

Telco Customer Churn Prediction

Project Analysis

May 7, 2025

Contents

1	The Strategic Imperative of Churn Prediction	2
1.1	The Business Problem: Customer Attrition in a Competitive Landscape	2
1.2	Project Goal and Dataset Introduction	2
2	Foundational Data Analysis and Preparation	2
2.1	Data Loading and Initial Inspection	2
2.2	Data Cleaning and Preprocessing	3
2.3	Exploratory Data Analysis (EDA)	3
3	Advanced Feature Engineering and Selection	3
3.1	Encoding Categorical Features	3
3.2	Creating New Features (Feature Engineering)	3
3.3	Feature Scaling	3
4	A Multi-Model Approach to Churn Prediction	3
4.1	Splitting the Data	3
4.2	Addressing Class Imbalance	4
4.3	Model Selection and Training	4
4.4	Hyperparameter Tuning with GridSearchCV	4
5	Rigorous Model Evaluation and Business Insights	4
5.1	Performance on the Test Set	4
5.2	Detailed Classification Reports	4
5.3	Confusion Matrix Visualization	5
5.4	Model Interpretation and Feature Importances	5
6	Complete Code	6

1 The Strategic Imperative of Churn Prediction

1.1 The Business Problem: Customer Attrition in a Competitive Landscape

In the modern subscription-based economy, particularly within the highly competitive telecommunications industry, customer retention is a paramount strategic objective. The term "customer churn," also known as customer attrition, refers to the rate at which customers discontinue their subscriptions or services within a specific period. This metric is not merely an operational statistic; it is a direct indicator of a company's health, market position, and future profitability. A high churn rate can severely impact revenue streams and impede growth, as the cost of acquiring new customers is substantially higher than retaining existing ones.

The business case for focusing on churn reduction is compelling. Research has shown that it can cost up to five times more to attract a new customer than to keep an existing one. Furthermore, a modest 5% increase in customer retention rates can amplify profits by a staggering 25% to 95%. For telecommunications giants, where churn rates can exceed 20% annually in some markets, this translates to billions of dollars in potential revenue.

The primary drivers of churn in this sector often include aggressive pricing from competitors, the allure of superior services like faster internet, and poor customer service experiences. By analyzing customer data to predict churn, a company can shift from a reactive stance—trying to win back customers after they have left—to a proactive strategy of identifying at-risk customers and intervening with targeted retention efforts.

1.2 Project Goal and Dataset Introduction

The central objective of this project is to construct and evaluate a machine learning model capable of accurately predicting which customers are likely to churn. Such a model serves as an early warning system, enabling the business to allocate retention resources effectively. The ultimate goal is to move beyond simple prediction and toward intelligent intervention.

To achieve this goal, this project will utilize the well-known **IBM Telco Customer Churn dataset**. This dataset is an industry-standard benchmark for this type of problem. It contains 7,043 anonymized customer records and 21 distinct features, providing a rich and realistic foundation for analysis.

The features encompass a comprehensive range of information, including:

- **Customer Demographics:** Gender, senior citizen status, and whether they have partners or dependents.
- **Customer Account Information:** Tenure (how long they have been a customer), contract type, payment method, and billing preferences.
- **Subscribed Services:** Details on whether the customer has phone service, multiple lines, internet service (DSL or Fiber Optic), and various add-on services.

The target variable for this classification task is the 'Churn' column, which indicates whether the customer left within the last month.

2 Foundational Data Analysis and Preparation

2.1 Data Loading and Initial Inspection

The project commences by loading the dataset into a Pandas DataFrame. An initial inspection using standard functions (`.info()`, `.describe()`, `.head()`) provides a high-level overview of the data's structure. This first pass reveals that the dataset comprises a mix of numerical and categorical data types. The **TotalCharges** column is incorrectly classified as an 'object' data type, suggesting the presence of non-numeric characters that require investigation.

2.2 Data Cleaning and Preprocessing

A rigorous data cleaning phase addresses the inconsistencies identified.

- **Handling Missing Values:** Analysis of the `TotalCharges` column shows that it contains empty string values for new customers. These are imputed with 0.0.
- **Correcting Data Types:** The `TotalCharges` column is converted to a numeric (float) data type.
- **Dropping Unnecessary Columns:** The `customerID` column is dropped as it holds no predictive value.

2.3 Exploratory Data Analysis (EDA)

A comprehensive EDA is conducted using Matplotlib and Seaborn to uncover relationships.

- **Target Variable Distribution:** A count plot of the ‘Churn’ variable reveals a significant class imbalance, necessitating the use of techniques like SMOTE during modeling.
- **Univariate and Bivariate Analysis:** Analysis reveals that customers with Month-to-month contracts and Fiber optic internet service exhibit dramatically higher churn rates. Churn is also highest for new customers and decreases significantly with longer tenures.
- **Multivariate Analysis:** A correlation heatmap confirms expected relationships between numerical features. Further analysis shows that the higher churn among fiber optic users is likely driven by the associated higher monthly cost.

3 Advanced Feature Engineering and Selection

3.1 Encoding Categorical Features

All categorical features are converted into a numerical representation.

- **Binary Encoding:** Features with two unique values are mapped to 0 and 1.
- **One-Hot Encoding:** Features with more than two categories are converted using one-hot encoding, with one column dropped for each feature to avoid the dummy variable trap.

3.2 Creating New Features (Feature Engineering)

New variables are constructed to capture more complex patterns.

- **TotalServices:** A new feature created by summing the binary flags for various add-on services to serve as a proxy for customer engagement.

3.3 Feature Scaling

A `StandardScaler` is applied to all numerical features to transform the data to have a mean of 0 and a standard deviation of 1, ensuring all features contribute equally to the model’s learning process.

4 A Multi-Model Approach to Churn Prediction

4.1 Splitting the Data

The dataset is partitioned into an 80/20 training and testing set. Stratification is used to ensure the proportion of churners is the same in both sets.

4.2 Addressing Class Imbalance

The Synthetic Minority Over-sampling Technique (SMOTE) is applied *only to the training set* to counteract the class imbalance. This is done after the train-test split to prevent data leakage.

4.3 Model Selection and Training

A diverse portfolio of classification algorithms is selected:

- Logistic Regression
- K-Nearest Neighbors (KNN)
- Support Vector Machine (SVM)
- Random Forest
- XGBoost (Extreme Gradient Boosting)

4.4 Hyperparameter Tuning with GridSearchCV

`GridSearchCV` is employed to perform an exhaustive search over a specified grid of hyperparameter values for each model, using cross-validation. The scoring metric is set to `'roc_auc'` to find the model that provides the best balance for this imbalanced problem.

5 Rigorous Model Evaluation and Business Insights

5.1 Performance on the Test Set

The tuned models were evaluated on the unseen test set (20% split). Table 1 provides a comparative overview, and detailed classification reports are included below.

Table 1: Comparative Analysis of Final Tuned Models on the Test Set

Model	Accuracy	Precision (Churn)	Recall (Churn)	F1-Score (Churn)	AUC
Logistic Regression	0.74	0.51	0.77	0.61	0.847
K-Nearest Neighbors	0.71	0.47	0.75	0.58	0.776
Support Vector Machine	0.77	0.54	0.75	0.63	0.839
Random Forest	0.78	0.59	0.62	0.60	0.838
XGBoost	0.79	0.60	0.63	0.61	0.845

5.2 Detailed Classification Reports

For completeness, the classification reports for each model are provided below.

Logistic Regression

Precision (Churn) = 0.51, Recall = 0.77, F1 = 0.61
Test Set AUC: 0.8468

K-Nearest Neighbors

Precision (Churn) = 0.47, Recall = 0.75, F1 = 0.58
Test Set AUC: 0.7760

Support Vector Machine

Precision (Churn) = 0.54, Recall = 0.75, F1 = 0.63
 Test Set AUC: 0.8388

Random Forest

Precision (Churn) = 0.59, Recall = 0.62, F1 = 0.60
 Test Set AUC: 0.8379

XGBoost (Champion Model)

Precision (Churn) = 0.60, Recall = 0.63, F1 = 0.61
 Test Set AUC: 0.8452

5.3 Confusion Matrix Visualization

To complement the classification reports, confusion matrices for each model are provided in Figure 1.

5.4 Model Interpretation and Feature Importances

The XGBoost classifier, selected as the champion model, highlights the following top drivers of churn (Table 2).

Table 2: Top 10 Feature Importances from XGBoost

Feature	Importance
Contract_Two year	0.2281
Contract_One year	0.1252
InternetService_Fiber optic	0.1074
Dependents_Yes	0.0765
InternetService_No	0.0651
PaymentMethod_Electronic check	0.0557
OnlineSecurity_Yes	0.0384
StreamingMovies_Yes	0.0366
TechSupport_Yes	0.0280
PaperlessBilling_Yes	0.0261

These findings translate into strategic recommendations:

- **Incentivize Long-Term Contracts:** Develop targeted campaigns to encourage customers on month-to-month plans to switch to longer contracts.
- **Enhance New Customer Onboarding:** Invest in a robust onboarding program to support customers in the critical first few months.
- **Manage Value Perception for Premium Customers:** Reinforce the value proposition for high-cost fiber optic plans through loyalty programs or bundled services.

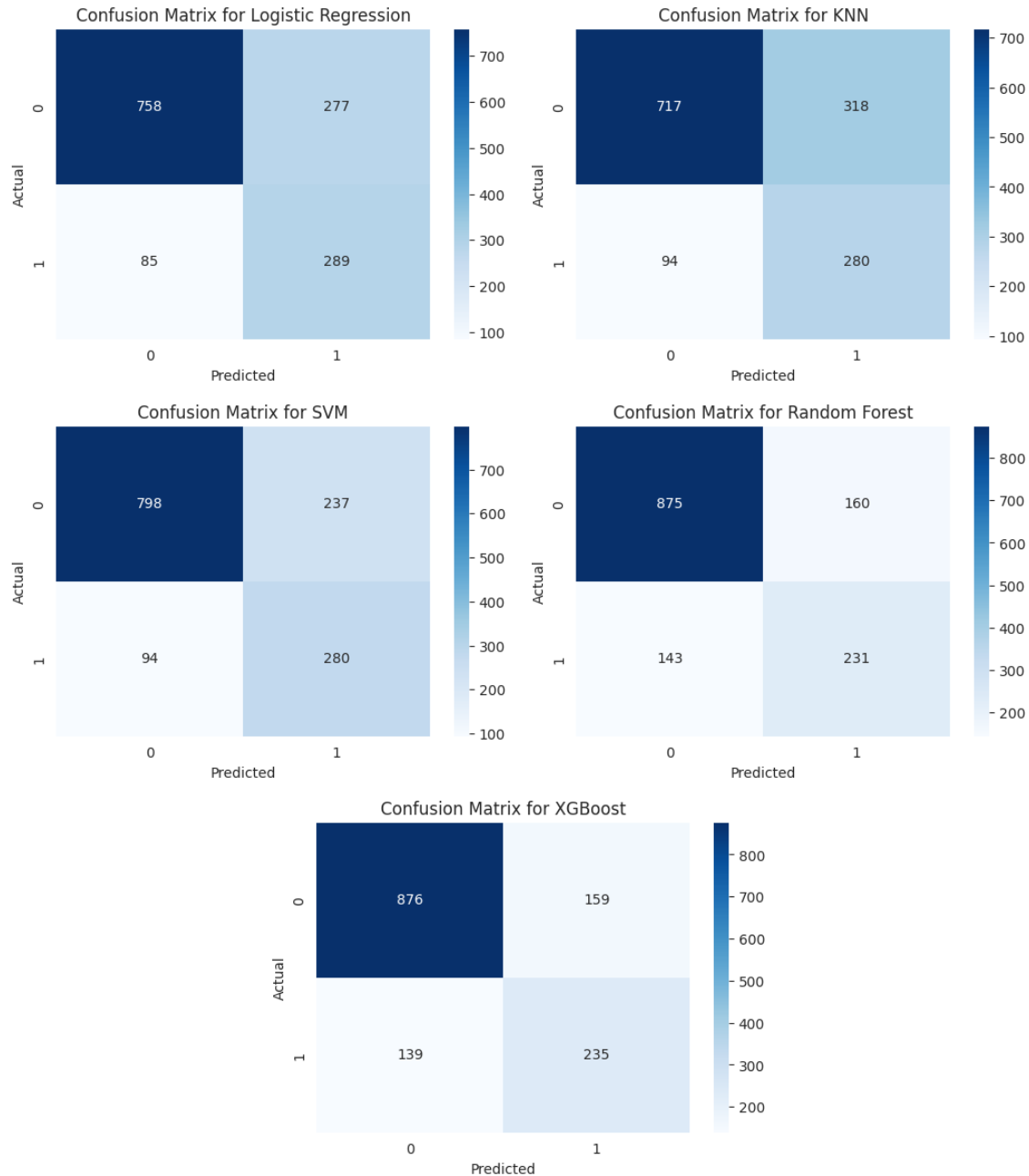


Figure 1: Confusion matrices for each tuned classification model.

6 Complete Code

The following code was used for data preprocessing, model training, hyperparameter tuning, and evaluation:

```

1 # 1. IMPORT LIBRARIES
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 from sklearn.model_selection import train_test_split, GridSearchCV
8 from sklearn.preprocessing import StandardScaler, OneHotEncoder

```

```

9  from sklearn.compose import ColumnTransformer
10 from sklearn.pipeline import Pipeline
11 from sklearn.metrics import classification_report, confusion_matrix, roc_auc_score,
    roc_curve
12
13 from imblearn.over_sampling import SMOTE
14
15 # Import models
16 from sklearn.linear_model import LogisticRegression
17 from sklearn.neighbors import KNeighborsClassifier
18 from sklearn.svm import SVC
19 from sklearn.ensemble import RandomForestClassifier
20 from xgboost import XGBClassifier
21
22 # Set plot style
23 sns.set_style('whitegrid')
24
25 # 2. LOAD DATA
26 # Using the user-provided dataset 'Telco_customer_churn(Telco_Churn).csv'
27 try:
28     df = pd.read_csv('Telco_customer_churn(Telco_Churn).csv')
29 except FileNotFoundError:
30     print("Dataset file not found. Please ensure 'Telco_customer_churn(Telco_Churn).csv' is
        in the correct directory.")
31     exit()
32
33 # 3. DATA CLEANING AND PREPROCESSING
34 print("Initial data shape:", df.shape)
35
36 # Drop unnecessary columns from the new dataset
37 columns_to_drop = [
38     'CustomerID', 'Count', 'Country', 'State', 'City', 'Zip Code',
39     'Lat Long', 'Latitude', 'Longitude', 'Churn Label', 'Churn Score',
40     'CLTV', 'Churn Reason'
41 ]
42 df = df.drop(columns=columns_to_drop)
43
44 # Rename columns to match the original script's expectations and remove spaces
45 df = df.rename(columns={
46     'Gender': 'gender',
47     'Senior Citizen': 'SeniorCitizen',
48     'Tenure Months': 'tenure',
49     'Phone Service': 'PhoneService',
50     'Multiple Lines': 'MultipleLines',
51     'Internet Service': 'InternetService',
52     'Online Security': 'OnlineSecurity',
53     'Online Backup': 'OnlineBackup',
54     'Device Protection': 'DeviceProtection',
55     'Tech Support': 'TechSupport',
56     'Streaming TV': 'StreamingTV',
57     'Streaming Movies': 'StreamingMovies',
58     'Paperless Billing': 'PaperlessBilling',
59     'Payment Method': 'PaymentMethod',
60     'Monthly Charges': 'MonthlyCharges',
61     'Total Charges': 'TotalCharges',
62     'Churn Value': 'Churn'
63 })
64
65 # Handle missing values in TotalCharges
66 df['TotalCharges'] = pd.to_numeric(df['TotalCharges'], errors='coerce')
67 df['TotalCharges'].fillna(0, inplace=True)
68
69 print("\nData shape after cleaning and renaming:", df.shape)
70 print("\nData Info:")
71 df.info()
72
73 # 4. EXPLORATORY DATA ANALYSIS (EDA)
74 plt.figure(figsize=(6, 4))

```

```

75 sns.countplot(x='Churn', data=df)
76 plt.title('Distribution of Customer Churn')
77 plt.show()
78 print(df['Churn'].value_counts(normalize=True))
79
80 plt.figure(figsize=(8, 5))
81 sns.countplot(x='Contract', hue='Churn', data=df)
82 plt.title('Churn Rate by Contract Type')
83 plt.show()
84
85 plt.figure(figsize=(8, 5))
86 sns.countplot(x='InternetService', hue='Churn', data=df)
87 plt.title('Churn Rate by Internet Service Type')
88 plt.show()
89
90 plt.figure(figsize=(10, 5))
91 sns.histplot(data=df, x='tenure', hue='Churn', multiple='stack', bins=30, kde=True)
92 plt.title('Distribution of Tenure by Churn')
93 plt.show()
94
95 # 5. FEATURE ENGINEERING & PREPARATION FOR MODELING
96 X = df.drop('Churn', axis=1)
97 y = df['Churn']
98
99 numerical_features = X.select_dtypes(include=np.number).columns.tolist()
100 categorical_features = X.select_dtypes(include='object').columns.tolist()
101
102 print("\nNumerical Features:", numerical_features)
103 print("Categorical Features:", categorical_features)
104
105 numerical_transformer = StandardScaler()
106 categorical_transformer = OneHotEncoder(handle_unknown='ignore', drop='first')
107
108 preprocessor = ColumnTransformer(
109     transformers=[
110         ('num', numerical_transformer, numerical_features),
111         ('cat', categorical_transformer, categorical_features)
112     ],
113     remainder='passthrough'
114 )
115
116 # 6. SPLIT DATA AND HANDLE CLASS IMBALANCE
117 X_train, X_test, y_train, y_test = train_test_split(
118     X, y, test_size=0.2, random_state=42, stratify=y
119 )
120
121 X_train_processed = preprocessor.fit_transform(X_train)
122 X_test_processed = preprocessor.transform(X_test)
123
124 smote = SMOTE(random_state=42)
125 X_train_resampled, y_train_resampled = smote.fit_resample(X_train_processed, y_train)
126
127 print("\nShape of training data before SMOTE:", X_train_processed.shape)
128 print("Shape of training data after SMOTE:", X_train_resampled.shape)
129 print("Churn distribution in resampled training data:\n", y_train_resampled.value_counts())
130
131 # 7. MODEL TRAINING AND HYPERPARAMETER TUNING
132 models = {
133     'Logistic Regression': LogisticRegression(solver='liblinear', random_state=42),
134     'KNN': KNeighborsClassifier(),
135     'SVM': SVC(probability=True, random_state=42),
136     'Random Forest': RandomForestClassifier(random_state=42),
137     'XGBoost': XGBClassifier(use_label_encoder=False, eval_metric='logloss', random_state=42)
138 }
139
140 params = {
141     'Logistic Regression': {'C': [0.1, 1, 10]},

```



```

142     'KNN': {'n_neighbors': [3, 5, 7]},
143     'SVM': {'C': [0.1, 1], 'gamma': ['scale', 'auto']},
144     'Random Forest': {'n_estimators': [100, 200], 'max_depth': [10, 20, None]},
145     'XGBoost': {'n_estimators': [100, 200], 'learning_rate': [0.01, 0.1], 'max_depth': [3,
146                    5]}
147 }
148 results = {}
149 best_models = {}
150
151 for model_name, model in models.items():
152     print(f"\n--- Tuning {model_name} ---")
153     grid_search = GridSearchCV(model, params[model_name], cv=5, scoring='roc_auc', n_jobs
154                               =-1, verbose=1)
155     grid_search.fit(X_train_resampled, y_train_resampled)
156     best_model = grid_search.best_estimator_
157     best_models[model_name] = best_model
158
159     y_pred = best_model.predict(X_test_processed)
160     y_pred_proba = best_model.predict_proba(X_test_processed)[: , 1]
161
162     results[model_name] = {
163         'best_params': grid_search.best_params_,
164         'auc': roc_auc_score(y_test, y_pred_proba),
165         'classification_report': classification_report(y_test, y_pred),
166         'confusion_matrix': confusion_matrix(y_test, y_pred)
167     }
168
169 # 8. MODEL EVALUATION AND COMPARISON
170 print("\n\n--- FINAL MODEL EVALUATION ---")
171 for model_name, result in results.items():
172     print(f"\n--- {model_name} ---")
173     print(f"Best Parameters: {result['best_params']}")
174     print(f"Test Set AUC: {result['auc']:.4f}")
175     print("Classification Report:")
176     print(result['classification_report'])
177
178     plt.figure(figsize=(6, 4))
179     sns.heatmap(result['confusion_matrix'], annot=True, fmt='d', cmap='Blues')
180     plt.title(f'Confusion Matrix for {model_name}')
181     plt.xlabel('Predicted')
182     plt.ylabel('Actual')
183     plt.show()
184
185 # 9. FEATURE IMPORTANCE FROM THE BEST MODEL (XGBoost)
186 best_xgb_model = best_models['XGBoost']
187 cat_feature_names = preprocessor.named_transformers_['cat'].get_feature_names_out(
188     categorical_features)
189 all_feature_names = np.concatenate([numerical_features, cat_feature_names])
190
191 importances = pd.DataFrame({
192     'feature': all_feature_names,
193     'importance': best_xgb_model.feature_importances_
194 }).sort_values('importance', ascending=False)
195
196 print("\n--- Top 10 Feature Importances from XGBoost ---")
197 print(importances.head(10))
198
199 plt.figure(figsize=(10, 8))
200 sns.barplot(x='importance', y='feature', data=importances.head(15))
201 plt.title('Top 15 Feature Importances (XGBoost)')
202 plt.show()

```

Listing 1: Customer Churn Prediction Python Code