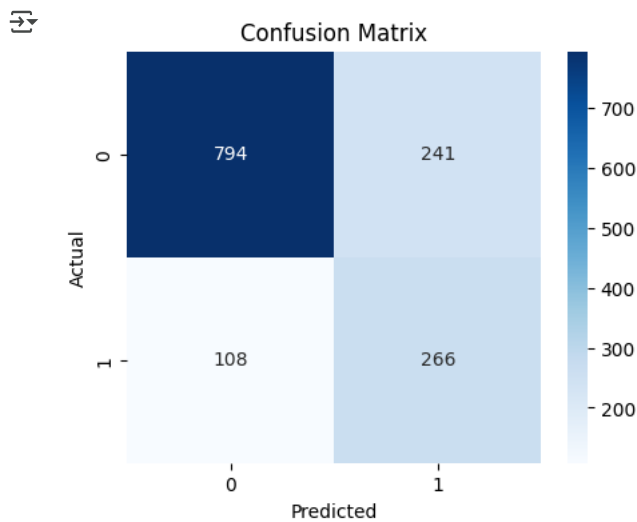
Confusion Matrix:

```
[[1197 355]
 [ 142 419]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.89	0.77	0.83	1552
1	0.54	0.75	0.63	561
accuracy			0.76	2113
macro avg	0.72	0.76	0.73	2113
weighted avg	0.80	0.76	0.77	2113

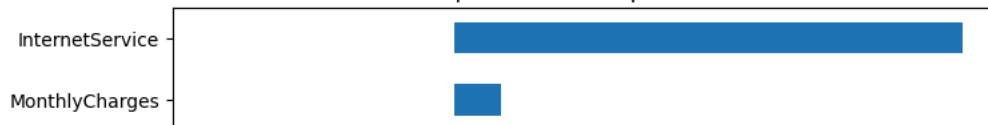
```
1 cm = confusion_matrix(y_test, y_pred)
2 plt.figure(figsize=(5,4))
3 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
4 plt.xlabel('Predicted')
5 plt.ylabel('Actual')
6 plt.title('Confusion Matrix')
7 plt.show()
```



```
1 # Get feature importance from logistic regression coefficients
2 feature_importance = pd.Series(model.coef_[0], index=X.columns).sort_values(ascending=False)
3
4 plt.figure(figsize=(8,6))
5 feature_importance.head(10).plot(kind='barh')
6 plt.title('Top 10 Feature Importances')
7 plt.xlabel('Coefficient Value')
8 plt.gca().invert_yaxis()
9 plt.show()
```



Top 10 Feature Importances



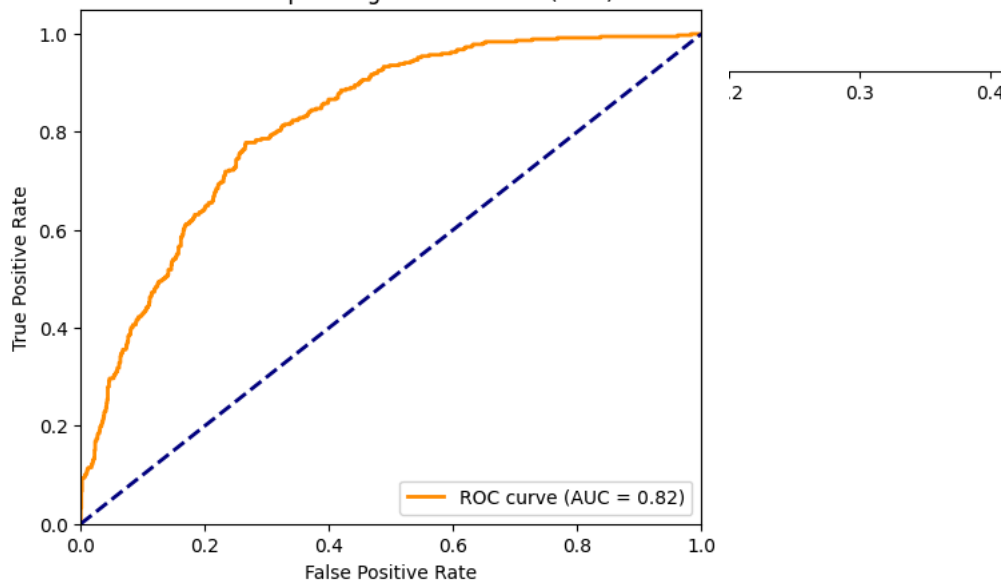
```

1 y_pred_proba = model.predict_proba(X_test)[: ,1]
2 fpr, tpr, thresholds = roc_curve(y_test, y_pred_proba)
3 roc_auc = auc(fpr, tpr)
4
5 plt.figure(figsize=(6,5))
6 plt.plot(fpr, tpr, color='darkorange', lw=2, label='ROC curve (AUC = %0.2f)' % roc_auc)
7 plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
8 plt.xlim([0.0, 1.0])
9 plt.ylim([0.0, 1.05])
10 plt.xlabel('False Positive Rate')
11 plt.ylabel('True Positive Rate')
12 plt.title('Receiver Operating Characteristic (ROC)')
13 plt.legend(loc="lower right")
14 plt.show()

```



Receiver Operating Characteristic (ROC)



```

1 # Print final accuracy and AUC
2 print(f"Final Accuracy: {accuracy_score(y_test, y_pred):.4f}")
3 print(f"Final ROC AUC: {roc_auc:.4f}")

```



Final Accuracy: 0.7523
Final ROC AUC: 0.8180

```

1 from sklearn.preprocessing import StandardScaler
2 from sklearn.linear_model import LogisticRegression
3
4 # Scale the data
5 scaler = StandardScaler()
6 X_train_sm_scaled = scaler.fit_transform(X_train)
7 X_test_scaled = scaler.transform(X_test)
8
9 # Fit the model
10 model = LogisticRegression(max_iter=2000, random_state=80)
11 model.fit(X_train_sm_scaled, y_train)
12
13 # Predict
14 y_pred = model.predict(X_test_scaled)

```

1 Start coding or [generate](#) with AI.