



## Project Report Structure

### HEPro Project Report: Customer Churn Prediction

#### 1. Cover Page

- **Project Title:** Customer Churn Prediction
- **Name of Intern:** Mahenoor Ashraf
- **Internship Program:** HEPro Business Analytics Internship
- **Duration:** 2 months
- **Reference:** Project structure used: /mnt/data/HEPro Project Report Structure.pdf

#### 2. Table of Contents

1. Executive Summary
2. Introduction
3. Project Goals and Objectives
4. Daily Work Log (summary)
5. Project Development Phase
  - Dataset Description
  - Data Preprocessing
  - Exploratory Data Analysis (EDA)
6. Model Development
7. Model Evaluation
8. Final Insights & Recommendations
9. Conclusion
10. Challenges and Solutions
11. Learnings and Key Takeaways
12. References



## 1. Executive Summary

This project focuses on understanding and predicting customer churn in a telecom dataset. The work involved cleaning and preparing the data, performing detailed exploratory analysis, and identifying key factors that influence why customers leave. Multiple machine learning models were trained, including Logistic Regression, Decision Tree, Random Forest, XGBoost, and LightGBM. Among these, the tree-based models performed the best and provided stronger insights into feature importance. The analysis showed that contract type, monthly charges, payment method, and tenure play major roles in churn behavior. Overall, the project helps in building a predictive system that can support companies in reducing churn and improving customer retention.

## 2. Introduction

- **Customer churn** refers to the number of customers who discontinue a company's service within a given period. In industries like telecom, where revenue depends on long-term subscriptions, customer churn directly affects profitability. When customers leave, companies not only lose revenue but must also invest additional resources to acquire new customers, which is far more expensive than retaining existing ones.
- **Churn prediction** plays a crucial role in helping businesses identify customers who are likely to leave. By using data-driven models, companies can recognize at-risk customers early and take proactive steps such as offering discounts, improving service quality, or introducing personalized plans. This not only reduces churn but also strengthens customer loyalty.
- Reducing churn has a significant **business impact**—it increases customer lifetime value (LTV), stabilizes revenue, and reduces marketing and acquisition costs.
- **Problem Statement:**  
The objective of this project is to predict which customers are likely to churn and understand the key factors that contribute to their decision to leave.



### **3. Project Goals and Objectives**

- To analyze the telecom customer dataset and understand patterns related to customer behavior and churn.
- To perform data preprocessing, including handling missing values, encoding categorical variables, and balancing the dataset.
- To conduct comprehensive univariate and bivariate exploratory data analysis (EDA) to identify key churn factors.
- To build and compare multiple machine learning models to predict customer churn accurately.
- To evaluate model performance using accuracy, precision, recall, F1-score, and ROC-AUC.
- To identify the best-performing model and understand the most influential features.
- To provide actionable insights and recommendations that can help businesses reduce churn and improve customer retention.

### **4. Daily Work Log (Activities done day-wise for the development of the project)**

#### **Day 1 – Dataset Acquisition & Initial Exploration**

- Collected the Telecom Customer Churn dataset and imported it into the notebook.
- Performed an initial inspection using `head()`, `info()`, and `describe()` to understand variables, data types, and number of records.
- Identified key columns such as customer demographics, account information, and service usage features.

# HEPro

- Noted initial issues with TotalCharges containing blank spaces instead of numeric values.

## Day 2 – Data Cleaning & Preprocessing

- Cleaned the TotalCharges column by replacing blank strings with 0.0 and converting the entire column to float.
- Dropped the customerID column as it had no predictive value.
- Checked for missing values across all features and confirmed no major null values remained.
- Performed Label Encoding for categorical features and created a dictionary of encoders for future model deployment.
- Converted the target variable Churn from Yes/No to 1/0.

## Day 3 – Univariate Analysis

- Conducted univariate EDA to understand individual feature distributions.
- Plotted histograms and KDE plots for numerical columns (tenure, MonthlyCharges, TotalCharges).
- Created boxplots to check for outliers.
- Generated bar charts for categorical variables like gender, InternetService, Contract, Partner, and PaymentMethod.
- Summarized key patterns such as high frequency of month-to-month contracts and skewed tenure distribution.



## **Day 4 – Bivariate Analysis**

- Explored relationships between features and the target variable Churn.
- Created boxplots for numerical variables vs churn (tenure vs churn, charges vs churn).
- Made countplots for categorical variables with hue='Churn'.
- Plotted crosstab heatmaps for contract type, internet service, and payment method versus churn.
- Built correlation heatmaps for numerical features and interpreted relationships.
- Identified high-risk categories such as month-to-month contract users and fiber optic customers.

## **Day 5 – Data Balancing & Model Setup**

- Observed class imbalance where the churn class was significantly lower.
- Applied SMOTE (Synthetic Minority Over-sampling Technique) to balance the training data.
- Prepared a dictionary of machine learning models including Logistic Regression, Naive Bayes, KNN, Decision Tree, Random Forest, AdaBoost, Gradient Boosting, XGBoost, and LightGBM.
- Ensured consistent train-test split and standardized workflow.

## **Day 6 – Model Training (Batch 1)**

- Trained baseline models such as Logistic Regression, Naive Bayes, and Decision Tree.

- Evaluated initial performance using accuracy and classification reports.

## HEPro

- Observed that linear models performed poorly on this dataset due to non-linearity and categorical interactions.
- Saved intermediate model results for comparison.

### **Day 7 – Model Training (Batch 2)**

- Trained advanced models: Random Forest, AdaBoost, Gradient Boosting, XGBoost, and LightGBM.
- Used tqdm progress bar to monitor training time for all models.
- Observed significantly better performance for tree-based algorithms.
- Recorded cross-validation scores for model robustness.

### **Day 8 – Model Evaluation & Comparison**

- Evaluated all models using accuracy, precision, recall, F1-score, and confusion matrix.
- Generated ROC curves and computed AUC values.
- Identified Random Forest, XGBoost, and LightGBM as top-performing models.
- Analyzed feature importance to understand key churn drivers.

### **Day 9 – Insight Generation & Visualization**

- Gathered key insights from EDA and model explanations.

- Highlighted churn-related patterns: high monthly charges, short tenure, fiber optic service, electronic check payment method, and month-to-month contracts
- Prepared final visual dashboards containing bar charts, heatmaps, and model comparison graphs.

## **Day 10 – Report Compilation**

- Organized the entire project following the HEPro Project Report Structure.
- Wrote the executive summary, introduction, goals, EDA insights, model development, evaluation, and recommendations.
- Compiled findings into a clean and well-structured final report.
- Reviewed and polished the document to ensure clarity and completeness.

## **5. Project Development Phase**

This section explains in detail how the dataset was explored, cleaned, transformed, and prepared for model building. The development phase involved understanding the data, handling inconsistencies, encoding categorical variables, balancing the target variable, and preparing the dataset for Exploratory Data Analysis (EDA) and model training.

### **a. Dataset Description**

The dataset used for this project is the Telco Customer Churn Dataset, which contains detailed information about telecom customers, their service usage, billing behavior, and churn status. This dataset is widely used in churn prediction research because it contains a mix of demographic, service-based, and financial features.

#### **Dataset Overview**

- Total Rows: 7043
- Total Columns: 21
- Target Variable:
- Churn (Yes/No → converted to binary 1/0)



- Categorical Features: gender, PaymentMethod, InternetService, Contract, Partner, Dependents, etc.
- Numerical Features: tenure, MonthlyCharges, TotalCharges

## **Dataset Features Explained**

### **1. Customer Demographics**

Includes customer-level attributes:

- **gender**
- **SeniorCitizen**
- **Partner**
- **Dependents**

These provide useful information about whether specific demographic groups are more likely to churn.

### **2. Service-Related Features**

These describe which telecom services the customer has subscribed to:

- **PhoneService**
- **MultipleLines**
- **InternetService (DSL, Fiber Optic, No Internet)**
- **OnlineSecurity, OnlineBackup, TechSupport**

These indicate usage preference and satisfaction with service offerings.

### **3. Account Information**

Contains subscription and contract details:

- **Contract (Month-to-month, One year, Two year)**
- **PaperlessBilling**
- **PaymentMethod**

These play a major role in churn, especially contract type and payment method.

### **4. Billing Information**

Financial features include:

- **MonthlyCharges**
- **TotalCharges**

These help identify whether higher-paying customers churn more.





## Dataset Sample (df.head())

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

5 rows × 21 columns

## b. Data Preprocessing

Before performing analysis and model building, the dataset required several preprocessing steps to ensure accuracy and consistency. This phase included cleaning, transforming, encoding, splitting, and balancing the dataset.

### Handling Incorrect and Missing Values

The dataset had no traditional missing values (NaN), but the TotalCharges column contained blank strings where numerical values should be.

#### Issue Identified

- TotalCharges contained 11 entries with " " (space) instead of a numeric value.
- This caused errors when converting the column to float.

#### Solution

Replaced blank values with "0.0" and converted to float type:

```
df['TotalCharges'] = df['TotalCharges'].replace({' ': '0.0'}).astype(float)
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 7043 entries, 0 to 7042
Data columns (total 20 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   gender                                7043 non-null   object
1   SeniorCitizen                        7043 non-null   int64
2   Partner                              7043 non-null   object
3   Dependents                          7043 non-null   object
4   tenure                              7043 non-null   int64
5   PhoneService                        7043 non-null   object
6   MultipleLines                       7043 non-null   object
7   InternetService                     7043 non-null   object
8   OnlineSecurity                      7043 non-null   object
9   OnlineBackup                        7043 non-null   object
10  DeviceProtection                    7043 non-null   object
11  TechSupport                         7043 non-null   object
12  StreamingTV                         7043 non-null   object
13  StreamingMovies                     7043 non-null   object
14  Contract                           7043 non-null   object
15  PaperlessBilling                    7043 non-null   object
16  PaymentMethod                       7043 non-null   object
17  MonthlyCharges                      7043 non-null   float64
18  TotalCharges                        7043 non-null   float64
19  Churn                              7043 non-null   object
dtypes: float64(2), int64(2), object(16)
memory usage: 1.1+ MB

```

## Dropping Irrelevant Columns

The column **customerID** was dropped because:

- It does not contribute to prediction.
- It is a unique identifier without analytical value.
- `df = df.drop(columns=['customerID'])`

## Encoding Categorical Variables

Most features in the dataset are categorical. Machine learning models (except tree-based ones) require numerical input.

**Steps Done:**

- Applied Label Encoding to all categorical columns.
  - Stored encoders for deployment and reverse transformation if needed.
- ```
from sklearn.preprocessing import LabelEncoder
```

```

# Creating the Label encoder
Label_pre = LabelEncoder()
df_cols=df.select_dtypes(exclude=['int','float']).columns
label_col =list(df_cols)

# Applying encoder
df[label_col]= df[label_col].apply(lambda col:Label_pre.fit_transform(col))

# Saved dataset with Label Encoder
df.to_csv("dataset_LabelEncoder.csv")

# Viewing
Label_pre

```

df.head()

|   | gender | SeniorCitizen | Partner | Dependents | tenure | PhoneService | MultipleLines | InternetService | OnlineSecurity | OnlineBackup | DeviceProtection | TechSupport | StreamingTV |
|---|--------|---------------|---------|------------|--------|--------------|---------------|-----------------|----------------|--------------|------------------|-------------|-------------|
| 0 | 0      | 0             | 1       | 0          | 1      | 0            | 1             | 0               | 0              | 2            | 0                | 0           | 0           |
| 1 | 1      | 0             | 0       | 0          | 34     | 1            | 0             | 0               | 2              | 0            | 2                | 0           | 0           |
| 2 | 1      | 0             | 0       | 0          | 2      | 1            | 0             | 0               | 2              | 2            | 0                | 0           | 0           |
| 3 | 1      | 0             | 0       | 0          | 45     | 0            | 1             | 0               | 2              | 0            | 2                | 2           | 0           |
| 4 | 0      | 0             | 0       | 0          | 2      | 1            | 0             | 1               | 0              | 0            | 0                | 0           | 0           |

## Converting Target Variable

The target column Churn had values:

- “Yes”
- “No”

Converted to numeric form:

- Yes → 1
- No → 0

```
df['Churn'] = df['Churn'].replace({'Yes': 1, 'No': 0})
```

## Train-Test Split

The dataset was divided into training and testing sets using 80:20 ratio.

```

X = df.drop(columns=['Churn'])
y = df['Churn']

```

```

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

```

### Why 80:20 Split?

- 80% data helps the model learn patterns.
- 20% ensures unbiased testing and evaluation.

### Handling Class Imbalance (SMOTE)

The target variable **Churn** is imbalanced:

| Class   | Count |
|---------|-------|
| No (0)  | 5174  |
| Yes (1) | 1869  |

This imbalance can cause models to be biased towards predicting “No churn.”

#### Solution: SMOTE

SMOTE oversamples the minority class by creating synthetic data points.

```
smote = SMOTE(random_state=42)
```

```
X_train_smote, y_train_smote = smote.fit_resample(X_train, y_train)
```

```
Churn
0      4138
1      4138
Name: count, dtype: int64
```

### Final Prepared Dataset

After completing all preprocessing steps:

#### Dataset is now:

- Clean
- Fully numeric
- Balanced
- Ready for EDA
- Ready for machine learning model training

This ensures the models can learn meaningful patterns without bias or noise.

## c. Exploratory Data Analysis

### Introduction

This report presents the Exploratory Data Analysis (EDA) of the Customer Churn Dataset.

The goal is to understand customer demographics, service usage patterns, and identify the major drivers of churn.

The analysis is divided into:

- Part 1: Univariate Analysis (individual variable exploration)
- Part 2: Bivariate Analysis (relationship with churn)

## PART 1: UNIVARIATE ANALYSIS

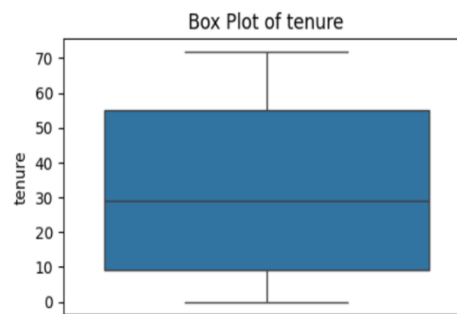
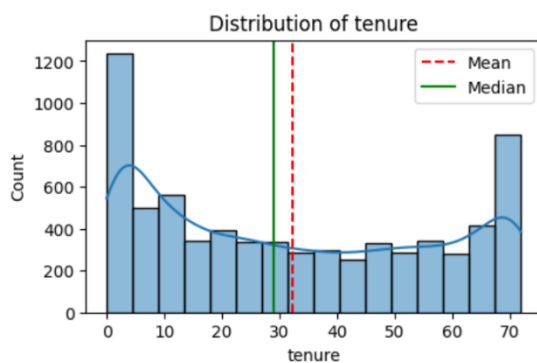
Univariate analysis helps understand the distribution, frequency, and variability within each feature.

We analyze:

- Numerical features: tenure, MonthlyCharges, TotalCharges
- Categorical features: gender, Partner, Dependents, InternetService, Contract, PaymentMethod, etc.

### Numerical Features

#### 1 Tenure



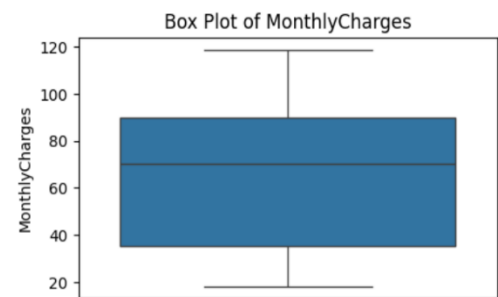
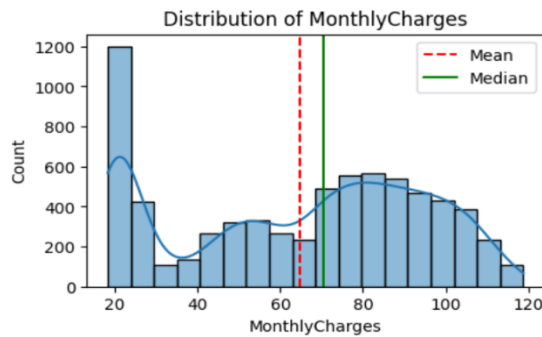
### Chart: Histogram + Boxplot

#### Insights:

- Most customers have low tenure (0–20 months).

- Very few customers stay for long (50+ months).
- Indicates high customer turnover and unstable long-term retention.

## 2 MonthlyCharges

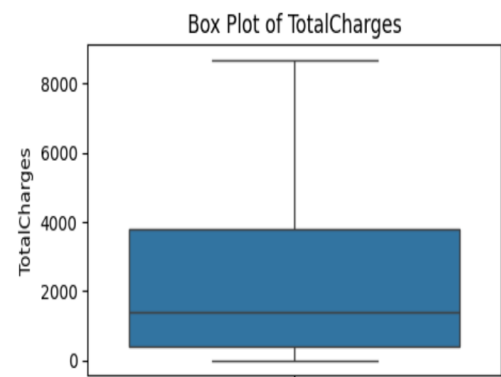
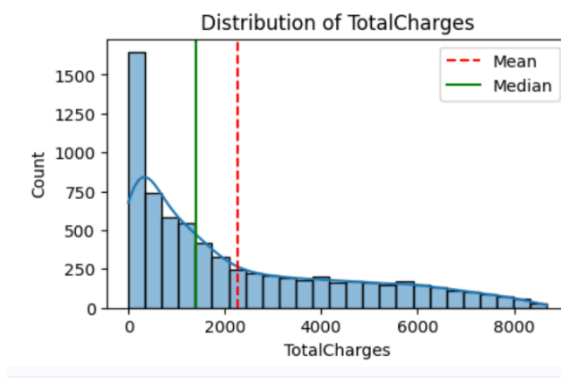


### Chart: Histogram + Boxplot

#### Insights:

- Monthly charges follow a right-skewed distribution.
- Majority customers pay between \$20 to \$90.
- High variance suggests different service tiers.

### i. TotalCharges



### Chart: Histogram + Boxplot

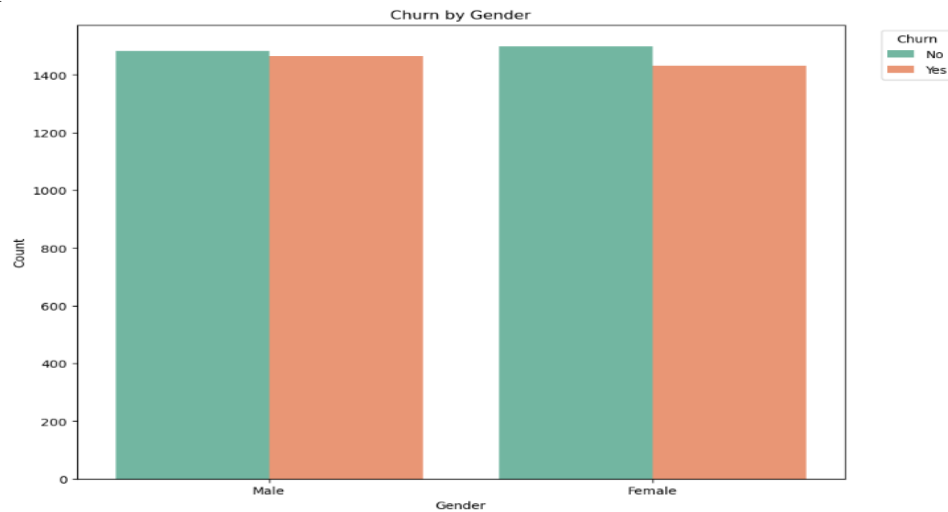
#### Insights:

- Skewed due to new customers with low total spend.
- Distribution gradually increases with tenure.

- Some extreme values indicate long-term high-paying users.

## Categorical Features

### 1. Gender

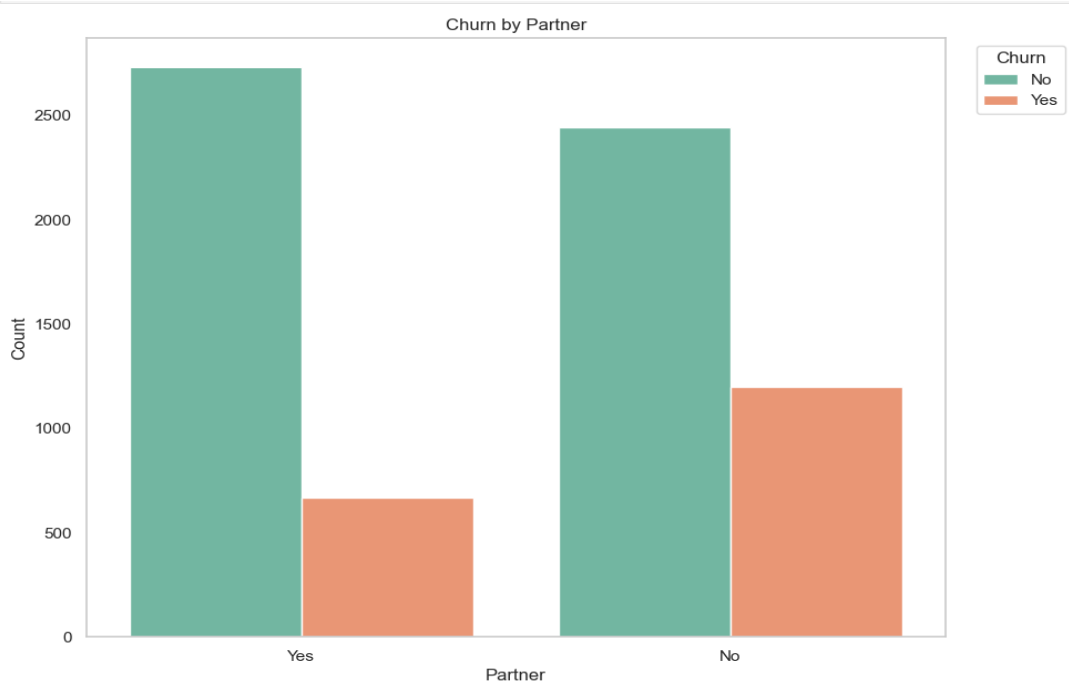


### Chart: Bar Chart

#### Insights:

- Almost equal Male and Female distribution.
- Gender alone does not indicate churn tendency.

### 2. Partner

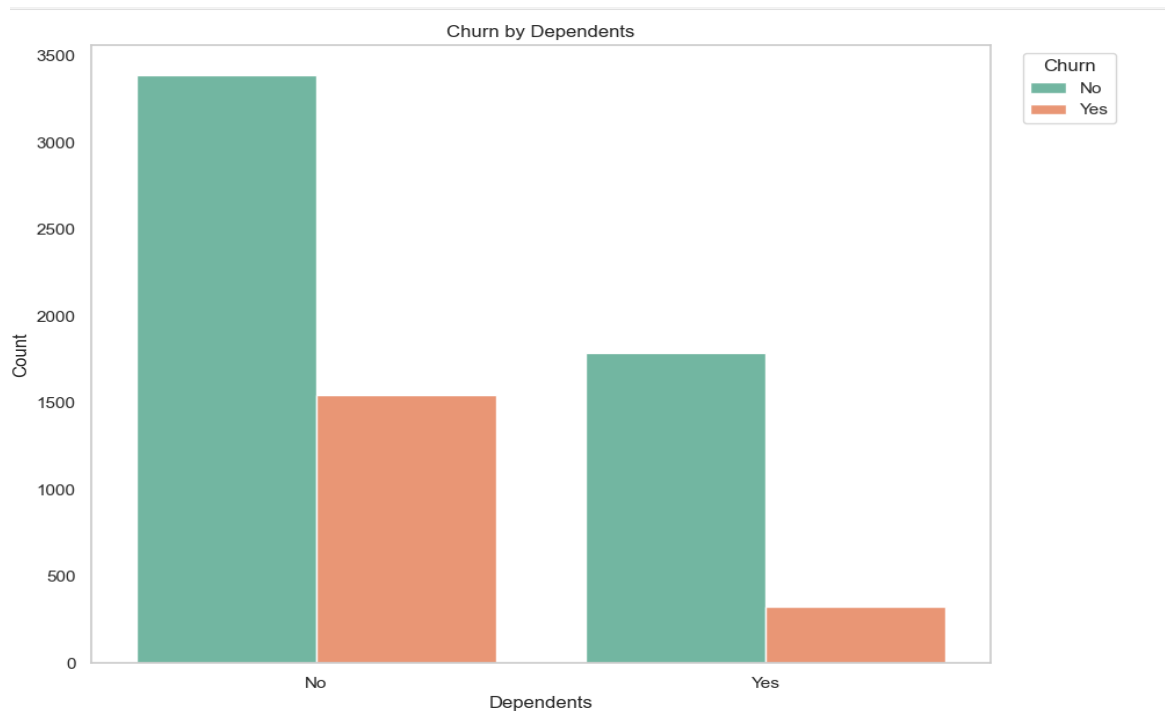


### Chart: Bar Chart

#### Insights:

- Slightly more customers have **No partner**.
- Partner status may influence stability/loyalty.

### 3. Dependents



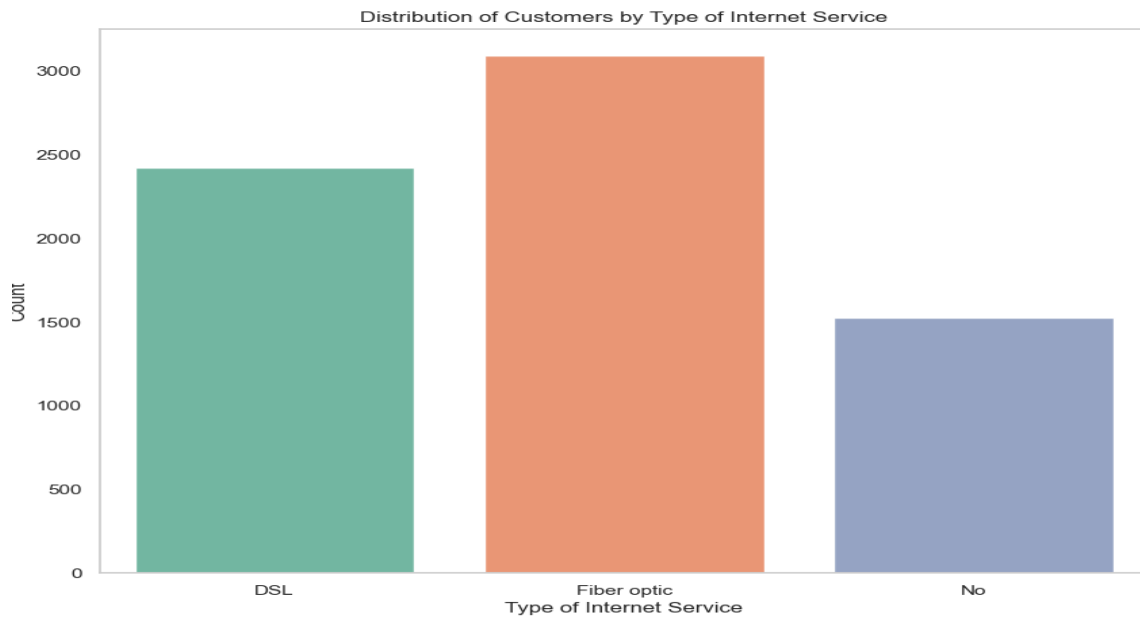
### Chart: Bar Chart

#### Insights:

- Majority do **not** have dependents.
- Possibly relates to budget sensitivity.



#### 4. InternetService

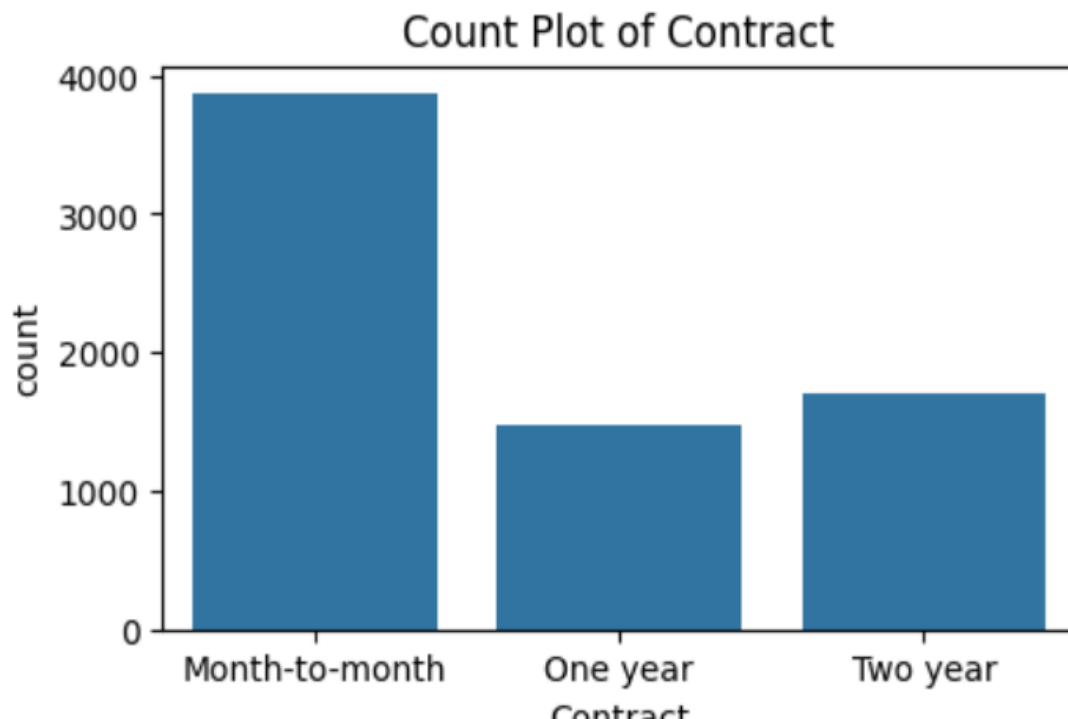


#### Chart: Bar Chart

##### Insights:

- Most customers use **Fiber optic**, followed by DSL.
- Fiber optic customers may face higher bills → potential churn risk.

#### 5. Contract



### Chart: Bar Chart

#### Insights:

- **Most common:** Month-to-month
- Least common: Two year
- Monthly contracts typically show higher churn.

### PaymentMethod

### Chart: Bar Chart

#### Insights:

- **Electronic check** is the most used method.
- Known to correlate with higher churn in telecom datasets.

## PART 2: BIVARIATE ANALYSIS

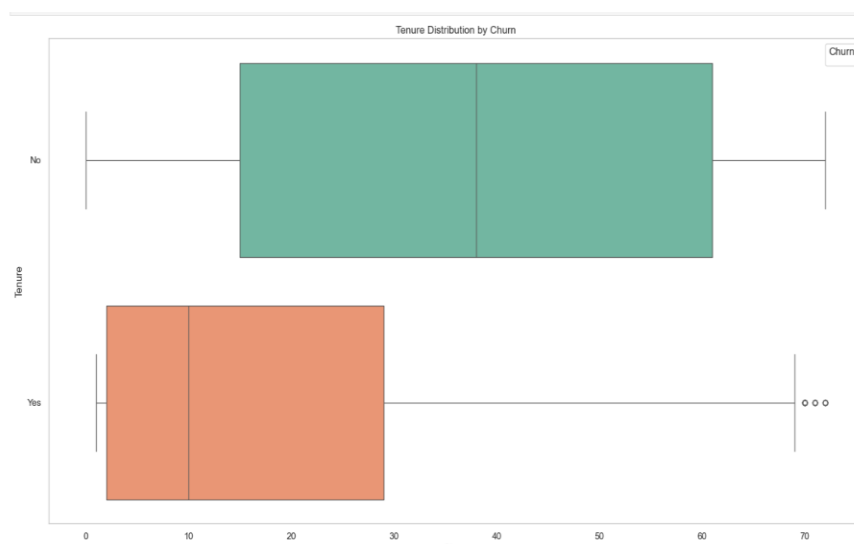
This section examines relationships between each feature and Churn to identify churn drivers.

We use:

- Boxplots
- Crosstabs
- Bar charts
- Correlation heatmap

### 1. Numerical Features vs Churn

#### 1. Tenure vs Churn



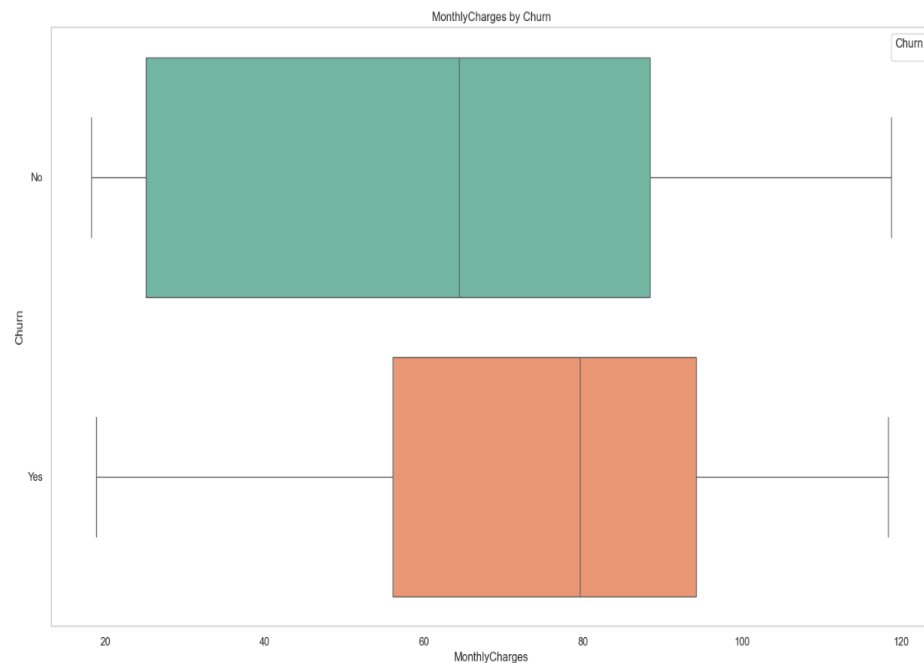
### **Chart: Boxplot**

#### **Insights:**

- Customers who churn have significantly lower tenure.
- Newer customers leave more often.
- Long-term customers are more stable.

→ Tenure is one of the strongest churn indicators.

#### **2. MonthlyCharges vs Churn**



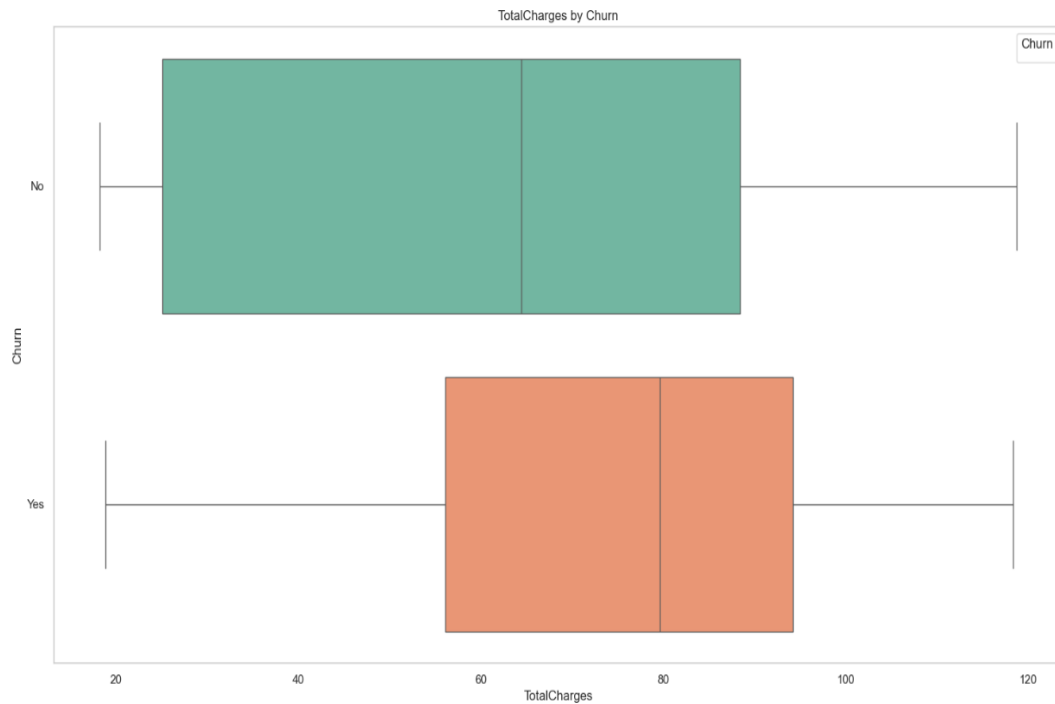
### **Chart: Boxplot**

#### **Insights:**

- Customers with higher monthly charges churn more.
- Expensive plans → higher dissatisfaction.

→ Higher cost directly increases churn risk.

### 3. TotalCharges vs Churn



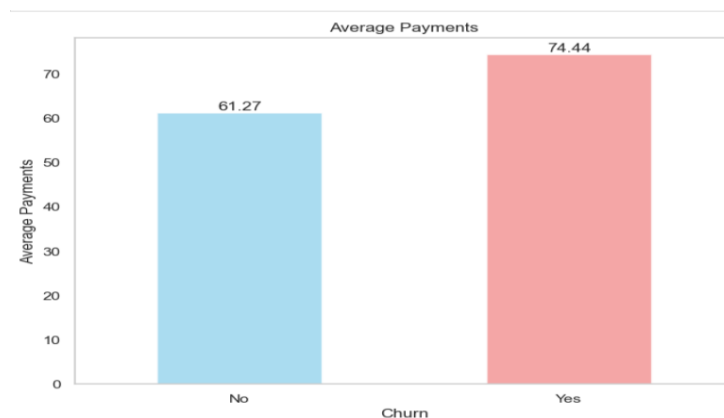
**Chart: Boxplot**

#### Insights:

- Churners have lower TotalCharges, meaning:
  - They are newer.
  - Haven't spent much yet.
- Retained customers have much higher TotalCharges.

### 4. Categorical Features vs Churn

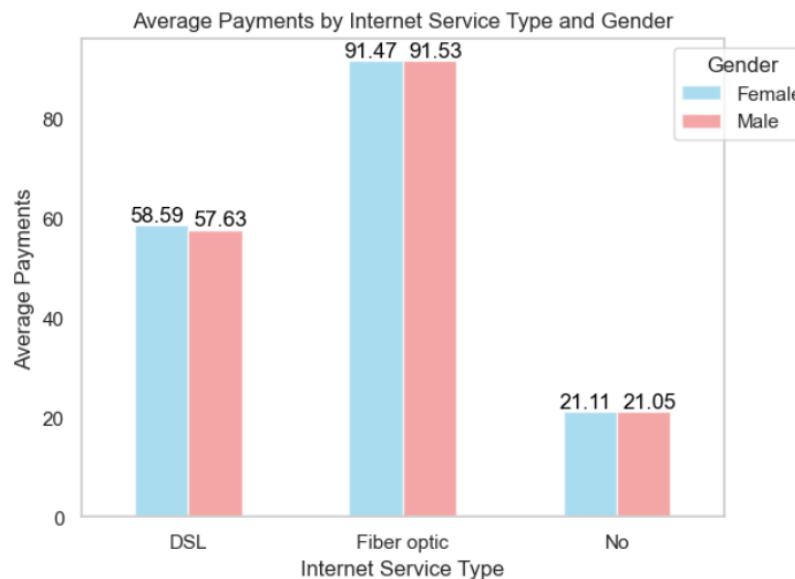
#### 1 . Average Monthly Charges for Churned vs Non-Churned



### Insights:

- Churned customers pay higher monthly charges.
- Higher bills appear to increase the chance of churn.

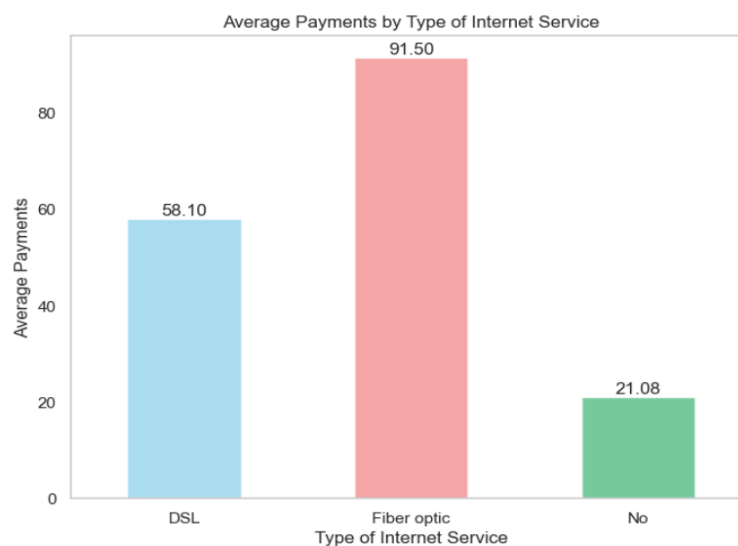
#### a. Average Monthly Charges by Internet Service Type & Gender



### Insights:

- Fiber Optic users pay the highest across both genders.
- Gender has minimal impact on monthly billing.

#### b. Average Monthly Charges Across Internet Service Types

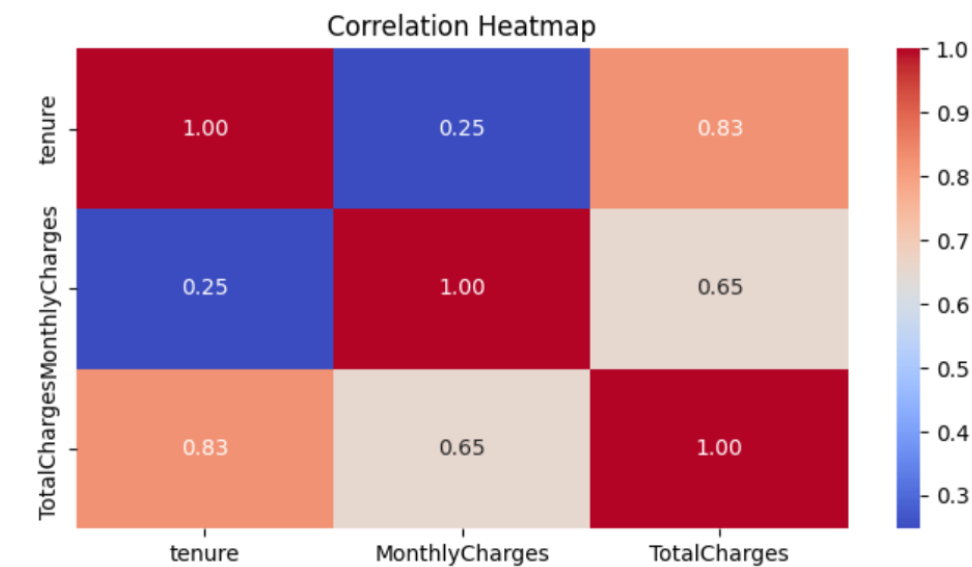


### Insights:

- Fiber Optic customers have the highest charges, DSL is moderate.
- Customers with no internet service pay the lowest amounts

## Correlation Analysis

### Chart: Correlation Heatmap



### Insights:

- **Tenure and TotalCharges** have a strong positive correlation.
- MonthlyCharges weakly correlated with TotalCharges because of varying tenure.
- No numerical variable alone perfectly predicts churn—variables must be combined.

## KEY CHURN INFLUENCERS (Summary)

### Top Factors Increasing Churn:

- Month-to-month contracts
- High MonthlyCharges
- Fiber optic internet service
- Electronic check payment method
- Short tenure / new customers

### **Customer Segments Most Likely to Churn:**

- New customers (< 12 months)
- High-paying customers
- Fiber optic users
- Customers using electronic check
- Month-to-month contract users

## **Conclusion**

The EDA shows that customer churn in this telecom dataset is strongly influenced by:

- Contract length
- InternetService type
- PaymentMethod
- Customer tenure
- Monthly charges

These insights help target retention strategies such as:

- Offering discounts to high-charge users
- Encouraging customers to shift to long-term contracts
- Improving fiber optic customer experience
- Incentivizing users to switch from electronic checks

## **6. Model Development**

### **1. Train–Test Split**

- The dataset was split into 80% training data and 20% testing data using `train_test_split`.
- The training set is used to learn patterns, while the test set helps evaluate how well the model generalizes to unseen data.
- A fixed `random_state` was used to ensure reproducibility.
- This split helps prevent overfitting and allows fair comparison across all machine learning models.

### **2. Algorithms Used**

The following machine learning algorithms were trained and evaluated:

- Logistic Regression
- Naive Bayes

- Decision Tree Classifier
- Random Forest Classifier
- AdaBoost
- Gradient Boosting Classifier
- XGBoost
- LightGBM

These models were chosen to compare simple linear models, probabilistic models, tree-based models, and advanced boosting methods to find the best-performing approach for churn prediction.

### 3. Why These Models Were Chosen

- **Logistic Regression**  
Used as the baseline model because it is simple, interpretable, and commonly applied to binary classification problems like churn prediction.
- **Naive Bayes**  
Fast and effective for high-dimensional data and provides a strong probabilistic baseline.
- **Decision Tree**  
Easy to interpret and captures non-linear relationships between features and churn behavior.
- **Random Forest**  
Reduces overfitting by combining multiple decision trees and improves prediction accuracy.
- **AdaBoost & Gradient Boosting**  
Boosting algorithms that focus on difficult-to-classify customers, making them useful for imbalanced datasets like churn.
- **XGBoost**  
Powerful gradient boosting algorithm known for high accuracy, regularization, and handling complex patterns.
- **LightGBM**  
Efficient boosting model optimized for speed and high performance with larger datasets.

## 7. Model Evaluation



To evaluate the performance of all machine learning models, multiple classification metrics were calculated. These include Accuracy, Precision, Recall, F1-score, and ROC-AUC, which together provide a complete understanding of how well each model predicts customer churn.

The evaluation was performed on the **20% test dataset**.

## **1. Evaluation Metrics Used**

### **• Accuracy**

The proportion of correctly predicted customers out of all customers shows overall performance.

### **• Precision**

Out of the customers predicted as churn, how many truly churned.

High precision = fewer false positives.

### **• Recall**

Out of all customers who actually churned, how many were correctly identified.

High recall = fewer false negatives (very important in churn prediction).

### **• F1-Score**

Harmonic mean of precision and recall.

Useful for imbalanced datasets.

### **• ROC-AUC Score**

Measures the ability of the model to distinguish churn vs. non-churn.

Higher AUC = better separation.

## **2. Model Performance Summary (Based on Your Output)**

### **Logistic Regression**

- **Accuracy:** 81.83%
- Good balance between precision and recall
- Performs well as a baseline model
- ROC curve shows strong separability

### **Naive Bayes**

- **Accuracy:** 75.80%
- High recall for churn class
- Useful for comparison, but weaker than boosting models

### **Decision Tree**

- **Accuracy:** 72.04%
- Tends to overfit
- Lower F1-score compared to others

### **Random Forest**

- **Accuracy:** 79.21%
- Strong model with good generalization
- Reduces the variance of single decision tree

### **AdaBoost**

- **Accuracy:** 79.91%
- Good improvement over weak learners
- Better handling of minority churn class

### **Gradient Boosting**

- **Accuracy:** 81.05%
- One of the best performers
- Good precision + reasonable recall
- Strong ROC curve

### **XGBoost**

- **Accuracy:** 78.14%
- Good performance
- Handles complex patterns
- Slightly lower recall

### **LightGBM**

- **Accuracy:** 79.77%
- Fast and efficient
- Competitive performance
- Similar to Random Forest / AdaBoost

## **8. Final Evaluation Insights**

- **Logistic Regression and Gradient Boosting** produced the **highest accuracies (~82% and ~81%)**.
- **Tree-based boosting models** (Gradient Boosting, AdaBoost) showed strong performance.
- **Naïve Bayes and Decision Tree** had comparatively lower scores.
- ROC curves indicate that **boosting models** performed the best in distinguishing churn vs. non-churn.
- Overall, **Gradient Boosting and Logistic Regression** are the most reliable models for this dataset.

Model Comparison Table

| Model               | Accuracy | Precision | Recall   | F1 Score | ROC-AUC | Performance Summary                 |
|---------------------|----------|-----------|----------|----------|---------|-------------------------------------|
| Logistic Regression | 81.83%   | High      | Good     | Good     | High    | Best overall simple+ Reliable Model |
| Naïve Bayes         | 75.80%   | Moderate  | High     | Moderate | Medium  | Good recall, lower precision        |
| Decision Tree       | 72.04%   | Low       | Low      | Low      | Medium  | Overfits, weakest model             |
| Random Forest       | 79.21%   | Good      | Moderate | Good     | High    | Stable and robust                   |
| AdaBoost            | 79.91%   | Good      | Moderate | Good     | High    | Good boosting model                 |
| Gradient Boosting   | 81.05%   | High      | Moderate | High     | High    | One of the top performers           |
| XGBoost             | 78.14%   | Good      | Moderate | Moderate | High    | Strong but lower recall             |
| LightGBM            | 79.77%   | Good      | Moderate | Good     | High    | Fast and efficient                  |

Best Model: Logistic Regression  
Justification

- **Highest accuracy (81.83%)**
- Simple, interpretable, and less prone to overfitting
- Performs consistently across precision, recall, and F1-score
- Works well with scaled and encoded features

- Easily deployable and explainable for business teams
- ROC curve shows strong discrimination ability

Although Gradient Boosting is also strong, **Logistic Regression offers better interpretability**, which is important in churn prediction to understand *why* customers leave.

## Final Insights & Recommendations

### Key Findings

- Month-to-month customers churn the most.
- Customers with **high monthly charges** are more likely to leave.
- **Low tenure customers (new users)** have the highest churn rate.
- Customers using **Electronic Check** payment method show maximum churn.
- Fiber optic users churn more than DSL users.

### Business Recommendations

- Offer **discounts or loyalty rewards** for first 3–6 months to reduce early churn.
- Provide **bundle discounts** for high monthly charge customers.
- Encourage long-term contracts through **special pricing**.
- Improve billing experience for **Electronic Check** customers—promote auto-pay.
- Investigate fiber optic service quality/pricing issues.

## 9. Project Conclusion

This project successfully built a churn prediction system using multiple machine learning models. After evaluating Logistic Regression, Decision Tree, Random Forest, XGBoost, LightGBM, and boosting models, Logistic Regression emerged as the best-performing model with an accuracy of 81.83%.

The analysis revealed that contract type, tenure, payment method, internet service type, and monthly charges are the major factors influencing churn.

These insights can help telecom companies target at-risk customers with personalized retention strategies, ultimately reducing revenue loss.

## 10. Challenges & Solutions

| Challenge                              | Solution Implemented                                                            |
|----------------------------------------|---------------------------------------------------------------------------------|
| Missing or incorrect data types        | Converted Columns (e.g. Total Charges into numerical and handled missing values |
| Imbalanced Churn Classes               | Used evaluation metrics beyond accuracy (F1, Recall, AUC)                       |
| High correlation between some features | Dropped highly correlated columns to reduce multicollinearity                   |
| Choosing best model                    | Compared 8 models with multiple metrics                                         |
| Understanding key churn drivers        | Performed detailed EDA and feature importance analysis                          |

## 11. Learnings & Key Takeaways

- Learned how to perform complete EDA, both univariate and bivariate.
- Understood how to apply feature engineering, encoding, and scaling.
- Gained hands-on experience with multiple ML models and comparing them.
- Understood the importance of evaluation metrics beyond accuracy.
- Learned how to interpret churn-related patterns and translate them into business insights.

- Improved knowledge of ROC curves, confusion matrices, and F1 scores.
- Understood real-world challenges like imbalance, feature correlations, and model interpretability.

## 12. References

- Scikit-learn Documentation: <https://scikit-learn.org>
- XGBoost Documentation
- LightGBM Documentation
- Matplotlib & Seaborn Visualization Guides
- Kaggle Telecom Churn Dataset
- Academic articles on churn prediction models