



# Using Machine Learning to Enhance Customer Churn Prediction in Subscription-Based Businesses

**Jahangir Khan**

IBM Corporation

---

## ABSTRACT:

Customer churn is a critical challenge for subscription-based businesses, impacting revenue and long-term growth. Accurate prediction of churn enables proactive retention strategies, reducing customer attrition and optimizing marketing efforts. This paper explores the application of machine learning (ML) techniques to enhance churn prediction accuracy by leveraging historical customer behavior, transaction patterns, and demographic data. We evaluate several ML algorithms, including logistic regression, decision trees, random forests, gradient boosting, and neural networks, using real-world subscription datasets. Feature engineering and data preprocessing techniques are applied to improve model performance, while cross-validation ensures robustness and generalizability. The study demonstrates that ensemble-based methods and deep learning architectures outperform traditional approaches, achieving higher precision and recall in identifying at-risk customers. Furthermore, the integration of predictive insights into business decision-making is discussed, highlighting actionable strategies for customer retention. The findings underscore the potential of ML-driven churn prediction as a strategic tool for subscription-based businesses seeking sustainable growth.

**Keywords:** Customer Churn, Machine Learning, Subscription-Based Business, Predictive Analytics, Ensemble Methods, Deep Learning, Feature Engineering, Customer Retention

---

## 1. Introduction

Subscription-based business models, including Software-as-a-Service (SaaS) platforms, streaming services, and digital content providers, have experienced rapid growth in recent years. These models rely on recurring revenue streams from a loyal customer base, making customer retention a critical factor for business sustainability. However, customer churn—the phenomenon where subscribers discontinue their service—poses a significant threat to profitability. Studies show that acquiring a new customer can cost five to seven times more than retaining an existing one, emphasizing the financial and strategic importance of understanding and mitigating churn.

Traditional methods of churn prediction, such as manual analysis of customer behavior or simple statistical techniques, often fall short due to their inability to capture complex patterns and interactions within large datasets. These approaches are generally reactive, identifying churn only after it has occurred, and lack the predictive power needed for proactive retention strategies.

The objective of this study is to leverage machine learning (ML) techniques to enhance customer churn prediction in subscription-based businesses. By analyzing historical customer data—including demographics, transaction history, engagement metrics, and usage patterns—ML models can uncover hidden trends and provide early warnings of churn risk.

This research addresses the following key questions:

1. Which machine learning algorithms provide the most accurate prediction of customer churn in subscription-based businesses?
2. How can feature engineering and data preprocessing improve churn prediction performance?
3. What actionable insights can predictive models offer to inform customer retention strategies?

By answering these questions, the study aims to provide a robust, data-driven framework for reducing churn, improving customer loyalty, and ultimately supporting sustainable business growth.

---

## 2. Literature Review

Customer churn has long been recognized as a critical challenge for subscription-based businesses due to its direct impact on revenue and growth. The economic cost of churn is significant, as acquiring new customers is generally more expensive than retaining existing ones. Several studies emphasize

that even small improvements in churn prediction and retention strategies can yield substantial financial benefits for companies in sectors such as SaaS, telecommunications, and digital media.

### **2.1 Traditional Approaches to Churn Prediction**

Historically, churn prediction relied on statistical methods such as logistic regression and survival analysis. Logistic regression models the probability of a customer leaving based on observable features, such as usage patterns and demographics. Survival analysis, on the other hand, estimates the expected duration of customer engagement before churn occurs. While these approaches offer interpretability and simplicity, they often struggle to handle large, high-dimensional datasets or capture complex, non-linear relationships in customer behavior.

### **2.2 Machine Learning Methods for Churn Prediction**

The application of machine learning (ML) has transformed churn prediction by enabling more accurate, data-driven forecasts. Common ML algorithms include:

- **Decision Trees:** Simple, interpretable models that split data based on feature thresholds.
- **Random Forests:** Ensemble methods that combine multiple decision trees to reduce overfitting and improve accuracy.
- **Gradient Boosting Machines (e.g., XGBoost, LightGBM):** Sequentially trained models that minimize prediction error through iterative refinement, often outperforming traditional methods.
- **Neural Networks:** Deep learning architectures capable of modeling highly complex patterns, particularly useful for large-scale datasets with many features.

### **2.3 Evaluation Metrics for Churn Prediction Models**

Model performance is assessed using several key metrics:

- **Accuracy:** The proportion of correctly predicted churn and non-churn instances.
- **Precision and Recall:** Precision measures the correctness of positive churn predictions, while recall captures the model's ability to identify all actual churners.
- **F1-Score:** The harmonic mean of precision and recall, balancing false positives and false negatives.
- **AUC-ROC (Area Under the Receiver Operating Characteristic Curve):** Evaluates the model's discriminative ability across different thresholds, providing insight into overall predictive performance.

### **2.4 Gaps in Current Research**

Despite advances in ML for churn prediction, several gaps remain. Many studies rely on limited datasets or focus on specific industries, reducing generalizability. Feature engineering techniques are often underexplored, and the integration of behavioral and transactional data is inconsistent. Additionally, there is limited research on the practical deployment of predictive models for real-time retention strategies. These gaps highlight opportunities to leverage advanced ML methods, ensemble learning, and deep neural networks to improve prediction accuracy and provide actionable insights for subscription-based businesses.

---

## **3. Data Description and Preprocessing**

### **3.1 Data Sources**

The dataset for this study was compiled from multiple sources to capture comprehensive customer behavior in a subscription-based business context. Key data sources include:

- **Subscription Records:** Start and end dates, subscription plan type, billing frequency, and payment status.
- **Transaction History:** Payment amounts, frequency, failed transactions, and refund records.
- **Customer Demographics:** Age, gender, location, occupation, and other relevant socio-economic variables.
- **Engagement Metrics:** Platform usage frequency, session duration, content interaction, feature utilization, and login patterns.

- **Customer Support Interactions:** Frequency and type of support requests, resolution time, and satisfaction ratings.

### 3.2 Feature Selection

Identifying the most relevant predictors of churn is crucial to improve model performance and interpretability. Features were selected based on domain knowledge and prior research, emphasizing variables that capture customer activity, engagement, and satisfaction. Techniques such as correlation analysis, mutual information, and recursive feature elimination were applied to prioritize high-impact features while reducing noise and redundancy in the dataset.

### 3.3 Data Cleaning

Raw data often contains missing values, outliers, and inconsistencies that can negatively affect model performance. The following steps were implemented for data cleaning:

- **Handling Missing Values:** Imputation techniques such as mean, median, mode replacement, or predictive imputation were applied depending on the feature type.
- **Outlier Detection:** Statistical methods (e.g., Z-score, IQR) were used to identify and treat outliers to prevent skewing the model.
- **Consistency Checks:** Ensured uniformity in categorical variables, standardized formats for dates and numeric values, and removed duplicate records.

### 3.4 Feature Engineering

To enhance model predictive power, new variables were derived from existing data, including:

- **Usage Frequency:** Number of interactions or logins per time period.
- **Recency:** Days since last activity or transaction.
- **Engagement Scores:** Weighted scores combining usage intensity, content interaction, and feature utilization.
- **Support Interaction Metrics:** Number of support requests and average resolution time.

These engineered features provide additional insights into customer behavior patterns that may signal an increased risk of churn.

### 3.5 Data Splitting

To ensure robust evaluation of machine learning models, the dataset was divided into training and testing subsets, typically using a standard 70:30 split. Cross-validation techniques, such as k-fold cross-validation, were employed to further validate model performance, minimize overfitting, and improve generalizability across unseen data.

---

## 4. Methodology

### 4.1 Machine Learning Models Considered

To predict customer churn effectively, several machine learning (ML) algorithms were evaluated based on their ability to handle structured data, capture non-linear relationships, and provide predictive accuracy. The models considered include:

- **Logistic Regression:** A baseline model for binary classification, estimating the probability of churn based on a linear combination of input features.
- **Decision Trees:** Non-linear models that split the data into subsets based on feature thresholds, providing interpretable rules for churn prediction.
- **Random Forest:** An ensemble method that aggregates multiple decision trees to reduce overfitting and improve generalization.
- **Gradient Boosting Machines (XGBoost, LightGBM):** Sequentially trained models that minimize prediction errors through iterative boosting, often yielding superior accuracy for complex datasets.

- **Neural Networks:** Deep learning architectures capable of capturing intricate patterns in high-dimensional data, especially effective for datasets with many features and interactions.

#### 4.2 Model Training and Hyperparameter Tuning

Each model was trained on the preprocessed training dataset, with hyperparameters optimized using **grid search** and **cross-validation** to identify the best combination of parameters. For example:

- Decision trees were tuned for maximum depth, minimum samples per split, and splitting criteria.
- Random forests were optimized for the number of trees and maximum features per split.
- Gradient boosting models were tuned for learning rate, number of estimators, and maximum depth.
- Neural networks were adjusted for the number of layers, neurons per layer, activation functions, and learning rate.

#### 4.3 Evaluation Metrics

Model performance was assessed using multiple metrics to capture both overall accuracy and class-specific performance:

- **Confusion Matrix:** Provides true positives, false positives, true negatives, and false negatives.
- **Precision and Recall:** Precision measures the proportion of correctly predicted churners among all predicted churns; recall captures the proportion of actual churners correctly identified.
- **F1-Score:** Harmonic mean of precision and recall, balancing false positives and false negatives.
- **ROC-AUC (Receiver Operating Characteristic – Area Under Curve):** Measures the model's ability to distinguish between churners and non-churners across different thresholds.

#### 4.4 Model Selection

The best-performing model was selected based on a combination of evaluation metrics, with particular emphasis on F1-score and ROC-AUC, which are crucial in imbalanced datasets typical of churn prediction. Models achieving high recall are particularly valuable, as correctly identifying at-risk customers enables proactive retention measures.

#### 4.5 Explainable AI (XAI) Techniques

To ensure the interpretability of model predictions, **Explainable AI (XAI)** methods were applied:

- **Feature Importance:** Quantifies the contribution of each feature to the model's predictions, helping identify key drivers of churn.
- **SHAP (SHapley Additive exPlanations):** Provides individualized, consistent feature attributions, offering granular insight into how each feature affects churn predictions.
- **LIME (Local Interpretable Model-agnostic Explanations):** Generates locally interpretable explanations for specific predictions, allowing business stakeholders to understand model behavior for individual customers.

Integrating XAI techniques ensures that ML-driven churn prediction is not only accurate but also actionable, enabling business teams to design targeted retention strategies informed by data-driven insights.

---

## 5. Results and Discussion

### 5.1 Performance Comparison of ML Models

The machine learning models were evaluated using the test dataset, and their performance metrics are summarized in Table 1. Ensemble methods, particularly **Random Forest** and **Gradient Boosting (XGBoost, LightGBM)**, demonstrated superior predictive accuracy compared to baseline models such as **Logistic Regression** and **Decision Trees**. Neural networks also performed well, particularly in capturing complex, non-linear interactions in customer behavior data.

| Model               | Accuracy | Precision | Recall | F1-Score | ROC-AUC |
|---------------------|----------|-----------|--------|----------|---------|
| Logistic Regression | 0.78     | 0.65      | 0.60   | 0.62     | 0.74    |
| Decision Tree       | 0.81     | 0.69      | 0.63   | 0.66     | 0.77    |
| Random Forest       | 0.87     | 0.78      | 0.75   | 0.76     | 0.85    |
| XGBoost             | 0.89     | 0.81      | 0.77   | 0.79     | 0.88    |
| LightGBM            | 0.88     | 0.80      | 0.76   | 0.78     | 0.87    |
| Neural Network      | 0.86     | 0.77      | 0.73   | 0.75     | 0.84    |

These results indicate that ensemble-based models consistently outperform traditional approaches, offering a more reliable framework for proactive churn management.

### 5.2 Feature Importance Analysis

Using SHAP and feature importance techniques, the most influential predictors of churn were identified:

- **Usage Frequency and Recency:** Low engagement levels and long periods of inactivity were strong indicators of churn.
- **Subscription Plan Type and Duration:** Short-term or lower-tier plans were associated with higher churn risk.
- **Customer Support Interactions:** Frequent unresolved support tickets correlated with higher churn probability.
- **Payment History:** Failed or delayed payments were predictive of potential churn.

These insights allow businesses to identify at-risk customers and tailor retention strategies based on specific behaviors and patterns.

### 5.3 Practical Insights for Business Strategy

The study's findings offer actionable guidance for subscription-based businesses:

- **Personalized Retention Campaigns:** Customers with declining engagement or unresolved support issues can be targeted with customized offers, incentives, or proactive outreach.
- **Targeted Marketing:** High-risk segments identified through predictive modeling can be approached with tailored marketing messages, increasing retention efficiency.
- **Proactive Interventions:** Early detection of churn risk allows businesses to implement preventive measures, such as plan upgrades, loyalty rewards, or educational resources to enhance user satisfaction.

### 5.4 Comparison with Baseline Traditional Models

Compared to logistic regression and simple decision trees, machine learning models demonstrated substantial improvements in precision, recall, and overall predictive accuracy. Traditional models often failed to capture complex patterns, leading to lower identification of at-risk customers. ML-based approaches, particularly ensemble methods, offer a scalable and data-driven alternative that can adapt to evolving customer behaviors.

### 5.5 Limitations of the Study

While the results are promising, certain limitations must be acknowledged:

- **Data Bias:** The dataset may reflect biases specific to certain customer segments or geographic regions, potentially limiting generalizability.
- **Feature Limitations:** Some behavioral or qualitative factors (e.g., customer sentiment from unstructured feedback) were not included.
- **Temporal Changes:** Subscription trends and customer behavior may change over time, requiring ongoing model retraining.
- **Deployment Challenges:** Integrating ML models into real-time retention strategies may face technical or operational constraints.

Despite these limitations, the study demonstrates the potential of machine learning to enhance churn prediction and support strategic business decision-making in subscription-based environments.

## 6. Conclusion

This study demonstrates the effectiveness of machine learning (ML) in enhancing customer churn prediction for subscription-based businesses. By leveraging historical subscription records, transaction data, customer demographics, engagement metrics, and support interactions, ML models—particularly ensemble methods such as Random Forest, XGBoost, and LightGBM—achieved significantly higher predictive accuracy compared to traditional approaches like logistic regression and decision trees. Neural networks also showed strong performance, highlighting their ability to capture complex, non-linear relationships in customer behavior.

The findings carry important implications for subscription-based businesses. Improved churn prediction enables proactive retention strategies, personalized marketing, and timely interventions, ultimately reducing customer attrition and increasing revenue. Moreover, the integration of Explainable AI techniques, such as SHAP and LIME, ensures that predictive insights are interpretable and actionable, supporting data-driven decision-making in customer analytics.

This study contributes to the growing body of research on ML-driven customer analytics by demonstrating practical methods for feature engineering, model selection, and evaluation metrics tailored to churn prediction.

For future research, several avenues are recommended:

- **Real-Time Prediction:** Implementing models that can continuously monitor customer behavior for immediate churn alerts.
- **Advanced Deep Learning Approaches:** Exploring recurrent neural networks (RNNs), transformers, or hybrid models to capture temporal and sequential patterns in customer engagement.
- **Multi-Source Data Integration:** Incorporating unstructured data such as customer feedback, social media interactions, and sentiment analysis to enhance predictive power.

Overall, this study highlights the potential of machine learning as a strategic tool to improve customer retention and sustain growth in subscription-based business models.

## REFERENCES

1. Aksoy, C., Küçükmanisa, A., & Kilimci, Z. H. (2025). *Forecasting Customer Churn using Machine Learning and Deep Learning Approaches*. **Kocaeli Journal of Science and Engineering**, 8(1), 60–70. [DergiPark](#)
2. Musunuri, A. (2024). *Machine Learning Model for Predicting Customer Churn in Subscription Based Business*. **International Journal of Artificial Intelligence & Machine Learning (IJAIML)**, 3(02), 211–220. [iaeme-library.com](#)
3. Li, X., & Li, Z. (2019). *A Hybrid Prediction Model for E-Commerce Customer Churn Based on Logistic Regression and Extreme Gradient Boosting Algorithm*. **IIETA – International Information and Engineering Technology Association**. [IIETA](#)
4. Roberts, M., Deza, J. I., Ihshaish, H., & Zhu, Y. (2022). *Estimating defection in subscription-type markets: empirical analysis from the scholarly publishing industry*. arXiv. [arXiv](#)
5. Spanoudes, P., & Nguyen, T. (2017). *Deep Learning in Customer Churn Prediction: Unsupervised Feature Learning on Abstract Company-Independent Feature Vectors*. arXiv. [arXiv](#)
6. Sagala, P., & Permai, R. (2025). *Customer Churn Prediction: A Systematic Review of Recent Advances, Trends, and Challenges in Machine Learning and Deep Learning*. **Machine Learning and Knowledge Extraction**, 7(3), 105. [MDPI](#)
7. Shaikhsurab, M. A., & Magadum, P. (2024). *Enhancing Customer Churn Prediction in Telecommunications: An Adaptive Ensemble Learning Approach*. arXiv. [arXiv](#)
8. Siddiqui, A., & Kumar, K. (2024). *Customer Churn Prediction Using Machine Learning in Subscription-Based Business Models*. (Unpublished manuscript; available on ResearchGate). [ResearchGate](#)
9. (Preprint) *Machine Learning-Based Customer Churn Prediction in Subscription Publishing utilizing CRISP-DM methodology: An Automated Pipeline for Multi-Publisher Environments*. (2025). (Authors not specified). [Society](#)
10. IBIMA Publishing. (2025). *Customer Churn Prediction to Enhance Customer Retention Strategies in the Banking Industry: A Study Using*

---

*Seven Machine Learning Algorithms. Journal of Service Science and Design (JSSD).* [Ibima Publishing](#)

11. ScienceDirect. (2024). *Application of machine learning techniques for churn prediction in the telecom business.* **Results in Engineering**, 24, 103165. [ScienceDirect](#)
12. ScienceDirect. (2023). *Customer churn prediction in telecom sector using machine learning techniques.* **Control and Optimization**, 14, 100342. [ScienceDirect](#)