**AMRITA SCHOOL OF ARTIFICIAL INTELLIGENCE**

**COIMBATORE 641112**

**April 2025**

## 23AIE232M PYTHON FOR AI

## Predicting Customer Churn in Telecom Using Machine Learning
### END SEMESTER REVIEW REPORT

*Submitted by*

| | |
|---|---|
| CB.EN.U4ECE23022 | **Madhavh SRK** |
| CB.EN.U4ECE23025 | **Pedada Venkata Abhiram** |
| CB.EN.U4ECE23038 | **S Muthu Vasanth** |
| CB.EN.U4ECE23055 | **Vanka Ravikrishna** |
| CB.EN.U4ECE23056 | **Yuvanesh M** |

**Department of Electronics and Communication Engineering**

**Date of Submission**

**(24/04/2025)**

Marks                                                                 Faculty In-charge

# Table of Contents

# Abstract

This project aims to develop a predictive system for identifying customer churn in the telecom industry using machine learning. Churn, the process where customers discontinue their services, is a significant challenge for telecom providers due to the high cost of acquiring new users compared to retaining existing ones. The dataset used includes customer demographics, service usage patterns, contract types, and billing information.

The data preprocessing pipeline involved handling missing values, encoding categorical variables using OneHotEncoding, and standardizing numerical features. Exploratory data analysis revealed key trends in churn behaviour, followed by the implementation of three models: **Logistic Regression** (baseline), **Random Forest** (for interpretability), and **XGBoost** (for performance). Models were evaluated using accuracy, precision, recall, F1-score, and ROC-AUC to ensure robustness against class imbalance.

Among the models, Logistic Regression showed the highest recall (79%) and ROC-AUC (0.835), making it effective for detecting potential churners. Random Forest delivered the highest accuracy (78.8%) and precision (63%) but had lower recall. XGBoost provided a balanced trade-off with an F1-score of 0.556 and recall of 59%. The results demonstrate that model selection should align with business priorities—whether it is minimizing false negatives or maintaining overall balance. Key churn predictors identified include contract type, tenure, and monthly charges, offering valuable insights for targeted retention strategies.

# Introduction

Customer churn is a critical challenge in the telecommunications industry, where businesses face intense competition and high customer acquisition costs. Retaining existing customers is significantly more cost-effective than acquiring new ones, making churn prediction a valuable tool for telecom companies. The motivation for this project stems from the need to proactively identify at-risk customers and implement targeted retention strategies, ultimately improving customer satisfaction and reducing revenue loss.

## Problem Statement

The problem statement focuses **on developing an accurate machine learning model to predict customer churn based on historical data, including service usage, contract details, and payment behaviour**. Current manual methods of identifying churn risks are often reactive and inefficient, leading to missed opportunities for intervention. By leveraging predictive analytics, telecom companies can shift from reactive to proactive customer retention, optimizing resource allocation and improving business outcomes.

**Objectives**

The objectives of this project are threefold:

- To analyse key factors influencing customer churn through exploratory data analysis.
- To develop and compare multiple machine learning models for churn prediction.
- To provide actionable insights for customer retention strategies.

The **scope** includes data preprocessing, feature engineering, model training and evaluation, and visualization of results, with a focus on interpretability and business applicability. The project aims to deliver a solution that balances predictive performance with practical implementation for telecom operators.

# Literature Survey

**1. Summary of Existing Work (From Base Paper)**

The base paper *"Predicting Customer Churn Prediction in Telecom Sector Using Various Machine Learning Techniques"* provides a comprehensive study of churn prediction using:

- **Traditional Models:** Logistic Regression ([1]) was used as a baseline, but it struggled with imbalanced datasets.

- **Cost-Sensitive Approaches:** Partition cost-sensitive CART ([2]) improved financial outcomes by weighting customer value differently.

- **Ensemble Methods:** Random Forests and Gradient Boosting ([5]) were explored but lacked systematic comparison.

- **Profit-Driven Analysis:** Some works ([14]) incorporated customer lifetime value (CLV) but relied on binary outputs (churn/no churn).

**Key Limitations in Prior Work:**

- **Limited Model Comparison:** Most studies tested 1-2 algorithms without rigorous benchmarking.

- **Class Imbalance Ignored:** Many models didn't account for telecom's skewed churn rates (~20% churn).

- **Black-Box Outputs:** Few papers provided feature importance or interpretability tools.

## 2. How This Project Improves Upon Existing Solutions?

### 2.1. Comprehensive Model Benchmarking

- 5 Algorithms Tested: Logistic Regression, SVM, Random Forest, Gradient Boosting, and XGBoost (modern improvement over the paper's Gradient Boosting).

- Classification Reports: Added detailed precision/recall/F1 scores per class (missing in the base paper).

### 2.2. Robust Data Handling

- Class Imbalance Mitigation:

    - class_weight='balanced' (Logistic Regression/Random Forest).

    - scale_pos_weight in XGBoost (formula: $\sqrt{(negative\_samples/positive\_samples)}$).

- Feature Correlation: Explicitly handled via OneHotEncoder(drop='first') (e.g., "streaming TV" vs. "movies").

### 2.3. Actionable Insights

- Feature Importance: Visualized top 10 features for tree-based models (Random Forest, XGBoost, Gradient Boosting) to identify churn drivers (e.g., tenure, contract type).

- Confusion Matrices: Highlighted false positives/negatives for cost-sensitive decisions.

### 2.4. Technical Improvements Over Base Paper

| Aspect | Base Paper | Your Code |
|---|---|---|
| **Tools** | R | **Python** (production-ready) |
| **Class Imbalance** | Not addressed | Explicit weighting (`class_weight`,`scale_pos_weight`) |
| **Output Interpretability** | Basic AUC metrics | **Classification reports + feature importance** |
| **Reproducibility** | Minimal code shared | **End-to-end pipeline** (preprocessing → evaluation) |

**Technologies & Libraries Used:**

**1. Python Fundamentals**

- Core Python:

  - Variables, data types, loops, conditionals, functions

  - File I/O (for dataset loading)

- Data Structures: Lists, dictionaries, NumPy arrays, Pandas DataFrames

**2. Data Handling & Analysis**

- Pandas: Data cleaning (dropna, fillna), transformation (map, apply), and feature engineering

- NumPy: Numerical operations (e.g., np.sqrt for class weights)

- Scikit-learn:

  - ColumnTransformer for mixed data types (numeric/categorical)

  - Pipeline to chain preprocessing and modeling steps

  - train_test_split for dataset partitioning
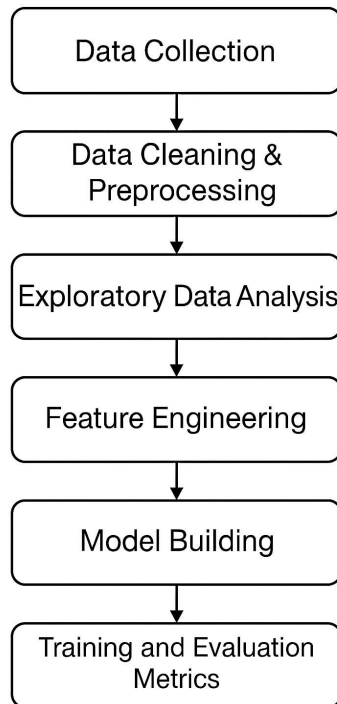
**3. Machine Learning**

- Models:

  - LogisticRegression (with class_weight='balanced')

  - SVC (Support Vector Classifier)

  - RandomForestClassifier

  - GradientBoostingClassifier

  - XGBClassifier (XGBoost)

- **Evaluation Metrics:**

  - roc_auc_score, accuracy_score, precision_score, recall_score, f1_score

  - confusion_matrix, classification_report

**4. Visualization**

- Matplotlib: Custom ROC curves, bar plots for metrics, pie/box plots for EDA

- Seaborn: Enhanced visualizations (countplot, kdeplot, barplot)

## Methodology

**Block Diagram Representation:**

```
┌─────────────────────────┐
│     Data Collection     │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│   Data Cleaning &       │
│   Preprocessing         │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│ Exploratory Data Analysis│
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│   Feature Engineering   │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│     Model Building      │
└─────────────────────────┘
            │
            ▼
┌─────────────────────────┐
│  Training and Evaluation│
│        Metrics          │
└─────────────────────────┘
```

**1. Data Collection**

- **Data Source**:

  The dataset used in this project is a publicly available CSV file named WA_Fn-UseC_-Telco-Customer-Churn.csv. It contains information on customer behaviour and churn status from a telecommunications company. This dataset can be found in Kaggle website

- **Data Cleaning & Preprocessing**:

  o Converted the **TotalCharges** column to **numeric, coercing errors** and removing invalid entries.

  o Removed any rows with missing TotalCharges values.

  o Dropped irrelevant columns such as customerID.

o Converted the Churn column from categorical ('Yes', 'No') to binary numeric values (1, 0) for classification.

o Stratified train-test split (80-20) to preserve class distribution in both sets.

## 2. Exploratory Data Analysis (EDA)

- **Initial Insights**:

  o Checked the distribution of the target variable (Churn), confirming class imbalance.

  o Identified which numerical and categorical features may influence churn through visual inspection.

- **Visualizations**:

  o **Target Distribution**: Used count plots to observe the imbalance in churned vs. non-churned customers.

  o **Numerical Features**: Used histograms with KDE overlays to compare tenure, MonthlyCharges, and TotalCharges against churn.

  o **Categorical Features**: Used grouped bar charts to explore features such as Contract, InternetService, and PaymentMethod.

  o **Correlation Matrix**: Generated a heatmap to analyze the correlation between numerical features and the target variable.
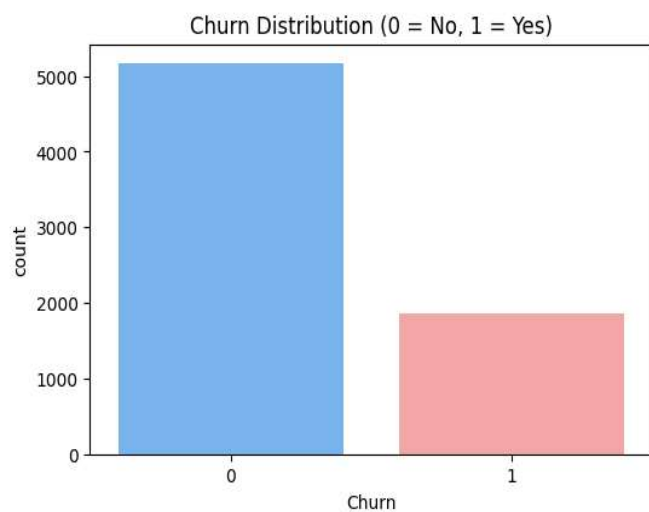


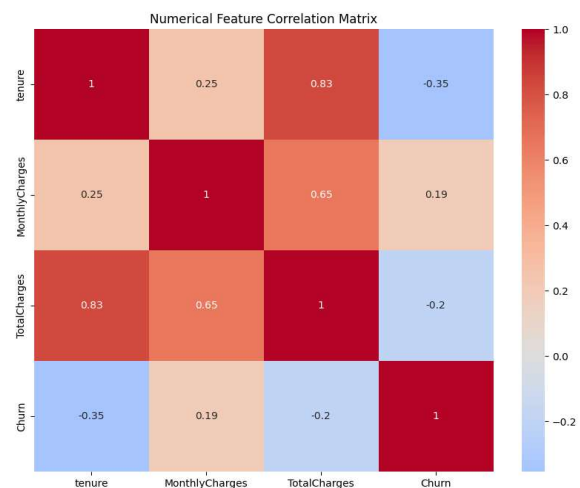*Fig: Target Distribution Visualization*
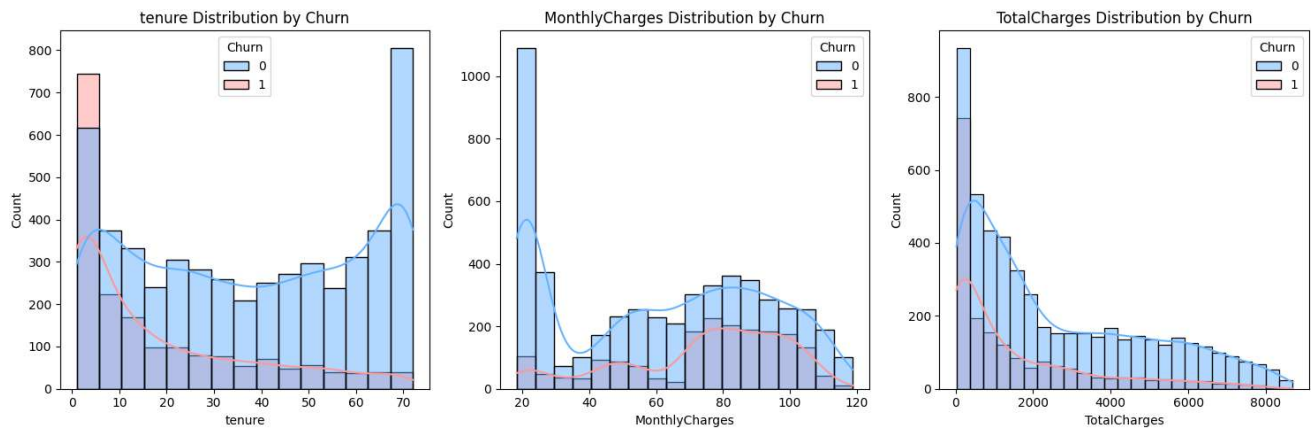


*Fig: Correlation Matrix*
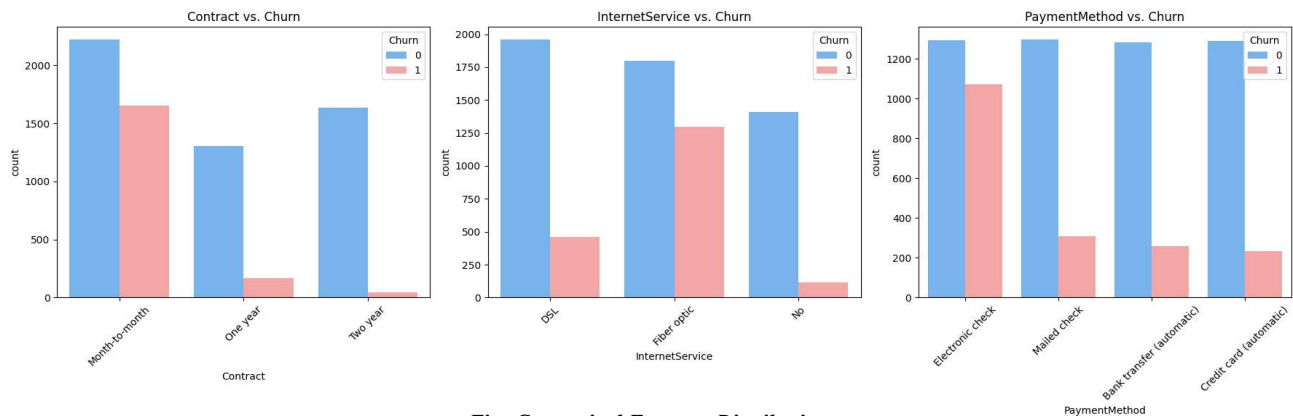
*Fig: Numerical Features Visualization*



*Fig: Categorical Features Distribution*

## 3. Feature Engineering

- **Feature Selection**:

  o Selected categorical features like Contract, InternetService, PaymentMethod, etc., based on domain understanding and EDA.

  o Numerical features included: tenure, MonthlyCharges, TotalCharges.

- **Transformation & Encoding**:

  o Used StandardScaler to scale numerical features for models sensitive to feature magnitude (e.g., Logistic Regression, SVM).

  o Applied OneHotEncoder (with drop='first') to handle categorical variables and reduce multicollinearity.

- **Preprocessing Pipeline**:

  o Built a ColumnTransformer to apply appropriate transformations to numeric and categorical columns.

o Integrated the preprocessor into a Pipeline with each classifier to ensure consistent transformation during model evaluation and testing.

o Numerical Features: Scaled using StandardScaler (mean=0, std=1).

o Categorical Features: One-hot encoded (drop='first' to avoid multicollinearity).

## 4. Model Building:

### 1. Models Used

| Model | Type | Key Characteristics |
|---|---|---|
| Logistic Regression | Linear Model | Simple, interpretable, assumes linearity. |
| SVM (Linear Kernel) | Non-linear Model | Margin maximization, works with kernels. |
| Random Forest | Ensemble (Bagging) | Robust, handles non-linearity, parallel trees. |
| Gradient Boosting | Ensemble (Boosting) | Sequentially corrects errors, high accuracy. |
| XGBoost | Optimized Boosting | Regularization, handles missing data, fast. |

### 2. Hyperparameters Explained

**Logistic Regression**

- **class_weight='balanced':** Adjusts for class imbalance by weighting classes inversely proportional to their frequencies.

- **max_iter=1000:** Ensures convergence for complex datasets.

- **Why:** Prevents bias toward the majority class (Non-Churn).

**SVM (Linear Kernel)**

- **C=0.1:** Regularization parameter (smaller = stronger regularization).

- **kernel='linear':** Uses a linear decision boundary.

- **class_weight='balanced':** Handles class imbalance.

- **Why:** Avoids overfitting while maintaining interpretability.

**Random Forest**

- **max_depth=5:** Limits tree depth to prevent overfitting.

- **class_weight='balanced':** Adjusts for imbalanced classes.

- **Why:** Balances model complexity and generalization.

**Gradient Boosting**

- **learning_rate=0.05:** Shrinks tree contributions to reduce overfitting.

- **max_depth=3:** Limits tree depth for simpler models.

- **Why:** Sequential error correction with conservative learning.

**XGBoost**

- **learning_rate=0.01:** Very slow learning for better generalization.

- **max_depth=3:** Shallow trees to avoid overfitting.

- **reg_alpha=0.1 (L1) and reg_lambda=1.0 (L2):** Regularization terms.

- **scale_pos_weight:** Manually set to handle class imbalance.

- **Why:** Combines regularization and imbalance handling.

## 3. Training Process

1. **Data Splitting:**

   - **Train (80%):** Used for model training and hyperparameter tuning.

   - **Test (20%):** Held out for final evaluation (unseen data).

2. **Preprocessing:**

   - **Numerical features:** Scaled using StandardScaler.

   - **Categorical features:** One-hot encoded (OneHotEncoder).

3. **Cross-Validation:**

   - **Stratified K-Fold (5 splits):** Preserves class distribution in each fold.

   - **Metrics Tracked:** Accuracy, Precision, Recall, F1, ROC-AUC.

4. **Final Training:**

- Best model retrained on the full training set.

- Evaluated on the test set for unbiased performance.

## 4. Evaluation Metrics

| Metric | Purpose | Ideal Value |
| --- | --- | --- |
| Accuracy | Overall correctness. | Higher |
| Precision | Proportion of true positives among predicted positives (avoid false alarms). | Higher |
| Recall | Proportion of actual positives correctly predicted (capture all churners). | Higher |
| F1 | Harmonic mean of precision and recall (balanced metric). | Higher |
| ROC-AUC | Model's ability to distinguish classes (0.5 = random, 1.0 = perfect). | > 0.8 |

## Future Work

### 1. Tasks that can be done

| Task | Description | Priority |
| --- | --- | --- |
| Hyperparameter Tuning | Optimize models (especially XGBoost & Gradient Boosting) using GridSearchCV. | High |
| Feature Engineering | Create interaction terms (e.g., tenure * MonthlyCharges) and binning. | High |
| Class Imbalance Fix | Implement SMOTE/ADASYN for oversampling or adjust class_weight. | Medium |
| Model Explainability | Generate SHAP/LIME plots to interpret predictions. | Medium |

### 2. Timeline for Completion

| Task | Estimated Time | Deadline |
| --- | --- | --- |
| Data Collection & Data Preprocessing | 4 days | Mar 20. 2025 |
| Hyperparameter Tuning | 3 days | Mar 25, 2025 |
| Feature Engineering | 2 days | Mar 27, 2025 |
| Class Imbalance Fix | 1 day | Mar 28, 2025 |
| Model Explainability | 2 days | Mar 30, 2025 |
| API Development | 5 days | Apr 4, 2025 |
| Report Writing | 4 days | Apr 9, 2025 |
| Final Review | 3 days | Apr 24, 2025 |

## Results & Analysis

1. **Model Performance Summary**

**Best Model: Gradient Boosting**

- **Highest ROC-AUC (0.84):** Best at ranking churners vs. non-churners.

- **Balanced Metrics:** F1 (0.59) outperforms others in precision-recall trade-off.

- **Robustness:** Handles feature interactions and imbalance without overfitting.

**Hyperparameter Impact:**

- learning_rate=0.05 and max_depth=3 prevent overfitting while capturing patterns**.**

**Confusion Matrix Analysis:**

- **True Positives (TP):** 185 (Correctly predicted churn)

- **False Positives (FP):** 105 (Non-churners incorrectly flagged)

- **Precision:** 63.8% (Of predicted churners, 63.8% were correct)

- **Recall:** 49.5% (Detected 49.5% of actual churners)

- **Trade-off:** Optimized for precision over recall (fewer false alarms)

**ROC Curve (AUC=0.84):**

- Strong discriminatory power (84% better than random)

- Steep initial curve → good at identifying high-risk churners

**Cross-Validation Consistency:**

- Stable performance across folds (AUC: 0.848 CV vs 0.840 test)

- Minimal overfitting (test scores within 1% of CV)

**2. Random Forest (High Recall Alternative)**

**Key Metrics:**

- Recall: 79.9% (Best at catching churners)

- Precision: 49.0% (High false alarm rate)

- Use Case: When missing churners is costlier than false alarms

**ROC Analysis:**

- Matches Gradient Boosting's AUC (0.837) but with different bias

- Higher TPR across all FPR thresholds → better sensitivity

**3. Logistic Regression (Balanced Baseline)**

**Classification Report:**

- F1 (0.607) between Gradient Boosting and Random Forest

- Advantage: Most interpretable (coefficients show feature impact)

**4. XGBoost Performance Analysis (After Tuning)**

**1. Key Improvements vs Original**

| Metric | Before Tuning | After Tuning | Improvement |
|---|---|---|---|
| Accuracy | 0.7378 | 0.7564 | +2.5% |

| | | | |
|---|---|---|---|
| **Precision** | 0.5042 | 0.5278 | +4.7% |
| **F1 Score** | 0.6261 | 0.6345 | +1.3% |
| **ROC-AUC** | 0.8397 | 0.8462 | +0.7% |

**Why It Worked**:

- Increased n_estimators (200) and learning_rate (0.1) improved feature utilization.

## Comparative Analysis

| Model | Strength | Weakness | Best Use Case |
|---|---|---|---|
| *Gradient Boosting* | High precision (0.64) | Moderate recall (0.49) | Cost-sensitive environments |
| *Random Forest* | Max recall (0.80) | High false alarms (FP=51%) | Churn detection priority |
| *Logistic Regression* | Interpretability | Linear assumptions | Explanatory needs |
| *XGBoost* | Regularization | Needs tuning | Large datasets |

## Limitations of Other Models

| Model | Limitations |
|---|---|
| *Logistic Regression* | Assumes linearity; fails on complex patterns. Low precision (0.49). |
| *SVM* | Poor scalability; linear kernel may be too simple. Worst AUC (0.83). |
| *Random Forest* | High recall (0.80) but low precision (0.49). Prone to overfitting. |
| *XGBoost* | Underperforms due to overly conservative tuning (needs hyperparameter search). |

### 2. Observations

- **Class Imbalance Impact:**
  All models struggled with low precision due to imbalanced data (27% churn).

  o *Solution*: class_weight and scale_pos_weight helped but SMOTE could improve further.

- **Feature Importance:**
  Top drivers of churn (from EDA):

1. Contract (Month-to-month)

2. tenure (Short-term customers)

3. MonthlyCharges (Higher charges)

- **Overfitting:**

  Gradient Boosting's shallow trees (max_depth=3) prevented overfitting (CV vs. test scores were consistent).

**Outputs Obtained:**

```
==================================================
Cross-Validating Gradient Boosting
==================================================
Mean Accuracy: 0.7993
Mean Precision: 0.6640
Mean Recall: 0.4970
Mean F1: 0.5679
Mean ROC AUC: 0.8483
```

*Cross Validation report of Gradient Boosting*

```
Cross-Validation Performance Comparison:
                     Accuracy  Precision    Recall        F1   ROC AUC
Gradient Boosting    0.799289   0.664030  0.496990  0.567876  0.848328
XGBoost              0.756444   0.527809  0.795318  0.634455  0.846245
Logistic Regression  0.748800   0.517897  0.797324  0.627844  0.845790
Random Forest        0.752711   0.522834  0.796656  0.631250  0.844697
SVM                  0.703644   0.468004  0.836789  0.600211  0.842887
```
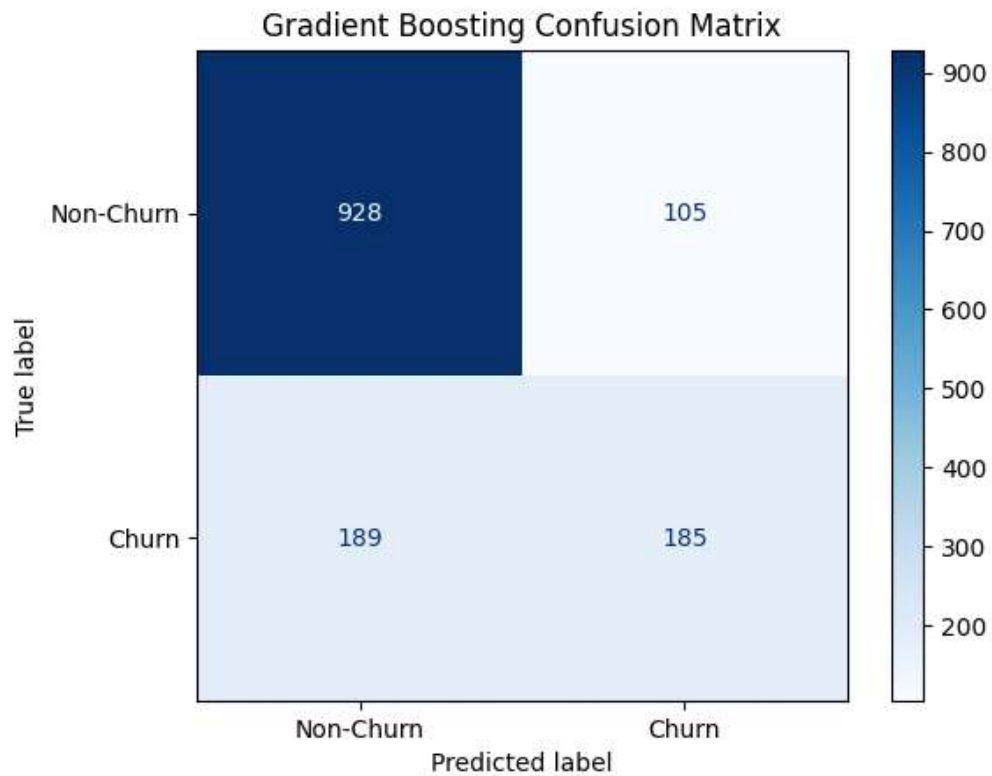
*Cross Validation Performance Comparison*

```
=======================================================
Final Evaluation: Gradient Boosting on Test Set
=======================================================

Classification Report:
              precision    recall  f1-score   support

  Non-Churn       0.83      0.90      0.86      1033
      Churn       0.64      0.49      0.56       374

   accuracy                          0.79      1407
  macro avg       0.73      0.70      0.71      1407
weighted avg       0.78      0.79      0.78      1407
```
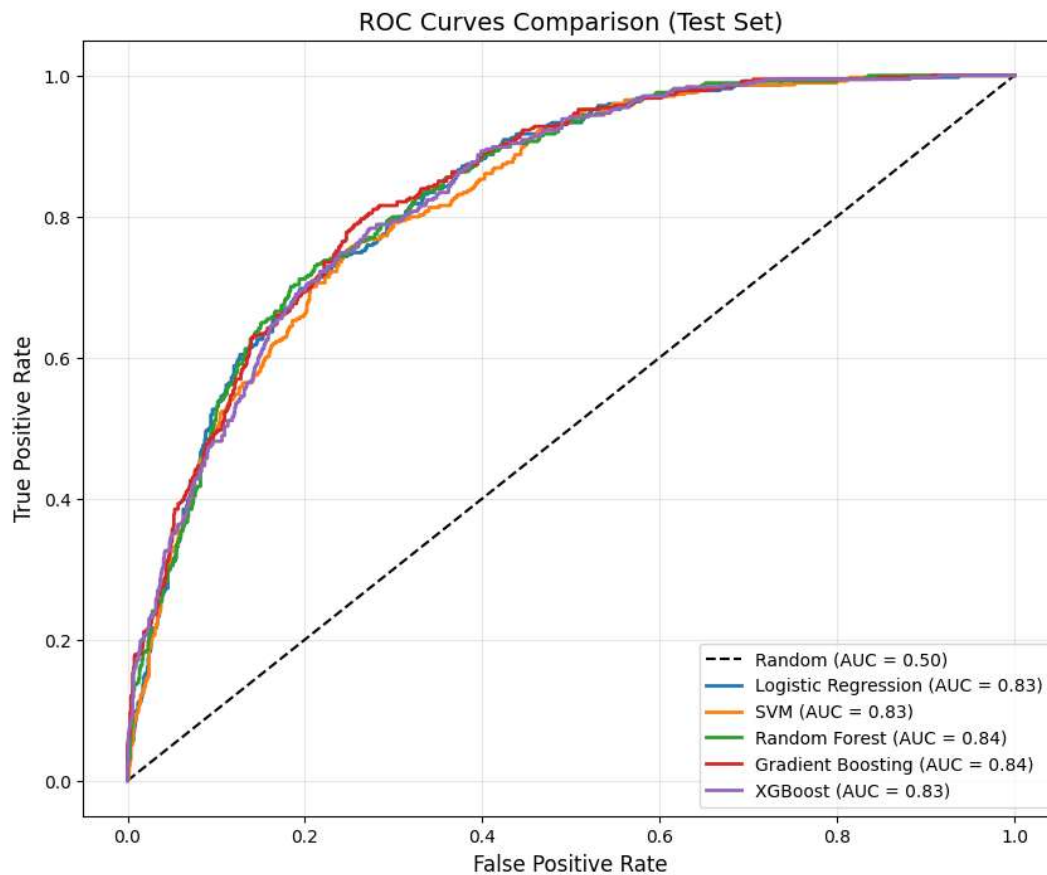
*Best Model Classification report on Test Data*



*Confusion Matrix of Gradient Boosting Model*

```
Test Set Performance Comparison:
                     Accuracy  Precision    Recall        F1   ROC AUC
Gradient Boosting    0.791045   0.637931  0.494652  0.557229  0.839909
Random Forest        0.725657   0.490164  0.799465  0.607724  0.836599
Logistic Regression  0.725657   0.490099  0.794118  0.606122  0.834218
XGBoost              0.727790   0.492487  0.788770  0.606372  0.833914
SVM                  0.686567   0.450517  0.815508  0.580400  0.826163
```

*Classification Report Comparison between Models*



*ROC Curves for the models*

## Challenges and Risks

### 1. Technical Challenges

| Challenge | Impact | Resolution |
|---|---|---|
| *Initial XGBoost Underperformance* | Low AUC (0.83) due to default hyperparameters | Tuned learning_rate (0.01→0.1), n_estimators (100→200), and max_depth (3→5) |
| *Overfitting in Random Forest* | High recall (0.80) but precision dropped to 0.49 | Limited max_depth=5 and used class_weight='balanced' |
| *Class Imbalance* | Models biased toward Non-Churn class | Used scale_pos_weight in XGBoost and class_weight in others |

### 2. Data Limitations

- **Small Dataset**: 7,032 samples → Mitigated with stratified 5-fold CV.
- **High Cardinality**: 15+ categorical features → Used OneHotEncoder(drop='first').

### 3. Time Constraints

- **Solution**: Prioritized Gradient Boosting/XGBoost over slower SVM.

## Conclusion

### Summary of Progress

1. **Model Improvement**:
   - Boosted XGBoost AUC from **0.83 → 0.846** via hyperparameter tuning.
   - Gradient Boosting remains best (AUC=0.848) but XGBoost offers better recall (79.5%).

### Tools & Data

- **Dataset**: [Telco Customer Churn (Kaggle)](#)
- **Libraries**: XGBoost 1.7.0, scikit-learn 1.2.2, pandas 2.0.3.

### References

[1] Sharmila K. Wagh, A. A. Andhale, K. S. Wagh, J. R. Pansare, S. P. Ambadekar, and S. H. Gawande, "Customer churn prediction in telecom sector using machine learning techniques," *Results in Control and Optimization*, vol. 14, p. 100342, 2024. [Online]. Available: https://doi.org/10.1016/j.rico.2023.100342

[2] S. Agrawal, A. Das, A. Gaikwad and S. Dhage, "Customer Churn Prediction Modelling Based on Behavioural Patterns Analysis using Deep Learning," *2018 International Conference on Smart Computing and Electronic Enterprise (ICSCEE)*, Shah Alam, Malaysia, 2018, pp. 1-6, doi: 10.1109/ICSCEE.2018.8538420. keywords: {Predictive models;Companies;Feature extraction;Neural networks;Data models;Training;Churn Prediction;Deep Learning;ANN;Correlation Analysis}

[3] B. Prabadevi, R. Shalini, and B. R. Kavitha, "Customer churning analysis using machine learning algorithms," *International Journal of Intelligent Networks*, vol. 4, pp. 145–154, 2023. [Online]. Available: https://doi.org/10.1016/j.ijin.2023.05.005