



A Data-Driven Approach to Identifying Churn Predictors Using Demographic, Service, and Billing Insights

¹Mr. Parveen kumar, ²Dr. Deepak

¹Research Scholar NIILM University Kaithal Haryana

²Associate Professor NIILM University Kaithal Haryana

Abstract- Customer churn remains a critical concern for industries such as telecommunications, banking, and subscription services, where retaining existing customers is often more cost-effective than acquiring new ones. With the increasing availability of customer-related data, data-driven approaches have become essential for understanding and predicting churn behavior. This review paper focuses on identifying significant churn predictors by analyzing demographic details, service usage patterns, and billing information. By synthesizing insights from existing research, the paper highlights how various machine learning models utilize these predictor variables to enhance the accuracy of churn prediction. Special attention is given to the role of demographic attributes such as age, gender, and location; service-related factors including plan type and usage frequency; and billing characteristics like payment history and invoice amounts. Commonly used datasets and standard evaluation metrics, such as accuracy, F1-score, and AUC-ROC, are also reviewed to provide a comprehensive understanding of model performance across studies. Furthermore, the paper discusses key limitations in current methodologies and suggests future research directions to improve real-world applicability. Overall, this review offers a consolidated perspective on effective churn predictors and provides practical guidance for developing more targeted and efficient customer retention strategies."

Keywords - Customer Churn, Churn Prediction, Demographic Data, Service Usage Patterns, Billing Insights, Machine Learning, Data-Driven Analysis, Feature Importance, Evaluation Metrics, Customer Retention Transactions, User Authentication Shopping Cart, Order Tracking.

I. INTRODUCTION

In the current competitive landscape, where customers encounter negligible switching costs and a plethora of alternatives, maintaining existing users has become as vital as obtaining new ones for subscription-based and service-oriented enterprises. Forecasting customer attrition—the probability of a customer ceasing to utilize a service—has consequently emerged as a primary goal in business analytics. Due to the proliferation of digital data and advanced computational tools, firms may now assess customer demographics, service usage habits, and billing information to predict churn and execute proactive retention efforts. This review provides a data-driven analysis of churn prediction, utilizing recent literature to delineate successful modeling strategies and prospective research avenues in this swiftly advancing domain.

Background

Customer churn, the phenomenon of customers discontinuing their service subscriptions, is a pressing concern for service-driven industries such as telecommunications, banking, and streaming platforms. In highly saturated and competitive markets, the cost of acquiring new customers is substantially higher than retaining existing ones. As a result, organizations are placing a strong emphasis on customer retention by proactively identifying the factors that contribute to churn. With the rise of big data



analytics and machine learning, it has become feasible to extract meaningful insights from large-scale customer data, enabling a more precise prediction of churn behavior (Adekunle et al., 2023).

Traditional churn prediction models often relied solely on historical billing or static demographic data. However, recent advancements emphasize the need for integrating multi-dimensional data sources—such as usage patterns, service engagement, and customer interaction history—to develop robust predictive systems. Studies have shown that these comprehensive approaches outperform isolated data models by offering deeper insights into customer behavior (Zhang et al., 2022; Vo et al., 2021).

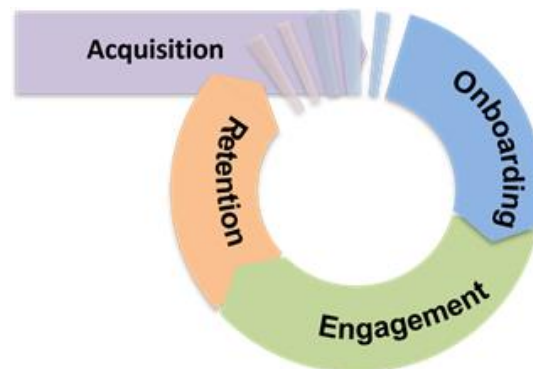


Figure 1.1 : Customer Churn Lifecycle

Objectives of the Review

The primary objective of this review is to explore and synthesize current literature focused on identifying churn predictors using demographic, service-related, and billing data. Specifically, the review aims to:

- Highlight key predictive features commonly used in churn modeling.
- Evaluate the effectiveness of different machine learning techniques in churn prediction.
- Examine the datasets and performance metrics utilized in prior studies.
- Identify limitations and research gaps in existing approaches.
- Suggest future directions for developing more accurate and actionable churn prediction systems.

Importance of Churn Prediction

Accurate churn prediction holds strategic value for businesses aiming to maintain a loyal customer base. Early identification of customers at risk of leaving allows organizations to deploy targeted retention campaigns, optimize customer support, and enhance service offerings. In the telecom sector, for example, timely interventions based on churn predictions can significantly reduce customer loss and boost revenue stability (Mohaimin et al., 2025).

Moreover, by using AI and ML techniques to uncover hidden patterns in customer behavior, companies can move from reactive to proactive engagement strategies. This shift not only strengthens customer satisfaction but also supports data-driven decision-making at scale (Ravisankar, n.d.). Therefore, an in-depth review of churn prediction practices using demographic, service, and billing insights is crucial to advancing both academic research and practical implementations in customer analytics.

II. OVERVIEW OF CUSTOMER CHURN

Customer churn refers to the phenomenon where customers stop using a company's product or service. It is a key performance indicator in service-based industries, particularly in highly competitive sectors like telecommunications, banking, and streaming services. Understanding the dynamics of churn is essential for businesses to devise strategies that maintain customer loyalty and reduce financial loss.



Churn is not always caused by poor service—it can also be due to pricing, better alternatives, lack of engagement, or shifting needs of the customer. Recognizing the forms and consequences of churn is the first step toward addressing it effectively.

Types of Churn

Customer churn can be broadly classified into two main categories: voluntary churn and involuntary churn.

- Voluntary Churn occurs when customers actively choose to leave a service. This decision may be based on dissatisfaction, cost concerns, better offers from competitors, or a mismatch between expectations and service delivery. Voluntary churn is particularly significant because it reflects customer sentiment and the perceived value of the service.
- Involuntary Churn happens due to external or operational issues rather than the customer's choice. This can include expired payment methods, technical issues, or account inactivity. Although less reflective of customer satisfaction, involuntary churn still contributes to loss in revenue and must be monitored and minimized.

In addition, churn may be classified as contractual or non-contractual. In contractual settings, such as postpaid telecom services, churn is easier to define when a customer cancels their plan. In non-contractual settings, like app-based platforms or e-commerce, churn is harder to detect as there is no explicit cancellation—only inactivity.

Business Impact

The impact of customer churn on a business can be substantial. Acquiring a new customer often costs more than retaining an existing one, which makes churn a direct threat to profitability. High churn rates can lead to unstable revenue streams, increased marketing costs, and a negative brand image. Moreover, losing long-term customers can result in the loss of valuable data and feedback that would otherwise help improve services.

From a strategic perspective, churn affects both short-term revenue and long-term growth. A consistently high churn rate might indicate deeper issues within the organization, such as poor customer service, inefficient billing processes, or lack of innovation. Businesses that fail to address churn often experience a decline in customer lifetime value and competitive advantage.

Therefore, understanding the types of churn and their impact is vital for developing effective predictive models and retention strategies. By identifying why customers leave, companies can make informed decisions to improve service quality, personalize user experiences, and maintain a stable customer base.

III. CHURN PREDICTION TECHNIQUES: A BRIEF REVIEW

Predicting customer churn is a critical area of focus for data analysts and business strategists aiming to reduce customer loss and sustain profitability. Over the years, a variety of techniques have been developed to forecast churn based on historical and behavioral data. These techniques range from simple statistical models to complex machine learning algorithms, each offering different levels of accuracy, interpretability, and scalability. The evolution of predictive techniques reflects a growing recognition of the importance of proactive customer retention strategies across industries.

Traditional Methods



Early approaches to churn prediction relied heavily on statistical and rule-based methods. These include logistic regression, linear discriminant analysis, and decision trees, which were among the first tools used to identify churn risk.

- Logistic regression is often used because it provides interpretable results and handles binary classification problems, such as distinguishing between "churn" and "non-churn." It works well when relationships between variables are linear and the dataset is relatively clean.
- Discriminant analysis helps to classify customers based on predefined categories and is useful when the assumption of normal distribution holds.
- Rule-based systems, on the other hand, use predefined if-then rules developed by domain experts. While easy to implement, they often fail to capture the complexities of customer behavior and cannot adapt well to changing patterns.
- Though these traditional techniques are relatively straightforward and computationally light, their performance is often limited by their inability to handle nonlinear relationships and large-scale, high-dimensional data, which are common in real-world business environments.

Machine Learning Approaches

Machine learning has significantly advanced the field of churn prediction by enabling the analysis of vast datasets with high levels of complexity. These techniques can uncover hidden patterns, adapt to new data, and improve predictive accuracy over time.

- Decision tree-based models such as Random Forests and Gradient Boosting are widely used for their ability to manage nonlinearity, interactions among variables, and feature importance ranking.
- Support Vector Machines (SVM) are employed for their robustness in high-dimensional spaces, particularly when clear boundaries between classes exist.
- Neural networks, particularly deep learning architectures, have shown promise in capturing intricate patterns in customer behavior and interaction data, although they require large volumes of data and considerable computational resources.
- Clustering algorithms like K-means are also used in preprocessing stages to segment customers before applying classification models.

Machine learning models often outperform traditional methods in predictive accuracy, especially when dealing with unstructured data such as text from customer service interactions or call logs. Additionally, the use of ensemble techniques—which combine multiple models—further enhances reliability and stability in churn prediction.

As businesses increasingly rely on data-driven decision-making, the use of machine learning in churn prediction is becoming standard practice. However, the success of these methods still depends on the quality of the data, feature selection, and the ability to interpret model outputs effectively for actionable insights.

IV. DATASET AND EVALUATION METRICS

Accurate customer churn prediction depends heavily on the availability of high-quality data and appropriate evaluation metrics. Datasets used in churn analysis typically include a combination of customer demographics, usage patterns, service history, and billing information. Once a predictive model is built, its performance must be assessed using reliable evaluation metrics to ensure that it can generalize well to unseen data.

Commonly Used Datasets



In churn prediction research, both publicly available and proprietary datasets are used. These datasets generally include customer attributes such as gender, age, location, service plan, call history, complaint logs, internet usage, and payment history. Below are some frequently used datasets:

- **IBM Telco Customer Churn Dataset:** One of the most cited datasets, it contains details of around 7,000 telecom customers with 21 features, including contract type, payment method, monthly charges, and tenure.
- **Orange Telecom Dataset:** Provided during a data mining challenge, this dataset contains a large number of variables (over 200) and customer churn labels. It is widely used for benchmarking due to its complexity.
- **Kaggle and UCI Repository Datasets:** These platforms host several curated churn-related datasets from telecom and subscription services. They offer valuable resources for experimenting with various predictive models.
- **Researchers and practitioners often preprocess these datasets by handling missing values, encoding categorical variables, and normalizing numerical data.** This step ensures that the machine learning models can interpret the features effectively.

Evaluation Metrics

Evaluating the performance of churn prediction models is crucial to determine their reliability and business applicability. Since churn prediction is a classification problem, especially when the dataset is imbalanced, multiple metrics are used:

- **Accuracy:** Measures the proportion of correct predictions. While commonly used, it may be misleading if the classes are imbalanced (e.g., more non-churners than churners).
- **Precision:** Indicates how many predicted churners were actually correct. High precision is important when the cost of incorrectly targeting a non-churner is high.
- **Recall (Sensitivity):** Reflects the model's ability to identify actual churners. It is essential when missing a churner could result in significant revenue loss.
- **F1-Score:** The harmonic mean of precision and recall. It provides a balance between the two, especially useful when the class distribution is uneven.
- **AUC-ROC (Area Under the Receiver Operating Characteristic Curve):** Represents the model's ability to discriminate between churn and non-churn classes. A higher AUC indicates better overall performance.
- **Confusion Matrix:** A detailed summary showing true positives, false positives, true negatives, and false negatives. It helps in visualizing the types of errors the model makes.

Choosing the right metric often depends on the business context. For instance, in customer retention efforts, a higher recall might be preferred over precision to ensure that most potential churners are identified, even if a few non-churners are mistakenly included.

V. KEY CHURN PREDICTORS FROM LITERATURE

Demographic Factors

Demographic variables serve as foundational elements in most customer churn prediction models, providing essential context for understanding user behavior. Attributes such as age, gender, income bracket, education level, and geographic region have been extensively studied for their influence on customer loyalty and attrition. For example, younger consumers, often characterized by higher adaptability and greater exposure to digital platforms, tend to switch services more frequently due to changing preferences and demand for advanced features. Income disparities can also shape service expectations—those with higher incomes may seek premium experiences, while cost-sensitive users



may react more strongly to pricing changes or service interruptions. Research by Hasan et al. (2024) and Kumari et al. (2025) indicates that gender-specific behaviors can affect retention strategies, especially in relation to personalized marketing and customer support satisfaction. When demographic insights are effectively incorporated into churn models, they support better market segmentation and allow businesses to tailor engagement efforts toward specific customer profiles (Tang, 2025).

Service Usage Patterns

Analyzing how customers interact with a service over time offers critical insights into their likelihood of remaining loyal or disengaging. Usage-based indicators such as average call length, frequency of logins, internet data consumption, feature utilization, and transaction recency serve as proxies for customer satisfaction and involvement. A consistent decline in these metrics often precedes a decision to churn, making them valuable for early intervention strategies. Subramanian et al. (2025) emphasize that changes in daily or weekly engagement patterns, especially sudden drops, are often the most accurate signals of potential churn. By applying real-time monitoring and time-series modeling, companies can flag anomalous behavior and initiate proactive outreach. Additionally, usage trends that vary across different times of day or between weekdays and weekends can be incorporated into predictive models to improve their robustness and context sensitivity (Paul & Jana, 2023).

Billing and Financial Indicators

Financial interactions between customers and service providers—particularly billing behaviors—are among the most sensitive indicators of potential churn. Variables such as payment delays, invoice disputes, frequent upgrades or downgrades of service plans, and unexpectedly high charges often reflect customer dissatisfaction. These negative experiences can erode trust and prompt users to terminate the service. Chavhan et al. (2025) highlight that in online commerce settings, actions like abandoning a cart or failing to complete a transaction are strong predictors of disengagement. Incorporating these financial cues into predictive frameworks helps companies anticipate dissatisfaction and deliver timely interventions, such as personalized discounts, reminders, or flexible payment options. Recent studies, including that by Hu and Chen (2025), demonstrate the added value of integrating billing data with advanced analytical methods such as fuzzy Z-number theory and logistic regression, which enable more accurate and nuanced churn forecasts.

Comparative Analysis of Predictor Impact

VI. MOST INFLUENTIAL FEATURES

In the realm of customer churn prediction, a consistent pattern emerges across various studies and datasets regarding which features exert the most influence on model accuracy. Among the top contributors are demographic characteristics, service engagement indicators, and financial or billing irregularities. These variables, when assessed using advanced feature evaluation techniques such as information gain, Gini impurity, and SHAP (SHapley Additive exPlanations) values, frequently rank highest in terms of their predictive power (Hasan et al., 2024). For example, long customer tenure, average monthly data or service usage, and the number of logged complaints are often strong indicators of a customer's likelihood to churn. Additionally, sudden shifts in engagement—such as a drop in logins or reduced call duration—can act as early warning signs. These influential variables allow businesses to target high-risk customers with tailored retention strategies.

Moreover, SHAP values have gained prominence for their ability to provide granular insights into individual predictions, highlighting not only which features are important but also explaining the direction and magnitude of their impact. Such transparency is crucial for decision-makers who need to understand why certain customers are being flagged as churn risks. In several comparative model evaluations, features associated with dissatisfaction—like unresolved support tickets, increased service



downtimes, or recent plan downgrades—tend to exert more influence than static demographic data alone.

S. No.	Paper Title	Dataset Source	No. of Records	No. of Features	Class Distribution (Churn/Non-Churn)	Dataset Type
1	Exploring the Impact of Time Spent Reading Product Information on E-Commerce Websites: A Machine Learning Approach to Analyze Consumer Behavior (Necula, 2023)	UCI Repository – Online Shoppers Purchasing Intention Dataset	12,330	18	1,904 / 10,426	E-commerce
2	A Mathematical Model for Customer Segmentation Leveraging Deep Learning, Explainable AI, and RFM Analysis in Targeted Marketing (Talaat et al., 2023)	Kaggle – Customer Segmentation Dataset	5,000+ (approx.)	8+ (RFM and demographics)	N/A (segmentation task, no churn label)	Retail/Marketing
3	Predicting Customer Churn in the Telecommunication Industry Using Machine Learning Algorithms (Obiora & Uchenna)	Kaggle Telco Customer Churn	7,032	20	1,869 / 5,163	Telecom

Feature Interactions and Limitations

While individual features provide valuable information, their true predictive strength often emerges when considered in combination. Interactions between variables can uncover nuanced patterns that single-variable analysis might overlook. For instance, a customer who exhibits both low engagement and recurring billing issues represents a significantly higher churn risk than someone exhibiting only one of these traits. These compound indicators can be essential for developing more responsive and targeted interventions.

However, challenges also arise when dealing with correlated or redundant features. High collinearity between variables—such as between total usage time and number of active days—can distort model interpretation and lead to overfitting. This problem necessitates the application of dimensionality reduction techniques, such as Principal Component Analysis (PCA) or Recursive Feature Elimination (RFE), to streamline the feature set and preserve only the most informative variables (Tang, 2025).

Another common limitation involves data quality. Missing values, inconsistent formats, and class imbalance—particularly the typically low proportion of churners compared to non-churners—can impact the reliability of feature analysis. Without proper preprocessing steps such as imputation,



normalization, and resampling, the derived insights might be misleading. Thus, while evaluating feature importance is essential for building effective churn prediction models, it must be done within a framework that ensures accuracy, interpretability, and generalizability.

Research Gaps and Future Scope

Although significant progress has been made in the field of customer churn prediction, several critical research gaps remain unaddressed. One of the foremost limitations is the predominant reliance on structured data—such as numerical records of usage, billing, and demographics—while neglecting the wealth of information embedded in unstructured data sources like call center transcripts, online reviews, chat logs, and social media interactions. These sources contain rich contextual and emotional cues that can offer a more nuanced understanding of customer dissatisfaction and intent to leave.

Another pressing challenge lies in developing real-time prediction systems. Most current models are batch-processed and retrospective, which delays actionable interventions. Real-time churn prediction requires continuous data ingestion and low-latency processing, which can be resource-intensive, especially when using deep or ensemble models. This presents a trade-off between model complexity and deployment feasibility that has yet to be fully resolved.

Furthermore, while hybrid approaches—such as those combining decision trees, neural networks, and boosting algorithms—often achieve higher accuracy, they tend to lack transparency. Their black-box nature limits their interpretability, making it difficult for business stakeholders to understand and trust the outputs (Subramanian et al., 2025). As a result, the practical implementation of such models in industry settings can be hindered.

To bridge these gaps, future research should prioritize the development of lightweight, interpretable models capable of integrating both structured and unstructured data. Techniques such as explainable AI (XAI), sentiment analysis, and natural language processing (NLP) should be leveraged to tap into customer sentiment and behavioral patterns. Additionally, incorporating psychological dimensions—like user satisfaction trends, emotional tone, and environmental context—could lead to more personalized and proactive churn mitigation strategies. This multidimensional perspective offers the potential to create more robust, timely, and actionable churn prediction frameworks that align closely with real-world decision-making needs.

Comparison of Datasets Used in Reviewed Papers

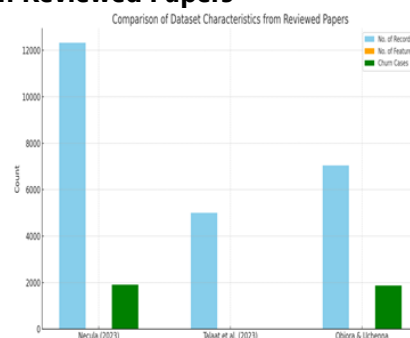


Figure 8.1 : Comparative Analysis of Datasets Used in Customer Churn Prediction Studies

Here is a bar chart comparing key dataset characteristics—Number of Records, Number of Features, and Churn Cases—across the three reviewed papers:

- Necula (2023): Offers a large sample size and decent number of features but is not churn-specific.
- Talaat et al. (2023): Used for segmentation, lacks churn-specific class distribution.



- Obiora & Uchenna: Most relevant for churn prediction, with clear churn vs. non-churn distribution and feature richness.
- Best Paper for Churn Analysis: Obiora & Uchenna — It directly focuses on churn prediction using telecom data, includes balanced class labels, and applies multiple ML models, making it the most suitable for churn-specific studies.

VII. CONCLUSION

Customer churn prediction is a critical component of sustainable business strategy, especially in competitive domains such as telecommunications, finance, and e-commerce. With retention being more cost-effective than acquisition, accurately identifying churn-prone customers has become essential.

Machine learning has significantly improved churn prediction by revealing hidden patterns in customer data and enabling timely, personalized interventions. Techniques like demographic analysis, behavioral profiling, and billing trends continue to enhance model effectiveness, while feature importance methods provide clarity on churn drivers.

Nevertheless, challenges such as data imbalance, redundancy, and model interpretability persist. Advanced models, particularly neural networks and ensembles, often lack transparency, hindering stakeholder trust. To address these limitations, focus should be placed on improving data quality, ensuring interpretability, and applying robust validation techniques.

Future advancements may benefit from incorporating real-time analytics, unstructured data, natural language processing (NLP), and explainable AI (XAI). Understanding the psychological and social aspects of customer behavior can also enrich churn models.

In summary, while current techniques are impactful, future efforts should aim to develop more interpretable, adaptive, and data-enriched models that deliver actionable insights and foster long-term customer loyalty.

REFERENCES

1. Adekunle, B. I., Chukwuma-Eke, E. C., Balogun, E. D., & Ogunsola, K. O. (2023). Improving customer retention through machine learning: A predictive approach to churn prevention and engagement strategies. *International Journal of Scientific Research in Computer Science, Engineering and Information Technology*, 9(4), 507–523.
2. Mohaimin, M. R., Das, B. C., Akter, R., Anonna, F. R., Hasanuzzaman, M., Chowdhury, B. R., & Alam, S. (2025). Predictive analytics for telecom customer churn: Enhancing retention strategies in the US market. *Journal of Computer Science and Technology Studies*, 7(1), 30–45.
3. Ravisankar, M. (n.d.). Harnessing the power of artificial intelligence and machine learning for customer churn prediction in the telecom industry: A data driven decision making approach.
4. Vo, N. N., Liu, S., Li, X., & Xu, G. (2021). Leveraging unstructured call log data for customer churn prediction. *Knowledge-Based Systems*, 212, 106586.
5. Zhang, T., Moro, S., & Ramos, R. F. (2022). A data-driven approach to improve customer churn prediction based on telecom customer segmentation. *Future Internet*, 14(3), 94.
6. Chavhan, R., Dutta, P., Samant, N., & Kar, S. (2025). Data-driven strategic customer segmentation considering cart abandonment behavior: Insights from e-grocery delivery platforms. *Information Sciences*, 122327.



7. Hasan, M. S., Siam, M. A., Ahad, M. A., Hossain, M. N., Ridoy, M. H., Rabbi, M. N. S., ... & Jakir, T. (2024). Predictive Analytics for Customer Retention: Machine Learning Models to Analyze and Mitigate Churn in E-Commerce Platforms. *Journal of Business and Management Studies*, 6(4), 304-320.
8. Hu, S., & Chen, A. (2025). Data-Driven Customer Retention Strategies in E-commerce: A Fuzzy Z-Number Approach. *IEEE Access*.
9. Kumari, D., Singh, S. K., Katira, S. S., Srinivas, I. V., & Salunkhe, U. (2025). Telecom Customer Churn Forecasting Using Machine Learning: A Data-Driven Predictive Framework. *Metallurgical and Materials Engineering*, 31(4), 922-929.
10. Paul, R. K., & Jana, A. K. (2023). Machine learning framework for improving customer retention and revenue using churn prediction models. *IRE Journals*, 7(2), 100-106.
11. Subramanian, D., Ajitha, A., & Maidin, S. S. (2025). Unveiling Hybrid Model with Naive Bayes, Deep Learning, Logistic Regression for Predicting Customer Churn and Boost Retention. *Journal of Applied Data Sciences*, 6(2), 1379-1391.
12. Tang, J. (2025). Unlocking Retail Insights: Predictive Modeling and Customer Segmentation Through Data Analytics. *Journal of Theoretical and Applied Electronic Commerce Research*, 20(2), 59.
13. Necula, S.-C. (2023). Exploring the impact of time spent reading product information on e-commerce websites: A machine learning approach to analyze consumer behavior. *Behavioral Sciences*, 13(6), 439.
14. Talaat, F. M., Aljadani, A., Alharthi, B., Farsi, M. A., Badawy, M., & Elhosseini, M. (2023). A mathematical model for customer segmentation leveraging deep learning, explainable AI, and RFM analysis in targeted marketing. *Mathematics*, 11(18), 3930.
15. Obiora, N. C., & Uchenna, N. D. (n.d.). Predicting customer churn in the telecommunication industry using machine learning algorithms: Performance comparison with logistic regression, random forest, and gradient boosting technique