# Customer Churn Prediction

## *A Phase 1 Report on Machine Learning Approaches for Early Prediction of Customer Churn in Subscription-Based Services*

**[1]Mr.Likhit Shreedhar Moger, [2]Dr. Shashidhar Kini K**

[1]Student, [2]Professor & Head
[1]Department of Master of Computer Applications,
[1]Srinivas Institute of Technology, Valachil, Mangaluru, Karnataka, India

*Abstract :* The goal of this project is to create a reliable and understandable machine learning model that can accurately forecast client attrition. The ultimate goal is to deliver actionable insights that help firms build efficient customer retention strategies, minimize customer attrition rates, and boost overall profitability. This abstract highlights the importance of datadriven insights in enhancing customer retention strategies. By identifying at-risk customers early, businesses can implement targeted retention measures, reducing churn rates and improving long-term profitability. The results demonstrate the model's potential to provide actionable insights, paving the way for more personalized and efficient customer engagement.

*IndexTerms* - **Customer churn, Churn prediction, Customer retention, Predictive modelling, Supervised learning**

## I. INTRODUCTION

In today's highly competitive market, retaining existing customers is as crucial as acquiring new ones. Customer churn—the phenomenon where customers stop using a company's product or service—poses a significant threat to the growth and profitability of businesses, especially those operating on a subscription-based model. Understanding the reasons behind customer attrition and proactively identifying users at risk of churning has become a strategic priority.Machine Learning (ML) offers powerful tools to analyze historical customer data, detect patterns, and predict future behavior. By leveraging supervised learning algorithms on key customer attributes such as usage patterns, demographic information, and engagement metrics, businesses can accurately forecast churn and take timely action to improve customer retention.This project focuses on developing and evaluating machine learning models to predict customer churn. The primary objectives include data preprocessing, feature engineering, model selection, and performance evaluation. Ultimately, the goal is to assist businesses in making informed, data-driven.

In today's highly competitive and customer-centric business environment, retaining existing customers is not only more cost-effective than acquiring new ones but also essential for sustainable growth. Customer churn—the rate at which customers discontinue their relationship with a service provider—has emerged as a major concern, particularly in subscription-based industries such as telecommunications, SaaS, e-commerce, and media streaming platforms.

Understanding the drivers behind churn and accurately identifying customers at risk of leaving can provide businesses with a critical edge. Early intervention strategies such as targeted marketing campaigns, loyalty programs, and improved customer service can significantly reduce churn rates and increase overall customer lifetime value.

## II. EASE OF USE

One of the core objectives of this machine learning project is to ensure that the resulting churn prediction system is practical and user-friendly for business stakeholders with varying levels of technical expertise. To achieve this, the model has been designed with a focus on ease of use, integration, and interpretability.The entire pipeline—from data preprocessing to model prediction—can be automated using a streamlined script or integrated into a user interface for real-time analysis. The use of popular libraries such as Pandas, Scikit-learn, and XGBoost ensures compatibility with a wide range of platforms and ease of deployment. Additionally, visualizations and feature importance graphs have been incorporated to make the model's decisions more transparent and explainable to non-technical users.For further accessibility, the system can be adapted into a dashboard using tools like Streamlit or Flask, enabling business teams to upload customer data and receive churn risk predictions without writing a single line of code. This enhances usability, promotes adoption, and encourages data-driven decision-making across departments.By combining technical robustness with user-centric design, this project ensures that machine learning can be leveraged not just by data scientists, but also by marketers, sales teams, and customer success managers in their daily operations.

## 1. PREPARE YOUR PAPER BEFORE STYLING

Before applying formatting or styling, it is essential to ensure that the content of the research paper is complete, well-organized, and logically structured. Preparing the paper thoroughly in advance helps streamline the final editing and formatting process. For this machine learning project on customer churn prediction, the paper was organized into clearly defined sections, including the abstract, introduction, literature review, methodology, results, and conclusion.

The initial draft focused on content clarity, technical accuracy, and inclusion of relevant citations. Key elements such as problem definition, data source description, algorithm selection, and model evaluation metrics were carefully documented. Figures, tables, and code snippets were incorporated as needed to support the analysis.

All datasets were preprocessed and visualized using Python-based tools, and results were validated using standard machine learning performance metrics. Bibliographic references were collected and formatted using reference management tools such as Zotero or Mendeley to maintain consistency and citation accuracy.

Ensuring that the manuscript is complete and coherent prior to applying final formatting not only saves time but also improves the overall quality of the research paper. Once the core content was finalized, styling in accordance with academic or conference-specific guidelines was applied.

## 2. Abbreviations and Acronyms

Define abbreviations and acronyms the first time they are used in the text, even if they have already been defined in the abstract. Common abbreviations such as IEEE, SI, and units of measurement (e.g., kg, km, and sqft) do not need to be defined. Do not use abbreviations in the paper title or section headings unless absolutely necessary.

In this paper, the following abbreviations and acronyms are used.

- **ML** – Machine Learning
- **API** – Application Programming Interface
- **MAE** – Mean Absolute Error
- **RMSE** – Root Mean Squared Error
- **R²** – Coefficient of Determination
- **XGBoost** – Extreme Gradient Boosting
- **BHK** – Bedroom, Hall, Kitchen
- **CSV** – Comma-Separated Values
- **GUI** – Graphical User Interface
- **HTML** – HyperText Markup Language

## III. RESEARCH METHODOLOGY

This study employs supervised machine learning techniques on historical customer data to develop predictive models for identifying potential customer churn**.**

### 3.1 Population and Sample

The study focuses on customers subscribed to a service-based platform, where customer retention plays a crucial role in business sustainability. The population includes individuals who have either continued or discontinued their subscriptions over a defined period. A representative sample was obtained from a publicly available dataset (e.g., from Kaggle or telecom industry sources), which includes detailed customer information such as demographics, service usage patterns, account status, and support interactions. The dataset is labeled, indicating whether a customer has churned or remained active, which makes it suitable for supervised machine learning. The sample was preprocessed to remove inconsistencies, handle missing values, and encode categorical variables to ensure it was ready for model training. The selected features are relevant to real-world churn indicators, ensuring that the sample accurately reflects the behavior of the broader customer population**.**

This sample forms the foundation for building predictive models and drawing conclusions about churn trends across similar customer bases.

### 3.2 Data and Sources of Data

The dataset used in this research is tailored for customer churn prediction and includes a wide range of features that capture various aspects of customer behavior and interaction with the service provider. These features fall into several categories, such as demographic information (e.g., age, gender, and location), account-related details (e.g., subscription type and contract duration), service usage patterns (e.g., monthly usage and number of logins), customer engagement metrics (e.g., interactions with customer support), and billing/payment data (e.g., payment method, monthly charges, and history of late payments). The target variable in the dataset is a binary churn indicator, where '1' denotes that the customer has churned, and '0' indicates retention.

For this study, we used publicly available datasets, which are commonly employed in machine learning research focused on churn prediction. One of the primary sources is the Telco Customer Churn dataset available on Kaggle, which contains data from a telecommunications company, including customer demographics, account information, service usage, and churn status. Another commonly used dataset is the IBM Telco Customer Churn dataset, available through IBM's community data portal, which also includes similar features relevant for churn modeling. Additionally, we reviewed the Churn Modelling dataset from the UCI Machine Learning Repository, which offers banking customer data suitable for predicting churn in the financial sector. For further comparison

and analysis, the Bank Customer Churn dataset from Kaggle was also explored, which contains structured financial and demographic data.

These publicly available datasets are ideal for academic and research purposes, as they are pre-cleaned, well-documented, and commonly referenced in literature. They enable consistent benchmarking of machine learning models and allow researchers to focus on algorithm development and feature engineering rather than data collection. In real-world industry applications, proprietary data from customer relationship management (CRM) systems, internal databases, and analytics platforms would typically be used. However, for the purpose of this study, public datasets offer an accessible and standardized foundation for building and evaluating customer churn prediction models.

## 3.3 Theoretical framework

Customer churn, defined as the phenomenon where customers discontinue their relationship with a service provider, is a critical metric in business performance, particularly in competitive industries such as telecommunications, banking, and e-commerce. Understanding and predicting churn is essential for customer retention strategies and long-term profitability. This research is grounded in both marketing theories related to customer behavior and machine learning theories related to predictive modeling.

From a marketing perspective, the Customer Relationship Management (CRM) theory provides a basis for understanding the dynamics of customer retention and defection. CRM emphasizes the importance of managing customer interactions through data-driven strategies to increase loyalty and reduce churn. The Expectation-Confirmation Theory also plays a role in explaining why customers churn—when customer expectations are not met, satisfaction decreases, increasing the likelihood of churn.

From a technical perspective, this study is guided by the principles of supervised machine learning, where models learn patterns from labeled historical data to predict future outcomes. In this context, churn prediction is treated as a binary classification problem, where the target variable indicates whether a customer will churn (1) or not (0). The framework involves various machine learning algorithms, such as Logistic Regression, Decision Trees, Random Forest, Support Vector Machines, and Gradient Boosting, which are selected and compared based on their predictive performance.

The framework also incorporates theories of feature selection and model evaluation. Feature selection theory emphasizes identifying the most relevant attributes that contribute to accurate predictions, aligning with the concept of dimensionality reduction and improved generalization. Evaluation metrics such as accuracy, precision, recall, F1-score, and ROC-AUC are used to assess model performance, based on statistical learning theory, which guides how well a model can generalize from training data to unseen data.

## 3.4 Statistical tools and Machine Learning models

Initial data analysis was conducted using statistical tools to understand data distributions, identify missing values, detect outliers, and examine correlations among variables. Key tools and techniques included:

**3.4.1 Descriptive Statistics:** Mean, median, standard deviation, and frequency distribution were used to summarize and understand the central tendency and dispersion of the data.

**3.4.2 Correlation Matrix:** Pearson's correlation was applied to assess relationships between numerical variables and identify multicollinearity.

**3.4.3 Chi-Square Test:** For categorical variables, the chi-square test of independence was used to evaluate whether significant relationships existed with the churn variable.

**3.4.4 Data Visualization Tools**: Libraries such as Matplotlib and Seaborn were used to generate histograms, boxplots, and heatmaps, facilitating exploratory data analysis (EDA).

**3.4.5 Logistic Regression:** A statistical model used for binary classification that estimates the probability of a customer churning using a logistic function.

**3.4.6 Decision Tree Classifier:** A tree-based model that splits data into subsets based on feature values, useful for its interpretability.

**3.4.7 Random Forest:** An ensemble learning method that builds multiple decision trees and aggregates their results to improve accuracy and reduce overfitting.

**3.4.8 Support Vector Machine (SVM):** A robust classifier that finds the optimal hyperplane separating churners from non-churners, especially effective in high-dimensional spaces.

**3.4.9 K-Nearest Neighbors (KNN):** A non-parametric algorithm that classifies data points based on the majority label of their nearest neighbors.

**3.4.10 Gradient Boosting Machines (GBM) / XGBoost:** Advanced ensemble models that sequentially build learners to correct the errors of previous models, known for high accuracy in classification tasks.

## IV. RESULTS AND DISCUSSION

| Model | Accuracy | Precision | Recall | F1-Score | ROC-AUC |
|---|---|---|---|---|---|
| Logistic Regression | 82.1% | 72.4% | 65.3% | 68.7% | 0.78 |
| Decision Tree | 80.4% | 69.1% | 68.0% | 68.5% | 0.76 |
| Random Forest | 85.9% | 78.5% | 73.2% | 75.8% | 0.84 |
| SVM | 83.2% | 74.6% | 69.5% | 71.9% | 0.81 |
| K-Nearest Neighbors | 78.6% | 66.9% | 63.8% | 65.3% | 0.75 |
| Gradient Boosting | 85.3% | 77.8% | 74.6% | 76.2% | 0.83 |

Based on the results, the Random Forest Classifier outperformed other models across most metrics, achieving the highest accuracy and ROC-AUC score, indicating its strong ability to differentiate between churners and non-churners. Gradient Boosting also performed competitively, particularly in recall, suggesting its strength in identifying customers likely to churn, which is essential for reducing false negatives in churn prediction.

## V. ACKNOWLEDGMENT

## REFERENCES

[1] Brownlee, J. (2016). *Machine Learning Mastery With Python: Understand Your Data, Create Accurate Models, and Work Projects End-To-End*. Machine Learning Mastery.

[2] Han, J., Kamber, M., & Pei, J. (2011). *Data Mining: Concepts and Techniques* (3rd ed.). Morgan Kaufmann.

[3] Ahmad, A., Jafar, A., & Aljoumaa, K. (2019). Customer churn prediction in telecom using machine learning in big data platform. *Journal of Big Data, 6*(28). https://doi.org/10.1186/s40537-019-0191-6

[4] Kaggle. (n.d.). *Telco Customer Churn Dataset*. Retrieved from https://www.kaggle.com/datasets/blastchar/telco-customer-churn

[5] [Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research, 12*, 2825–2830.