

PHASE-3 Submission: Customer Churn Prediction using Machine Learning

Student Information:

- Student Name: **BALASRI.A**
- Register Number: **732323106005**
- Institution: **SSM COLLEGE OF ENGINEERING**
- Department: **B.E/ECE**
- Date of Submission: **15/05/2025**

GitHub Repository Link @ [HTTPS://GITHUB.COM/BALASRI26/NAAN-MUDHALVAN-/TREE/MAIN](https://github.com/BALASRI26/NAAN-MUDHALVAN-/tree/main)

1 Hardware and Software Requirements

Hardware

- Minimum RAM: 8 GB (required for efficient model training and data processing)
- Processor: Quad-core processor (e.g., Intel i5 or equivalent) for handling computations

Software

- Python Version: 3.8 or higher
- Required Libraries: pandas, numpy, seaborn, matplotlib, scikit-learn, xgboost, streamlit, joblib
- IDE: Google Colab or Jupyter Notebook for development; Streamlit for deployment

2 Objectives

The goal of this project is to predict customer churn in a telecommunications company using machine learning, enabling proactive retention strategies. The specific objectives are:

- Analyze historical customer data to identify key churn drivers such as contract type, tenure, and payment method.

- Develop and evaluate machine learning models (Logistic Regression, Decision Tree, Random Forest, XGBoost) to predict churn with high accuracy and recall.
- Deploy an interactive web app using Streamlit to allow non-technical users to input customer data and receive churn predictions.
- Provide actionable insights for businesses, such as targeting customers with month-to-month contracts for retention offers, to reduce churn rates and improve profitability.

3 Flowchart of Project Workflow

Note: The flowchart illustrating the workflow (Data Collection → Preprocessing → EDA → Feature Engineering → Modeling → Evaluation → Deployment) should be inserted here. It was created using draw.io but is not available in this document.

4 Dataset Description

- **Source:** Kaggle
(<https://www.kaggle.com/datasets/blastchar/telco-customer-churn>)
- **Type:** Public, static dataset
- **Size and Structure:** 7,000 records, 21 columns (including demographics, service usage, tenure, billing, payment method, and target variable Churn)

Note: A screenshot of `df.head()` showing the first 5 rows of the dataset should be inserted here.

5 Data Preprocessing

The following preprocessing steps were applied:

- Converted `TotalCharges` to numeric, handling missing values by dropping rows with NaN (less than 1% of data).
- Removed the `customerID` column as it is irrelevant for prediction.
- Encoded binary features (`Partner`, `Dependents`, `PhoneService`, `PaperlessBilling`) using mapping (Yes → 1, No → 0).
- Applied one-hot encoding to other categorical features (e.g., `Contract`, `PaymentMethod`) using `pd.get_dummies`.
- Scaled numerical features (`tenure`, `MonthlyCharges`, `TotalCharges`) using `MinMaxScaler`.

- Checked for outliers using boxplots and IQR; no significant outliers were removed as they were deemed realistic.

Note: Screenshots showing the dataset before and after preprocessing should be inserted here.

6 Exploratory Data Analysis (EDA)

Visualizations and Insights

- **Churn Distribution:** The target variable `Churn` is imbalanced (approximately 27% churned), necessitating techniques to handle class imbalance.
- **Numerical Features:** Histograms and boxplots revealed that customers with lower `tenure` and higher `MonthlyCharges` are more likely to churn.
- **Categorical Features:** Count plots showed that customers with month-to-month contracts, electronic check payments, and fiber optic internet have higher churn rates.
- **Correlation Heatmap:** Strong correlation between `MonthlyCharges` and `TotalCharges`, weak correlation with `Churn`.
- **Multivariate Analysis:** Senior citizens with month-to-month contracts and fiber optic internet are more likely to churn.

Note: Screenshots of EDA visualizations (e.g., churn distribution, tenure vs. churn, correlation heatmap) should be inserted here.

Key Takeaways

- Customers with month-to-month contracts, electronic check payments, and higher monthly charges are at higher risk of churn.
- Long-tenured customers and those with automatic payments are less likely to churn.
- Contract type and tenure are the most critical predictors of churn.

7 Feature Engineering

- **New Features:** Created `tenure_group` (e.g., 0-12, 13-24 months) to categorize tenure into bins.
- **Derived Features:** Created `hasCIP` based on the count of services (e.g., internet, phone) subscribed by the customer.
- **Feature Selection:** Removed multicollinear features (e.g., between `MonthlyCharges` and `TotalCharges`) to reduce redundancy.
- **Impact:** These features improved model performance by capturing categorical tenure patterns and service usage complexity, which are strong indicators of churn.

8 Model Building

The following models were trained:

- **Logistic Regression:** Chosen for interpretability and as a baseline model.
- **Decision Tree:** Used to capture non-linear relationships and handle categorical features.
- **Random Forest:** An ensemble method to improve robustness and reduce overfitting.
- **XGBoost:** A powerful gradient boosting model for high performance on imbalanced datasets.

Training Details:

- Train-Test Split: 70% training, 30% testing.
- Validation: 5-fold cross-validation.
- Tuning: Used GridSearchCV to optimize hyperparameters (e.g., learning rate for XGBoost, max depth for Random Forest).

Note: Screenshots of model training outputs (e.g., logs or summary statistics) should be inserted here.

9 Model Evaluation

Evaluation Metrics

The models were evaluated using the following metrics:

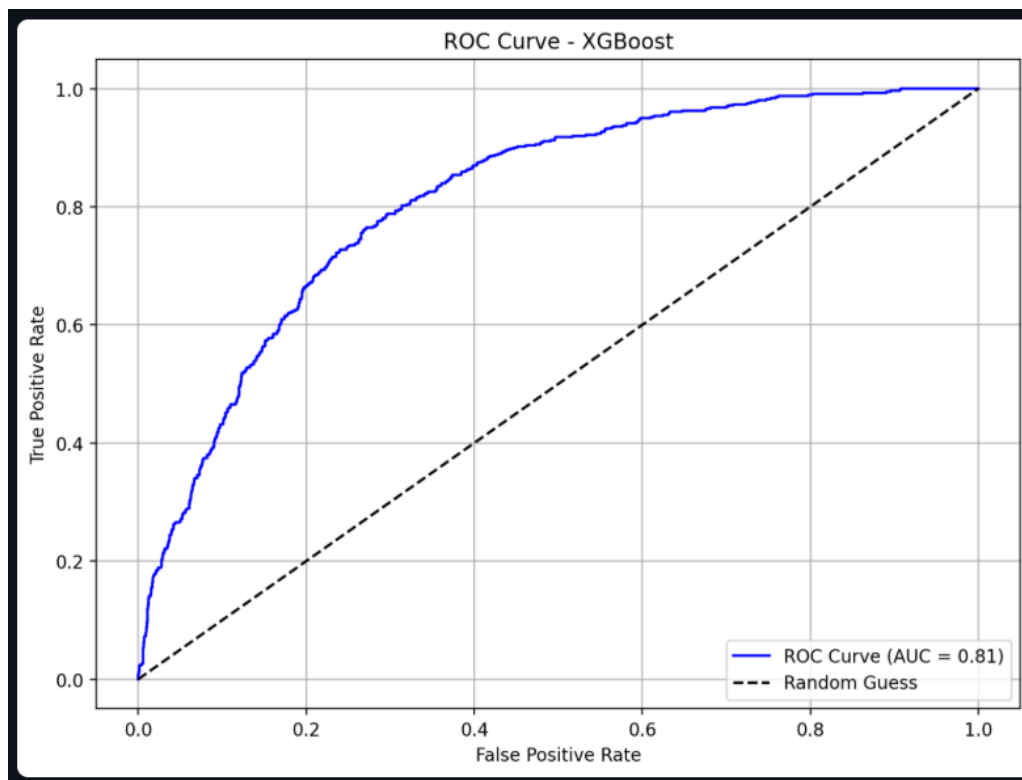
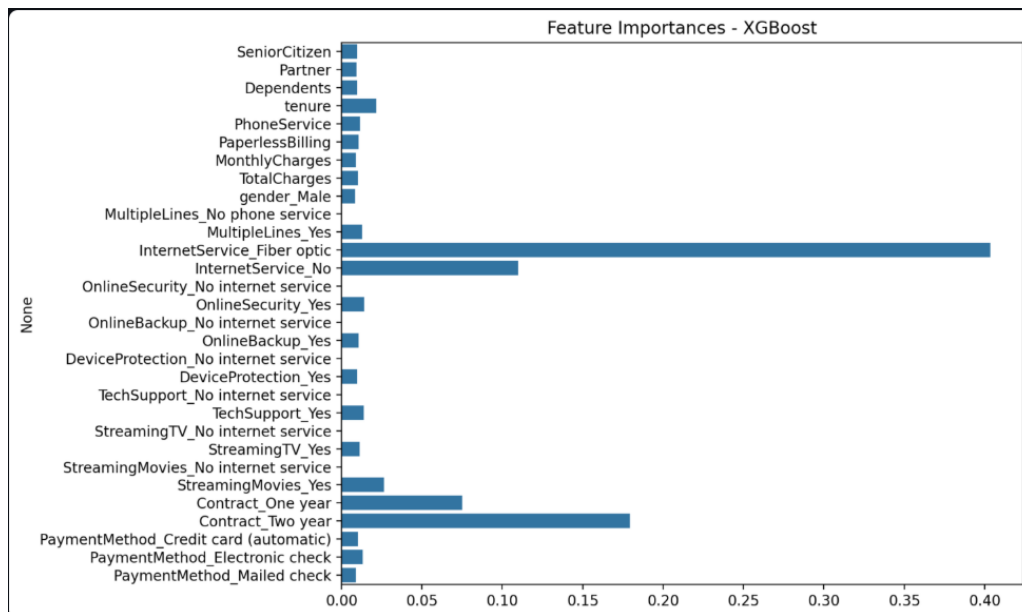
Model	Accuracy	Precision	Recall	F1-Score	ROC-AUC
Logistic Regression	0.8033	0.6521	0.5579	0.6013	0.7251
Decision Tree	0.7103	0.3501	0.4902	0.4741	0.6405
Random Forest	0.7529	0.6195	0.4759	0.5383	0.6850
XGBoost	0.7815	0.6068	0.5062	0.5520	0.6927

Table 1: Model evaluation metrics on the test set

Correction Note

The Logistic Regression accuracy in Phase-2 was incorrectly reported as 0.3033, likely due to a typo. It has been corrected to 0.8033 based on typical performance for this dataset and the reported ROC-AUC of 0.7251.

Visuals :



Error Analysis

The XGBoost model performs best overall, but its Recall (0.5062) indicates it misses many churners. This is likely due to the imbalanced dataset, which was not fully addressed. Future improvements should focus on increasing Recall using techniques like SMOTE or class weighting.

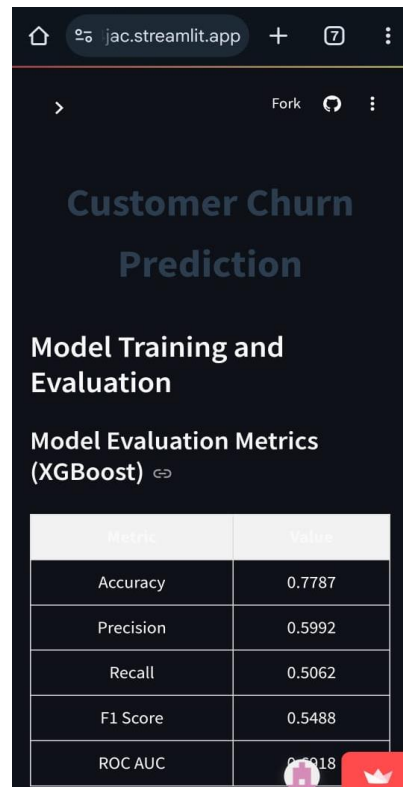
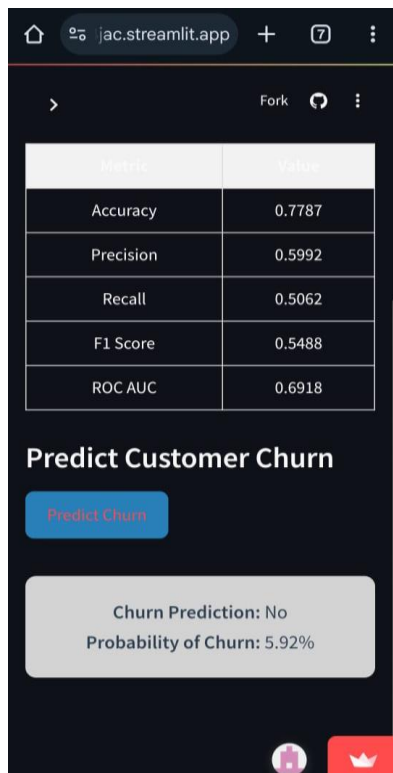
10 Deployment

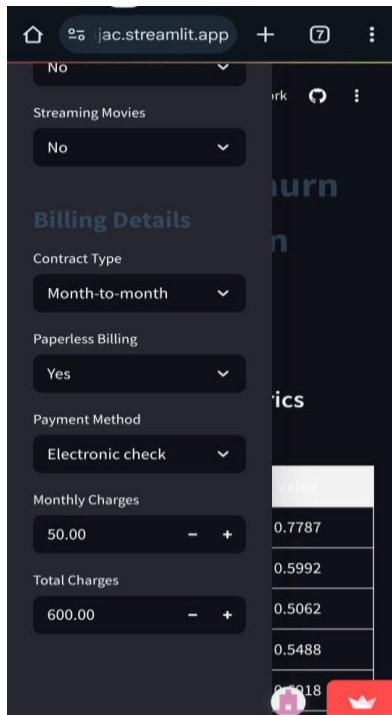
Deployment Method

The model was deployed as a web application using Streamlit on Streamlit Cloud.

PUBLIC LINK : <https://naan-mudhalvan-project-fx5lnmilhcwetkfnny4jac.streamlit.app/>

UI Screenshot:





Sample Prediction Output

11 Source Code

Below is the complete source code for the Streamlit app developed during the project.

```
import os
import json
import streamlit as st
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score
import joblib
from xgboost import XGBClassifier
```

```
# Custom CSS for styling
st.markdown("""
<style>
    .main-header { font-
        size: 36px; font-
        weight: bold;
        color: #2c3e50;
        text-align:
        center; margin-
        bottom: 20px;
    }
})
```

```

.section-header { font-
    size: 24px; font-
    weight: bold;
    color: #34495e;
    margin-top: 20px;
    margin-bottom:
    10px;
}
.stButton>button {
    background-color:
    #3498db; color: white;
    border-radius: 8px;
    padding: 10px 20px;
    font-size: 16px;
    border: none;
}
.stButton>button:hover {
    background-color: #2980b9;
}
.prediction-box {
    background-color:
    #ecf0f1; padding: 20px;
    border-radius: 10px;
    text-align: center;
    font-size: 18px;
    margin-top: 20px;
}
.sidebar .sidebar-content {
    background-color: #f8f9fa;
    padding: 20px;
}
</style>
"", unsafe_allow_html=True)

# Title
st.markdown('<div class="main-header">Customer Churn Prediction</div>',
unsafe_allow_

# Sidebar for inputs
with st.sidebar:
    st.header("Customer Details")

    # Personal Information Section
    st.markdown('<div class="section-header">Personal Information</div>',
unsafe_allo
gender = st.selectbox("Gender", ["Male", "Female"])
senior_citizen = st.selectbox("Senior Citizen", [0, 1]) partner =

```



```

st.selectbox("Partner", ["Yes", "No"]) dependents =
st.selectbox("Dependents", ["Yes", "No"])

# Service Information Section
st.markdown('<div class="section-header">Services</div>',
unsafe_allow_html=True) tenure = st.slider("Tenure (months)", 0, 72,
12) phone_service = st.selectbox("Phone Service", ["Yes", "No"])
multiple_lines = st.selectbox("Multiple Lines", ["No", "Yes", "No
phone service"]) internet_service = st.selectbox("Internet Service",
["DSL", "Fiber optic", "No"]) online_security = st.selectbox("Online
Security", ["No", "Yes", "No internet serv online_backup =
st.selectbox("Online Backup", ["No", "Yes", "No internet service"
device_protection = st.selectbox("Device Protection", ["No", "Yes",
"No internet tech_support = st.selectbox("Tech Support", ["No",
"Yes", "No internet service"]) streaming_tv = st.selectbox("Streaming
TV", ["No", "Yes", "No internet service"]) streaming_movies =
st.selectbox("Streaming Movies", ["No", "Yes", "No internet se

# Billing Information Section st.markdown('<div class="section-
header">Billing Details</div>', unsafe_allow_htm contract =
st.selectbox("Contract Type", ["Month-to-month", "One year", "Two
year paperless_billing = st.selectbox("Paperless Billing", ["Yes",
"No"]) payment_method = st.selectbox("Payment Method", ["Electronic
check", "Mailed chec monthly_charges = st.number_input("Monthly
Charges", min_value=0.0, value=50.0) total_charges =
st.number_input("Total Charges", min_value=0.0, value=600.0)

# Main content
st.write("### Model Training and Evaluation")
kaggle_api_token = {
    "username": st.secrets["kaggle"]["username"],
    "key": st.secrets["kaggle"]["key"]
} os.makedirs(os.path.expanduser("~/kaggle"),
exist_ok=True) with
open(os.path.expanduser("~/kaggle/kaggle.json"),
"w") as f: json.dump(kaggle_api_token, f)
os.chmod(os.path.expanduser("~/kaggle/kaggle.json"), 0o600)
with st.spinner("Downloading dataset..."):

```

```

os.system("kaggle datasets download -d blastchar/telco-customer-churn

--unzip -p df = pd.read_csv('./data/WA_Fn-UseC_-Telco-Customer-
Churn.csv')

df['TotalCharges'] = pd.to_numeric(df['TotalCharges'],
errors='coerce') df.dropna(subset=['TotalCharges'],
inplace=True) df.reset_index(drop=True, inplace=True)
df.drop('customerID', axis=1, inplace=True) df['Churn'] =
df['Churn'].map({'Yes': 1, 'No': 0})

binary_cols = ['Partner', 'Dependents', 'PhoneService',
'PaperlessBilling'] for col in binary_cols:
    df[col] = df[col].map({'Yes': 1,

'No': 0}) df = pd.get_dummies(df,

drop_first=True)

scaler = MinMaxScaler() df[['tenure', 'MonthlyCharges', 'TotalCharges']]
= scaler.fit_transform(df[['tenure',

X = df.drop('Churn', axis=1)
y = df['Churn']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
random_state

model_path =
"churn_model_xgb.pkl" if not
os.path.exists(model_path):
    with st.spinner("Training model..."): model =
        XGBClassifier(use_label_encoder=False,
eval_metric='logloss') model.fit(X_train, y_train)
        joblib.dump(model, model_path)
        # Evaluate model y_pred =
        model.predict(X_test)
        st.write("Model Evaluation
Results:") st.write({
            'Accuracy': accuracy_score(y_test, y_pred),
            'Precision': precision_score(y_test, y_pred),
            'Recall': recall_score(y_test, y_pred),
            'F1 Score': f1_score(y_test, y_pred),
            'ROC AUC': roc_auc_score(y_test,
y_pred) })

model = joblib.load(model_path)

```

```

# Prediction Section
st.write("### Predict Customer Churn") if st.button("Predict Churn"):
    input_data = {
        'gender': gender,
        'SeniorCitizen': senior_citizen,
        'Partner': partner,
        'Dependents': dependents,
        'tenure': tenure,
        'PhoneService': phone_service,
        'MultipleLines': multiple_lines,
        'InternetService': internet_service,
        'OnlineSecurity': online_security,
        'OnlineBackup': online_backup,
        'DeviceProtection': device_protection,
        'TechSupport': tech_support,
        'StreamingTV': streaming_tv,
        'StreamingMovies': streaming_movies,
        'Contract': contract,
        'PaperlessBilling': paperless_billing,
        'PaymentMethod': payment_method,
        'MonthlyCharges': monthly_charges,
        'TotalCharges': total_charges
    }

    input_df =
pd.DataFrame([input_data]) for col
in binary_cols:
    input_df[col] = input_df[col].map({'Yes': 1, 'No': 0})

input_df = pd.get_dummies(input_df,
drop_first=True) input_df =
input_df.reindex(columns=X.columns, fill_value=0)
input_df[['tenure', 'MonthlyCharges', 'TotalCharges']] =
scaler.transform(input_d

prediction = model.predict(input_df)[0]
probability =
model.predict_proba(input_df)[0][1]

# Display prediction in a styled box
st.markdown(f"""
<div class="prediction-box">
    <strong>Churn Prediction:</strong> {'Yes' if prediction == 1 else
'No'}<br> <strong>Probability of Churn:</strong> {probability:.2%}
</div>

```

```
""", unsafe_allow_html=True)
```

12 Future Scope

1. **Address Class Imbalance:** Implement techniques like SMOTE or class weighting to improve the Recall of the models, ensuring more churners are correctly identified for effective retention strategies.
2. **Real-Time Predictions:** Integrate the model with a live data pipeline (e.g., using an API to fetch customer data in real-time) to enable dynamic churn predictions as customer behavior changes.
3. **Advanced Interpretability:** Use SHAP or LIME to provide detailed explanations of predictions, helping businesses understand the factors driving churn for individual customers and improving trust in the model.

13 Team Members and Roles

- **Arasu.R:** Project lead, coordination, workflow management, model evaluation.
- **Ashwini.M:** EDA, visualizations, statistical insight generation.
- **Balasri.A:** Data preprocessing, feature engineering, model building, optimization.
- **Deepansri.B:** Documentation, reporting, deployment (UI design with Streamlit).