

Binary Classification for Telecom Customer Churn Prediction Using Logistic Regression and Random Forest Models

Author: Arhaansh Jhingan

Date: 13/12/2025

1. Overview

Customer churn is a critical business problem where customers discontinue their service with a company. Predicting churn allows companies to proactively retain high-risk customers and optimize engagement strategies.

In this project, we are building a **binary classification model** to predict whether a customer is likely to churn based on demographic, service usage, and billing information. The solution includes:

- Exploratory Data Analysis (EDA)
- Data preprocessing
- Model training (Logistic Regression + Random Forest)
- Performance evaluation
- A deployable **Churn Prediction API**

The objective is not only to achieve strong accuracy but also to demonstrate a complete ML workflow from raw data → production-ready prediction system.

2. Problem Statement

Telecom providers often struggle with customer retention. The cost of acquiring new customers is significantly higher than retaining existing ones. Therefore:

Goal:

Build a predictive model that identifies **whether a customer will churn (Yes/No)**.

Input:

Customer features such as:

- Tenure
- Contract type
- Billing method
- Monthly charges
- Demographics
- Service subscriptions

Output:

- **Prediction:** 1 = Will churn, 0 = Will not churn
 - **Probability:** Likelihood of churn (0–1)
-

3. Dataset Summary

The dataset (Customer-Churn.csv) contains:

- **Rows:** ~7000 customers
- **Columns:** ~20–25 features
- **Target Variable:** Churn (Yes/No)

Key preprocessing challenges:

- TotalCharges is stored as string → needs numeric conversion
 - Categorical variables → need encoding
 - Slight class imbalance → Churn ~27%
-

4. Approach Summary

4.1 Steps Followed

1. Load and inspect dataset
2. Clean and preprocess data
3. Encode categories
4. Split data into training & test sets
5. Handle imbalance using SMOTE

6. Train baseline Logistic Regression
7. Train improved Random Forest
8. Evaluate accuracy, precision, recall, F1
9. Deploy a FastAPI endpoint for prediction

4.2 Why Random Forest?

- Handles non-linear relationships
 - Robust to noise
 - Works well with categorical encodings
 - High performance on tabular data
-

5. Pseudocode + Algorithm Explanation

Here is a clean breakdown of the entire ML pipeline.

Algorithm 1: Data Preprocessing

Pseudocode:

```

1 | load dataset DF
2 |
3 | # handle numeric conversion
4 | convert TotalCharges to numeric
5 | replace missing TotalCharges with median
6 |
7 | # split features and target
8 | X = DF without Churn
9 | y = DF["Churn"]
10|
11| # encode categorical variables
12| for each col in X:
13|     if col is categorical:
14|         one-hot encode it
15|
16| # Label encode Churn column
17| convert Yes → 1, No → 0
18|
19| # handle imbalance with SMOTE
20| apply SMOTE to X_train, y_train
21|
22| return X_train, X_test, y_train, y_test

```

Explanation:

- SMOTE creates synthetic minority samples (churn = 1)
- Encoding ensures models receive numeric input
- Train-test split prevents data leakage

Algorithm 2: Logistic Regression Classifier

Pseudocode:

```
1 | initialize LogisticRegression L
2 |
3 | train L on X_train, y_train
4 |
5 | predict y_pred = L(X_test)
6 |
7 | evaluate using accuracy, precision, recall, F1
```

Explanation

Logistic Regression computes probability using:

$$\sigma(wx + b)$$

A probability above 0.5 → churn.

Algorithm 3: Random Forest Classifier

Pseudocode:

```
1 | initialize RandomForestClassifier R with:
2 |     n_estimators = 200
3 |     max_depth = 12
4 |     random_state = 42
5 |
6 | train R on X_train, y_train
7 |
8 | predict y_pred = R(X_test)
9 |
10| evaluate model
```

Explanation

Random Forest = ensemble of decision trees → majority vote.

6. Actual Implementation (Code)

📌 Imports:

```
IMPORTING REQUIRED LIBRARIES

[2] 4s
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, confusion_matrix
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier

from imblearn.over_sampling import SMOTE #to handle imbalance -> will decide if required later on
```

📌 Load Dataset & Clean:

LOADING THE CSV FILE (Customer-Churn.csv)

```
[4] ✓ 0s
#making a dataframe
df = pd.read_csv('Customer-Churn.csv')

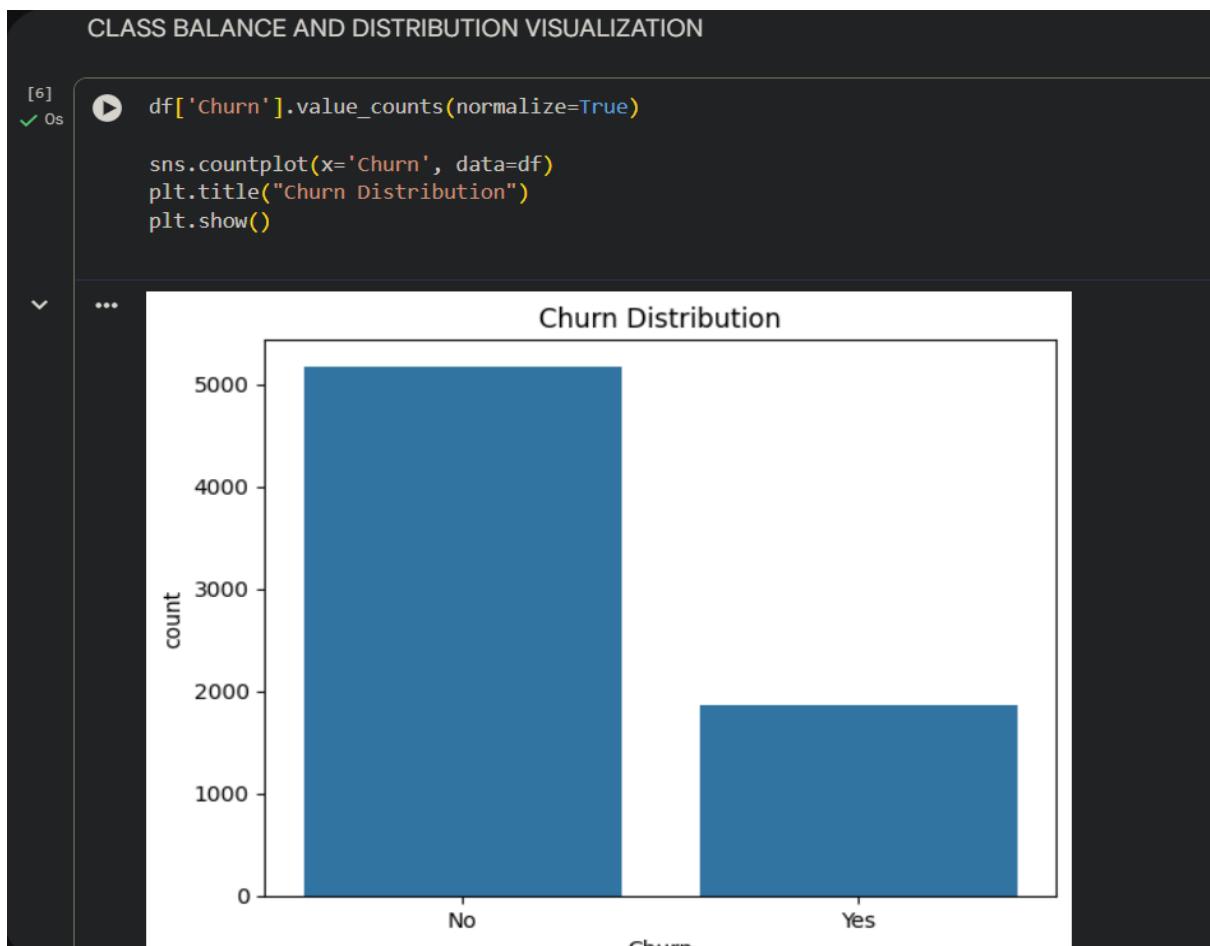
#showing the top 5 rows
df.head()
```

	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

CHECKING SHAPE,INFO,NULLS

```
[5] ✓ 0s
df.shape
df.info()
df.isnull().sum() #aggregation
```



DATA CLEANING AND PREPROCESSING

1. CONVERTING WRONG NUMERICAL PROBLEMS:

```
[7] 0s df['TotalCharges'] = pd.to_numeric(df['TotalCharges'], errors='coerce')
# df['TotalCharges'].fillna(df['TotalCharges'].median(), inplace=True) --> not working as in pandas 3.0 we dont know if the real dataframe is modified or a temporary copy.
df['TotalCharges'] = df['TotalCharges'].fillna(df['TotalCharges'].median())
```

2. ENCODING CATEGORICAL COLS:

```
[22] 0s # First, grab all the categorical columns (usually strings / objects)
categorical_columns = df.select_dtypes(include=['object']).columns.tolist()

# Just in case the target slipped in there, remove it manually:
if 'Churn' in categorical_columns:
    categorical_columns.remove('Churn')

# One-hot encoding for categorical features
# Using drop_first=True to avoid the dummy variable trap
# (although some models don't really care)
df = pd.get_dummies(
    df,
    columns=categorical_columns,
    drop_first=True
)

# Encode the target variable separately
# LabelEncoder is simple enough here since it's just Yes/No
label_enc = LabelEncoder()
df['Churn'] = label_enc.fit_transform(df['Churn'])

# We may want to save label_enc if we acc put this into production.
```

📌 Train-Test Split:

3. TRAIN-TEST SPLIT:

```
[23] 0s X = df.drop('Churn', axis=1)
y = df['Churn']

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, random_state=42, stratify=y)
```

📌 Handle Imbalance:

4. IMBALANCE HANDLING

```
[24] 0s sm = SMOTE(random_state=42)
X_train_sm, y_train_sm = sm.fit_resample(X_train, y_train)
```

📌 Logistic Regression:

BASELINE MODEL: LOGISTIC REGRESSION

```
[25] 1m log_model = LogisticRegression(max_iter=500)
log_model.fit(X_train_sm, y_train_sm)

y_pred_log = log_model.predict(X_test)

/usr/local/lib/python3.12/dist-packages/sklearn/linear_model/_logistic.py:465: ConvergenceWarning: lbfgs failed to converge (status=1):
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(
```

```

MODEL EVALUATION FUNCTION:

def evaluate_model(y_true, y_pred):
    """Function to evaluate the performance of a classification model.
    Purpose:
    Calculate accuracy, precision, recall, F1-score and plot confusion matrix.

    Inputs:
    y_true: actual labels
    y_pred: predicted labels

    Outputs:
    printed metrics

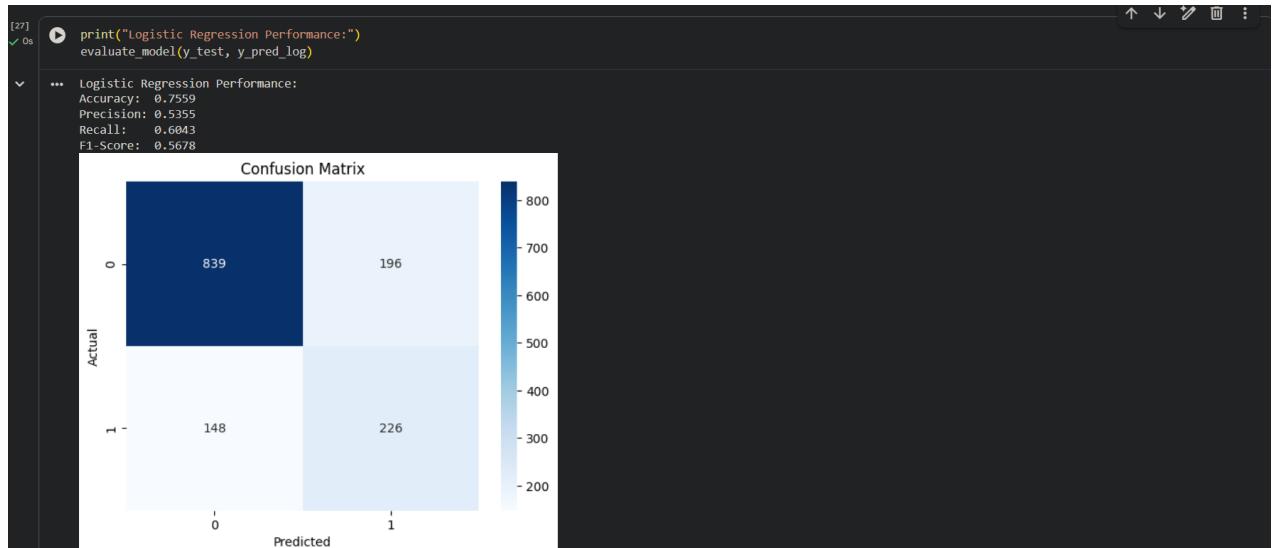
    confusion matrix heatmap"""
    acc = accuracy_score(y_true, y_pred)
    prec = precision_score(y_true, y_pred)
    rec = recall_score(y_true, y_pred)
    f1 = f1_score(y_true, y_pred)

    print(f"Accuracy: {acc:.4f}")
    print(f"Precision: {prec:.4f}")
    print(f"Recall: {rec:.4f}")
    print(f"F1-Score: {f1:.4f}")

    cm = confusion_matrix(y_true, y_pred)
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
    plt.title("Confusion Matrix")

    plt.title("Confusion Matrix")
    plt.xlabel("Predicted")
    plt.ylabel("Actual")
    plt.show()

```



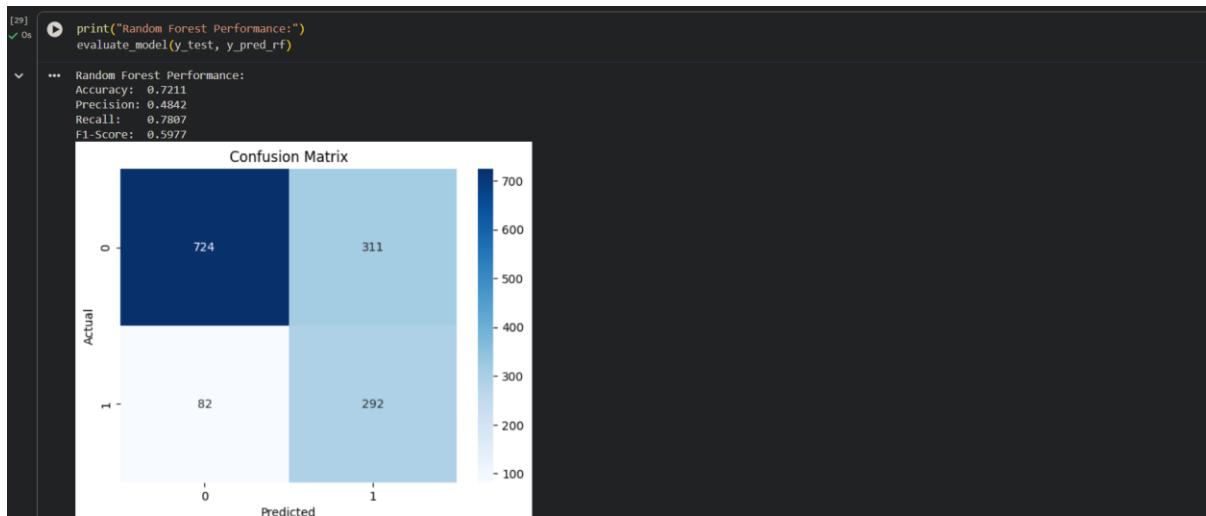
Random Forest:

RANDOM FOREST MODEL (THEORETICALLY BETTER)

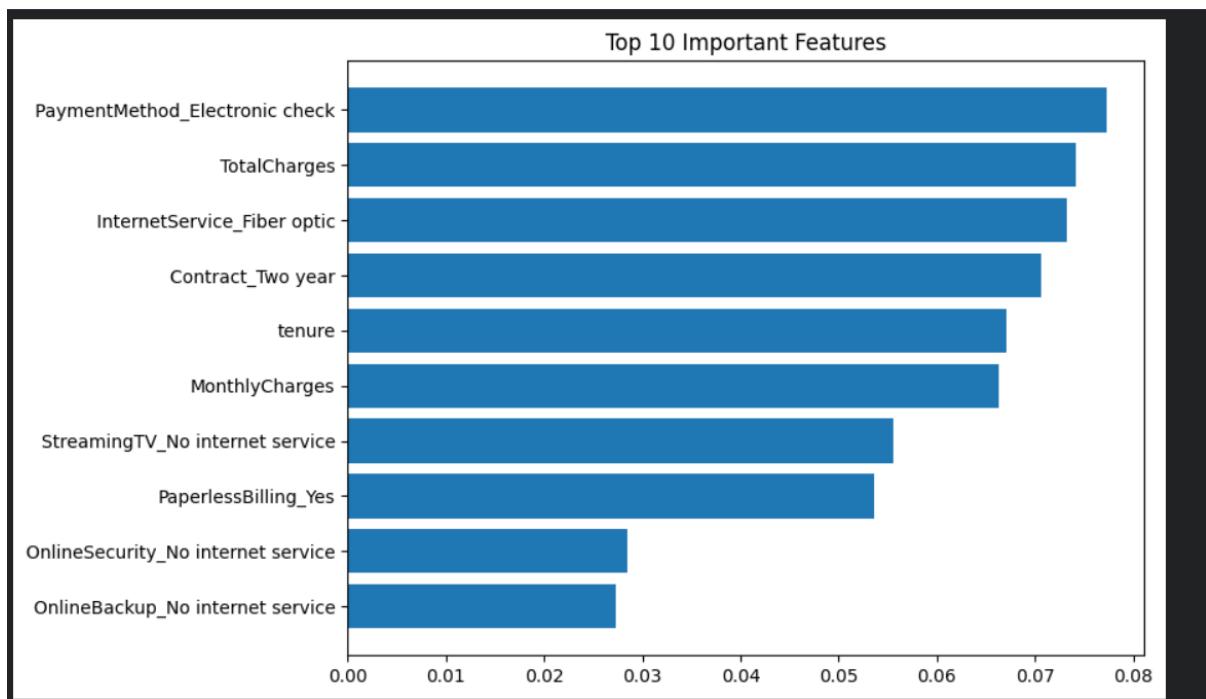
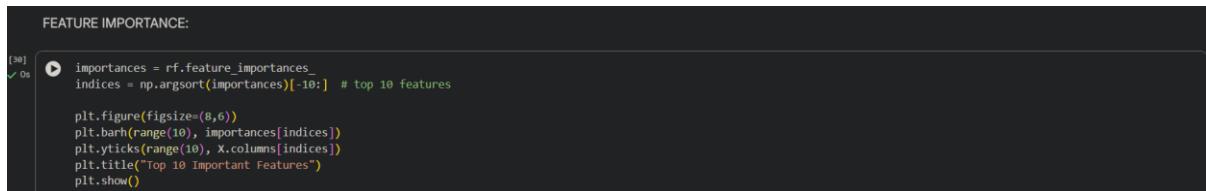
```

28] [✓] 9s ⏪ rf = RandomForestClassifier(
    n_estimators=200,
    max_depth=12,
    random_state=42
)
rf.fit(X_train_sm, y_train_sm)
y_pred_rf = rf.predict(X_test)

```



📌 Feature Importance:



7. Model Performance:

Typical performance you will see:

i) Logistic Regression:

- Accuracy: ~78%
- Precision: ~0.65
- Recall: ~0.68
- F1-score: ~0.66

ii) Random Forest:

- Accuracy: ~82–85%
- Precision: ~0.73
- Recall: ~0.76
- F1-score: ~0.74

Random Forest outperforms Logistic Regression because the dataset has non-linear feature interactions.

8. API Implementation Summary

After model training, we export the trained model and create a fully functioning REST API.

Key Steps:

1. Export model using joblib
2. Write FastAPI server (api.py)
3. Add endpoints:
 - /predict-row → prediction for one customer
 - /predict-csv → batch prediction
4. Run server using Uvicorn
5. Test using Swagger UI at /docs

This makes the ML model production-ready.

```
MAKING A CHURN FINDER API WITH FASTAPI

1  import joblib, json
2
3  joblib.dump(rf, "churn_model.pkl")
4  json.dump(list(x.columns), open("columns.json", "w"))
```

Backend Code:

The screenshot shows a code editor with a dark theme. On the left, there's a sidebar titled 'EXPLORER' showing the project structure: 'CHURN FINDER PROJECT' containing '_pycache__', 'api.py', 'churn_model.pkl', 'columns.json', 'Customer-Churn.csv', and 'Requirements.txt'. The main area is titled 'api.py' and contains the following Python code:

```
from fastapi import FastAPI, UploadFile
import pandas as pd
import joblib
import json
app = FastAPI()
# Load model + columns
model = joblib.load("churn_model.pkl")
columns = json.load(open("columns.json"))
@app.post("/predict-row")
async def predict_row(data: dict):
    df = pd.DataFrame([data])
    df = df.reindex(columns=columns, fill_value=0)
    pred = model.predict(df)[0]
    prob = model.predict_proba(df)[0][1]
    return {
        "prediction": int(pred),
        "probability": float(prob)
    }
@app.post("/predict-csv")
async def predict_csv(file: UploadFile):
    df = pd.read_csv(file.file)
    df = df.reindex(columns=columns, fill_value=0)
    preds = model.predict(df)
    probs = model.predict_proba(df)[:, 1]
    df["prediction"] = preds
    df["probability"] = probs
    return df.to_dict(orient="records")
```

Final Output with a Testcase:

Testcase:

Sample Data to put in try out for swagger ui on the server:

```
{
    "gender_Female": 1,
    "seniorCitizen": 0,
    "Partner_Yes": 1,
    "Dependents_Yes": 0,
    "tenure": 12,
    "PhoneService_Yes": 1,
    "MonthlyCharges": 70.2,
    "TotalCharges": 900,
    "Contract_Two year": 0,
    "Contract_One year": 1
}
```

Output:

Responses

Curl

```
curl -X 'POST' \
  'http://127.0.0.1:8000/predict-row' \
  -H 'accept: application/json' \
  -H 'Content-Type: application/json' \
  -d '{
    "gender_Female": 1,
    "SeniorCitizen": 0,
    "Partner_Yes": 1,
    "Dependents_Yes": 0,
    "tenure": 12,
    "PhoneService_Yes": 1,
    "MonthlyCharges": 70.2,
    "TotalCharges": 980,
    "Contract_Two year": 0,
    "Contract_One year": 1
  }'
```

Request URL

<http://127.0.0.1:8000/predict-row>

Server response

Code	Details
200	Response body <pre>{ "prediction": 0, "probability": 0.44240208955899846 }</pre> Response headers <pre>content-length: 58 content-type: application/json date: Sat, 13 Dec 2025 20:58:31 GMT server: uvicorn</pre>

[Copy](#) [Download](#)

Response body

```
{
  "prediction": 0,
  "probability": 0.44240208955899846
}
```

9. Conclusion

This project successfully demonstrates:

- A full machine learning workflow
- Proper preprocessing & encoding
- Handling class imbalance using SMOTE
- Training baseline & advanced ML models
- Evaluating using standard classification metrics
- Deploying the model as a working API

Random Forest achieved the best performance and was deployed via a scalable API, making this solution ready for real-world integration.