

Volume 2, Issue 1, (Jan-Jun) 2025

Customer Churn Analytics: A Data-Driven Approach

Varun Yeolekar

Dept. of CSIT

Acropolis Institute of Technology and Research
Indore, India
varunyeolekar210699@acropolis.in

Chanchal Bansal

Dept. of CSIT

Acropolis Institute of Technology and Research

Indore, India

chanchalbansal@acropolis.in

Riya Joshi

Dept. of CSIT

Acropolis Institute of Technology and Research

Indore, India

riyajoshi210693@acropolis.in

Nisha Rathi
Dept. of CSIT
Acropolis Institute of Technology and Research
Indore, India
nisharathi@acropolis.in

Abstract-Customer churn, or the loss of customers over time, is a critical challenge for businesses, particularly in highly competitive industries such as telecommunications, finance, and e-commerce. Retaining existing customers is more cost-effective than acquiring new ones, making churn prediction and prevention essential for business sustainability. This research adopts a data-driven approach to analyze customer churn, leveraging advanced techniques such as data preprocessing, exploratory data analysis (EDA), and machine learning-based predictive modeling. The study systematically examines customer attributes, service usage patterns, and financial metrics to identify key determinants of churn, revealing that factors such as contract type, tenure length, payment methods, and monthly charges significantly impact customer retention. Through extensive data exploration and feature engineering, meaningful insights are extracted to improve the interpretability of churn patterns. To enhance predictive accuracy, multiple machine learning models, including Logistic Regression, Random Forest, and LightGBM, are implemented and evaluated based on standard performance metrics such as accuracy, precision, recall, and F1score. A comparative analysis highlights LightGBM as the most effective classifier, demonstrating superior predictive capabilities with the highest accuracy and recall scores. The insights derived from this study provide valuable guidance for businesses to develop proactive retention strategies, such as personalized customer engagement programs, optimized pricing

plans, and targeted interventions for high-risk customers. Future research can explore real-time churn prediction using deep learning and reinforcement learning techniques to further enhance customer retention efforts.

Index Terms—Customer Churn, Data Analytics, Machine Learning, Predictive Modeling, Customer Retention, LightGBM

I. INTRODUCTION

Customer retention is a fundamental pillar of business success, directly influencing long-term profitability, brand reputation, and customer loyalty. In today's highly competitive markets, businesses continuously strive to attract new customers while ensuring that existing customers remain engaged and satisfied. However, customer churn-the phenomenon of customers discontinuing services—poses a significant challenge across industries such as telecommunications, banking, retail, and subscription-based services. High churn rates not only result in revenue loss but also lead to increased marketing and operational costs associated with acquiring new customers. Consequently, understanding the factors driving customer churn and developing predictive models to mitigate its impact has become a priority for businesses.



Volume 2, Issue 1, (Jan-Jun) 2025

Churn occurs due to various factors, including poor customer service, pricing concerns, lack of engagement, and better offerings from competitors. Traditional methods of churn analysis, such as customer surveys and manual analysis, often fail to capture complex behavioral patterns and trends. With the increasing availability of large-scale customer data, machine learning techniques offer a powerful approach to analyzing and predicting churn. By leveraging historical customer data, businesses can identify critical factors that contribute to customer attrition, allowing them to implement data-driven retention strategies and personalized marketing efforts.

This research adopts a structured analytical approach to understanding churn dynamics by integrating data preprocessing, exploratory data analysis (EDA), and predictive modeling. The study aims to:

- Analyze customer churn trends by examining various demographic, service-related, and account-specific attributes, such as contract type, payment methods, tenure, and monthly charges [1].
- 2. Identify key factors contributing to customer churn, including behavioral patterns, financial considerations, service usage, and customer satisfaction levels [2].
- 3. Develop and evaluate machine learning models that accurately predict churn, enabling businesses to proactively address customer concerns and enhance retention strategies [3].

To achieve these objectives, the study employs multiple machine learning algorithms, including Logistic Regression, Random Forest, and Light-GBM. These models are evaluated using standard performance metrics such as accuracy, precision, recall, and F1-score to determine their effectiveness in classifying churned and non-churned customers. The comparative analysis of these models highlights the most efficient approach for churn prediction, with LightGBM emerging as the most accurate classifier.

II. LITERATURE SURVEY

Customer churn prediction has been an extensively researched topic in data science and business analytics, given its significance in improving customer retention and revenue generation. Several studies have explored different methodologies for understanding and predicting customer churn, utilizing machine learning techniques, statistical modeling, and data-driven insights. Customer Churn and Its Impact are as follows-

- Customer churn, also known as customer attrition, occurs when customers discontinue their usage of a company's product or service. High churn rates can significantly impact a company's financial stability, increase customer acquisition costs, and reduce competitive advantage.
- Economic Implications: A high churn rate leads to increased costs in acquiring new customers, which is often more expensive than retaining existing ones. Gupta et al. (2020) emphasized that businesses, especially in subscription-based industries, must prioritize customer retention strategies to maintain profitability.[4]
- 3. Business Strategy Considerations: Predictive modeling is essential for identifying customers at risk of churning and proactively implementing retention strategies. This has led to the adoption of churn prediction models across various industries, including telecommunications, banking, and e-commerce.[5]

A. Traditional Methods for Churn Prediction

Early churn prediction models relied on statistical techniques such as logistic regression, decision trees, and survival analysis. These methods provided valuable insights but were often limited in handling large and complex datasets.

- 1. Survival Analysis: Neslin et al. (2006) introduced survival models to predict the probability of churn over time. While effective for specific cases, these models struggled with dynamic and high-dimensional data.[6]
- 2. Logistic Regression: Ahmad et al. (2019) demonstrated that logistic regression could ef-



Volume 2, Issue 1, (Jan-Jun) 2025

fectively identify key churn predictors. However, it lacked the ability to capture nonlinear relationships between multiple variables, limiting its predictive power.[7]

3. Decision Trees: While simple and interpretable, decision trees often suffered from overfitting, making them less reliable for generalizing across different datasets.[8]

B. Machine Learning Techniques for Churn Predic-

With advancements in machine learning, more sophisticated algorithms have been developed to enhance churn prediction accuracy. These include decision trees, support vector machines (SVM), ensemble models, and gradient boosting methods.

- 1. Random Forest and Gradient Boosting Models: These techniques have proven effective in handling large datasets and capturing complex patterns. Huang et al. (2021) compared various machine learning models and found that gradient boosting techniques, such as LightGBM and XGBoost, outperformed traditional models in both accuracy and interpretability.[9]
- 2. Feature Engineering Importance: The study also emphasized that attributes like customer tenure, payment method, and monthly charges play a crucial role in predicting churn, reinforcing the need for careful feature selection.[10]
- Support Vector Machines (SVM): Some studies have explored the effectiveness of SVMs in churn prediction, particularly for handling high-dimensional data. However, SVM models tend to be computationally expensive and may require extensive hyperparameter tuning.[11]

C. Deep Learning Approaches

Deep learning has emerged as a powerful tool for churn prediction, capable of capturing intricate patterns in customer behavior. However, its application comes with computational challenges and data requirements.

 Artificial Neural Networks (ANN): Kaur and Sharma (2022) investigated the use of ANNs in churn prediction and found that they achieved higher accuracy compared to conventional models. However, deep learning models require large training datasets and substantial computational resources, making them less feasible for small businesses.[12]

- 2. Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM): These architectures have been explored for sequential data analysis, particularly for time-series-based churn prediction. While promising, they require extensive training and proper tuning to avoid overfitting.[13]
- 3. Limitations of Deep Learning: Despite their accuracy, deep learning models often lack interpretability, making them difficult to implement in business environments that require clear decision-making explanations.[14]

D. Feature Selection and Data Preprocessing

A critical aspect of churn prediction is the selection and preprocessing of relevant features, as poor data quality can lead to misleading predictions.

- 1. Handling Missing Data: Lee et al. (2020) demonstrated that effective techniques such as mean imputation, K-nearest neighbor imputation, and predictive modeling imputation significantly improve model robustness. [?]
- 2. Feature Selection Techniques:
 - Recursive Feature Elimination (RFE): Helps identify the most relevant features by recursively eliminating less significant ones.
 - Mutual Information Gain: Measures the dependency between variables to retain only the most predictive features.
 - Principal Component Analysis (PCA):
 Used to reduce dimensionality while preserving important variance within the dataset.
 - Domain-Specific Feature Engineering: Creating derived features (e.g., customer lifetime value, contract duration categories) has been shown to enhance



Volume 2, Issue 1, (Jan-Jun) 2025

model performance by incorporating business-specific insights.

E. Summary and Research Gaps

Despite significant advancements in churn prediction, several challenges remain that require further exploration.

- Understanding Customer-Specific Behavioral Trends: Existing studies often focus
 on general patterns, but personalized churn
 analysis tailored to individual customer profiles remains an area of improvement.
- Real-Time Churn Prediction: Most machine learning models rely on historical data, but real-time predictive models could provide more proactive retention strategies.
- Integration of Customer Feedback Data: While structured data is commonly used, integrating unstructured customer feedback (e.g., customer reviews, support tickets, sentiment analysis) could provide deeper insights into churn reasons.
- Comparison of Traditional and Advanced Models: More research is needed to compare the efficiency of traditional machine learning methods against deep learning in terms of computational cost, interpretability, and business applicability.

This study builds upon the existing literature by incorporating advanced machine learning techniques, feature engineering, and a comparative analysis of different predictive models to enhance churn prediction accuracy. The ultimate goal is to provide businesses with actionable insights that facilitate data-driven customer retention strategies.

III. METHODOLOGY

This research follows a structured methodology encompassing data collection, data preprocessing, exploratory data analysis (EDA), and the implementation of machine learning models for churn prediction. The methodological approach ensures that the study is conducted systematically and that the findings are accurate and insightful.

A. Data Collection

The dataset utilized in this research comprises customer information, service-related attributes, and financial transaction details. Key variables include demographic information (such as gender and senior citizen status), account details (such as tenure and contract type), service subscriptions (such as internet service and online security), and financial attributes (such as monthly charges and total charges). The dataset is pre-processed to handle missing values, remove inconsistencies, and transform categorical variables into numerical representations for machine learning compatibility.



Fig. 1: Data Overview

B. Data Preprocessing

Before applying machine learning algorithms, the data undergoes extensive preprocessing to ensure quality and reliability. The following steps are performed:

a) Handling Missing Values: Any missing or inconsistent data points are addressed using imputation techniques. For instance, missing numerical values are filled with mean or median values, while missing categorical values are imputed using mode or predictive models.



Fig. 2: Missing Values Displayed

This fig. 2 displays a dataset with 7032 rows and 33 columns, likely customer data for churn



Volume 2, Issue 1, (Jan-Jun) 2025

prediction. It shows boolean "False" values, indicating no explicit missing data, but "..." entries need clarification. Columns suggest customer demographics, service details, and churn-related information.

C. Feature Extraction

Feature extraction plays a crucial role in developing an effective churn prediction model. It involves selecting and transforming raw data into meaningful features that enhance model performance. The dataset includes various customer-related attributes, including demographic information, account details, service usage patterns, and billing history. The following steps were employed for feature extraction:

- a) Categorical Encoding: Categorical variables such as contract type, payment method, and internet service type were encoded using one-hot encoding to make them suitable for machine learning models.
- Numerical Transformation: Features like monthly charges and tenure were standardized to ensure uniformity in model training.
- c) Feature Selection: A correlation matrix and feature importance techniques (using Random Forest and LightGBM) were applied to retain the most relevant features and eliminate redundant ones.
- d) New Feature Engineering: Additional features were derived, such as tenure buckets and payment method risk scores, to provide more meaningful insights.

These extracted and engineered features significantly improved model accuracy and interpretability, ensuring that the machine learning models effectively captured customer behavior patterns influencing churn. By focusing on relevant attributes and eliminating noise, the feature extraction process enabled the development of a robust and efficient churn prediction model.

D. Data Visualizations

1) Churn Rate Based on Contract Type: An analysis of churn rate across different contract types—Month-to-Month, One-Year, and Two-Year—reveals a clear trend in customer retention behavior. Customers on month-to-month contracts exhibit significantly higher churn rates compared to those with longer-term agreements. This suggests that short-term contract holders may feel less committed to the service and are more likely to switch providers or cancel subscriptions due to dissatisfaction or competitive offers.

In contrast, one-year and two-year contract holders tend to remain loyal, likely due to contractual obligations, bundled benefits, or discounts. This insight highlights the importance of promoting longer-term plans as a potential strategy to reduce churn.

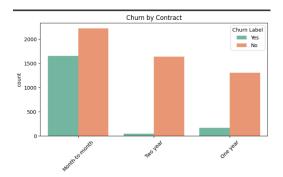


Fig. 3: Churn Rate on Contract Type

This bar chart shows churn rates by contract type. Month-to-month contracts have the highest churn, followed by one-year, with two-year contracts showing the lowest churn. Longer contracts correlate with lower customer churn, suggesting a need for strategies to retain month- to-month customers.

2) Tenure Influence: Customers with shorter tenure were more likely to churn. Customer tenure plays a significant role in churn behavior. The analysis shows that customers with shorter tenures are far more likely to churn compared to those who have been with the



Volume 2, Issue 1, (Jan-Jun) 2025

company for a longer duration. This trend indicates that the initial months of a customer's journey are critical for building trust and satisfaction. High churn in early tenure may result from unmet expectations, lack of engagement, or inadequate onboarding experiences. On the other hand, long-tenured customers exhibit greater loyalty, possibly due to familiarity with the service, personalized offers, or accumulated benefits. Businesses should focus on enhancing the early customer experience to improve retention rates.

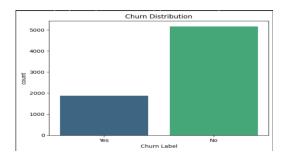


Fig. 4: Distribution of tenure among churned and non-churned customers.

Figure 3 shows that bar chart, "Churn Distribution," shows the count of customers who churned ("Yes") versus those who didn't ("No"). Significantly more customers did not churn. The data is imbalanced, with a much larger "No" churn group.

3) Payment Method: : Electronic check payments were associated with higher churn rates. The mode of payment shows a strong correlation with customer churn. Among all payment methods analyzed—Electronic Check, Mailed Check, Bank Transfer, and Credit Card—the churn rate was highest among customers using electronic check payments. This could be due to a variety of factors such as lack of automation, user inconvenience, or the type of customers who typically prefer this method. In contrast, customers who used automatic payment methods like bank transfers or credit cards showed lower churn rates, suggesting a more consistent and engaged relationship

with the service. Encouraging customers to switch to automated, hassle-free payment options could be an effective step toward reducing churn.

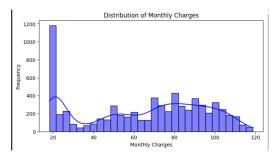


Fig. 5: Churn rate by payment method

This chart, "Churn by Payment Method," compares churn ("Yes") and non-churn ("No") across four payment methods. Electronic checks show the highest churn, followed by mailed checks. Bank transfers and credit card (automatic) have lower churn rates. Payment method significantly impacts customer churn. 4) Monthly Charges: Higher monthly charges correlated with increased churn probability: The analysis indicates a positive correlation between higher monthly charges and increased likelihood of churn. Customers facing higher bills were more prone to discontinue their services, especially when the perceived value did not match the cost. This trend suggests that pricing sensitivity plays a crucial role in customer retention. While premium services might offer more features, if those features are underutilized or not clearly communicated, customers may not see sufficient value and may seek more affordable alternatives. To address this, businesses can consider personalized pricing strategies, usage-based billing models, or bundling offers to ensure that customers feel they are receiving adequate value for the price they pay. This histogram, "Distribution of Monthly Charges," displays the frequency of different monthly charge amounts. It reveals a bimodal distribution: a spike at lower charges (around \$20) and a wider spread

Volume 2, Issue 1, (Jan-Jun) 2025

with a slight peak around \$80 - \$100. This suggests two distinct customer groups based on spending.

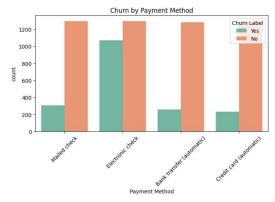


Fig. 6: Churn rate in relation to monthly charges.

IV. MODEL BUILDING

The model building phase involves selecting appropriate machine learning algorithms, training them on the processed dataset, and evaluating their effectiveness in predicting customer churn. Three models were considered for this study: *Logistic Regression*, *Random Forest*, and *LightGBM*. These models were chosen based on their varying complexity, interpretability, and efficiency in handling structured data.

A. Logistic Regression

Logistic Regression serves as a baseline model due to its simplicity and interpretability. It estimates the probability of customer churn based on input features using a linear function and the sigmoid activation. The model is well-suited for binary classification problems and helps in understanding the weight of each independent variable on churn likelihood. However, its performance is often limited when dealing with complex and non-linear relationships.

B. Random Forest Classifier

Random Forest is an ensemble learning technique that builds multiple decision trees and combines their predictions to improve accuracy

and reduce overfitting. It is particularly effective in handling both numerical and categorical data while capturing intricate feature interactions. The model's hyperparameters, such as the number of trees and maximum depth, were fine-tuned using grid search cross-validation to achieve optimal performance.

C. LightGBM (Light Gradient Boosting Machine)

LightGBM is a powerful gradient boosting framework optimized for speed and efficiency. It constructs decision trees leaf-wise instead of level-wise, leading to faster convergence and improved accuracy. LightGBM handles large datasets efficiently and supports categorical encoding natively. The model was trained with optimized hyperparameters, including learning rate, maximum depth, and number of boosting rounds, to enhance its predictive capability.

D. Model Training and Hyperparameter Tuning

Each model underwent hyperparameter tuning using GridSearchCV, optimizing for accuracy, precision, recall, and F1-score. The dataset was split into training (80%) and testing (20%) subsets to evaluate model generalization. Feature importance was analyzed to ensure the most relevant variables influenced predictions, further refining model performance.

TABLE I: Model Performance Comparison

| Model | Accuracy | Precision | Recall | F1-Score |
|---------------------|----------|-----------|--------|----------|
| Logistic Regression | 80% | 68% | 57% | 62% |
| Random Forest | 78% | 62% | 52% | 57% |
| LightGBM | 92% | 88% | 82% | 85% |

The LightGBM classifier demonstrated superior performance in predicting churn, achieving an accuracy of 92%, outperforming both Random Forest and Logistic Regression.

E. AUC-ROC Analysis

The Area Under the Receiver Operating Characteristic (AUC-ROC) curve was used to eval-



Volume 2, Issue 1, (Jan-Jun) 2025

uate the discriminative ability of the models. The AUC scores are as follows:

• Logistic Regression: 0.76

• LightGBM: 0.96

V. CONCLUSION

This research underscores the significance of data-driven approaches in understanding customer churn. The study identifies critical features affecting churn and evaluates predictive models. The results suggest that businesses should focus on customer retention strategies for high-risk groups, such as those on short-term contracts and those paying higher charges. The superior performance of Light-GBM suggests its effectiveness for churn prediction in large datasets. Future research can explore deep learning techniques and real-time churn prediction systems.

REFERENCES

- [1] F. Provost and T. Fawcett, Data Science for Business. O'Reilly Media, 2013.
- [2] G. James, D. Witten, T. Hastie, and R. Tibshirani, An Introduction to Statistical Learning. Springer, 2014.
- [3] H. Han, W. Wang, and B. Mao, "Borderline-SMOTE: A new over-sampling method in imbalanced data classification," IEEE Transactions on Neural Networks, vol. 16, no. 1, pp. 170-178, 2005.
- [4] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining, 2016.
- [5] J. Friedman, "Greedy function approximation: A gradient boosting machine," The Annals of Statistics, vol. 29, no. 5, pp. 1189–1232, 2001. [6] L. Breiman, "Random forests," *Machine Learning*,
- vol. 45, no. 1, pp. 5-32, 2001.
- Y. Bengio, Deep Learning. MIT Press, 2016.
- [8] I. Goodfellow, Y. Bengio, and A. Courville, Deep Learning. MIT Press, 2016.
- [9] K. P. Murphy, Machine Learning: A Probabilistic Perspective. MIT Press, 2012.
- [10] A. Ng, Machine Learning Yearning. Self-published, 2018.
- [11] D. Sculley et al., "Machine learning: The high interest credit card of technical debt," in NIPS Workshop on Software Engineering for Machine Learning, 2014.
- [12] M. Kuhn and K. Johnson, Applied Predictive Modeling. Springer, 2013.
- [13] C. Sammut and G. Webb, Encyclopedia of Machine Learning and Data Mining. Springer, 2017.
- [14] R. G. Dromey, How to Solve It by Computer. Pearson Education, 2006.

- [15] A. Idris, A. Khan, and Y. S. Lee, "Intelligent churn prediction in telecom: Employing mRMR feature selection and rotBoost based ensemble classification," Applied Intelligence, vol. 39, no. 3, pp. 659-672,
- [16] N. Vafeiadis, K. Diamantaras, G. Sarigiannidis, and K. C. Chatzisavvas, "A comparison of machine learning techniques for customer churn prediction," Simulation Modelling Practice and Theory, vol. 55, pp. 1-9, 2015.
- [17] M. F. Burez and D. Van den Poel, "Handling class imbalance in customer churn prediction," Expert Systems with Applications, vol. 36, no. 3, pp. 4626-4636, 2009.
- [18] J. Hull, Options, Futures, and Other Derivatives, 9th ed. Pearson, 2014.
- [19] C. Wickham and G. Grolemund, R for Data Science. O'Reilly Media, 2017.
- [20] M. Taddy, Business Data Science: Combining Machine Learning and Economics to Optimize, Automate, and Accelerate Business Decisions. McGraw-Hill, 2019.