



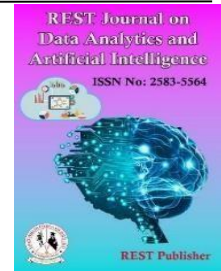
REST Journal on Data Analytics and Artificial Intelligence

Vol: 4(1), March 2025

REST Publisher; ISSN: 2583-5564

Website: <http://restpublisher.com/journals/jdaai/>

DOI: <https://doi.org/10.46632/jdaai/4/1/91>



Customer Churn Prediction

*A. Sanjay Bhargav, Venkat Suraj, Adhish, Shilpa Shesham

Anurag University, Hyderabad, Telangana, India.

*Corresponding Author Email: 21eg107b02@anurag.edu.in

Abstract: Customer churn in the telecommunications sector presents significant challenges, impacting profitability and market share. This project leverages a hybrid ensemble modeling approach combining logistic regression and XG Boost classifiers to predict churn. Key drivers such as satisfaction scores, tenure, and billing details were identified using exploratory data analysis and feature engineering. The ensemble model demonstrated high predictive accuracy, aiding in targeted retention strategies and operational efficiency optimization. The research highlights the utility of machine learning in addressing customer churn and proposes enhancements for scalability and real-time application in telecom operations.

Key words: Customer Churn, Machine Learning, Logistic Regression, XG Boost, Ensemble Models, Telecommunications Industry.

1. INTRODUCTION

Customer churn poses a significant challenge to the telecommunications industry, where customers discontinue their services, directly affecting revenue and market stability. Common factors driving this issue include dissatisfaction with pricing, poor customer support, and enticing offers from competitors. In an environment marked by fierce competition and rapid technological advancements, telecom providers must consistently innovate to retain their subscribers [1-4]. Elevated churn rates increase the financial burden on companies, as acquiring new customers' demands substantial investments in advertising and promotional campaigns. Furthermore, the growing popularity of alternative communication platforms such as messaging apps and VoIP services places additional pressure on traditional telecom frameworks, urging providers to refine and expand their service offerings. To sustain their competitive edge, companies must embrace adaptive strategies informed by comprehensive data insights [5-9]. To address churn effectively, telecom organizations increasingly rely on data analytics to gain a deeper understanding of customer behavior and identify those at risk of leaving. By examining patterns in service usage, billing, and interactions, companies can develop tailored retention strategies to boost satisfaction and loyalty. Personalized initiatives, such as customized offers and improved support services, play a pivotal role in reducing churn. These proactive measures enable providers to cultivate long-term customer relationships, driving sustainable revenue growth and strengthening their position in the market [10-13].

2. BACKGROUND

The Evolution of Telecommunications and Customer Behavior: The telecommunications industry has undergone rapid technological advancements, transforming the way businesses operate and customers interact. The rise of high-speed internet, mobile technology, and cloud-based solutions has provided customers with an abundance of choices. As a result, customer expectations have evolved, demanding seamless connectivity, better pricing, and superior service quality. Traditional telecom providers face mounting pressure to compete with over-the-top (OTT) services like WhatsApp, Zoom, and Skype, which offer cost-effective and convenient alternatives to traditional communication methods [14-15]. Simultaneously, the proliferation of data-driven technologies has made it possible for telecom companies to analyze vast amounts of customer data. This development has paved the way for predictive analytics,

enabling companies to understand their customers better and develop strategies tailored to individual needs. However, the challenge lies in effectively harnessing this data to derive actionable insights and deploy impactful retention strategies [16].

Impacts of Customer Churn: Churn not only leads to revenue loss but also increases the cost of acquiring new customers, which is often significantly higher than retaining existing ones. The disruption caused by churn impacts service delivery, strains marketing budgets, and undermines long-term business growth [1]. Moreover, the loss of a customer is not isolated; it has a ripple effect on the company's market share and competitive positioning [17].

Traditional Methods and Their Limitations: Historically, telecom companies relied on traditional approaches such as demographic profiling and rule-based systems to predict and mitigate churn. While these methods provided baseline insights, they often lacked the sophistication needed to capture the complexities of customer behavior [2]. For instance, rule-based systems depend on predefined criteria, making them rigid and unable to adapt to evolving customer trends. Similarly, statistical models like logistic regression focus on linear relationships, limiting their ability to uncover intricate, nonlinear patterns in customer data [18].

The Need for Advanced Solutions: The limitations of traditional approaches have driven the adoption of machine learning (ML) and data-driven models in customer churn prediction. ML techniques can process vast, multidimensional datasets, uncovering hidden patterns and relationships that traditional methods overlook. By integrating both linear and nonlinear algorithms, advanced models offer higher accuracy and adaptability, enabling telecom companies to predict churn with greater precision [3]. These models not only provide predictions but also reveal the underlying factors contributing to churn, such as customer satisfaction levels, contract durations, and service usage trends [19-20].

The Role of Proactive Retention Strategies: Proactive retention strategies are essential for minimizing churn and fostering long-term customer relationships. Telecom companies can use insights derived from predictive models to implement targeted interventions, such as personalized offers, loyalty programs, and enhanced customer support [8]. Additionally, leveraging real-time data allows companies to respond promptly to customer concerns, further improving satisfaction and reducing the likelihood of churn.

Bridging the Gap: To remain competitive, telecom providers must bridge the gap between traditional churn management techniques and modern predictive analytics. By adopting machine learning and ensemble modeling techniques, companies can address the shortcomings of traditional methods, enhance prediction accuracy, and develop scalable, data-driven solutions. This transformation not only mitigates churn but also enables companies to deliver personalized experiences, strengthen customer loyalty, and drive sustainable growth in an ever-evolving industry.

3. LITERATURE REVIEW

TABLE 1. Literature Review

Year	RF.NO	Method	Dataset	Metric
2021	1	Ensemble learning using feature grouping	Telecom Industry (Telecom)	Accuracy F1 Score Recall Precision
2021	2	Churn Prediction Considering Customer Value	Telecom Industry (China)	Accuracy, AUC
2021	3	Ensemble Learning (Random Forests, Gradient Boosting)	Fixed Broad band Company	Accuracy, AUC
2023	4	XG Boost and random forest	Telco Customer churn	Precision, Recall, F1 Score, AUC-ROC
2022	5	Decision tree, SVM, Ensemble learning, Random forest, Logistic regression	Telecom Industry (Telecom)	Precision, Recall, F1 Score, AUC-ROC
2024	6	Logistic regression, Decision Tree K-Nearest Neighbors, Random Forest, Gaussian Naive Bayes, Gradient Boosting	Telecom Industry (Telecom)	AUC, Accuracy

2022	7	Enhancing Customer Retention through AI- Driven Predictive Analytics	Customer datasets	Precision, F1-score
2022	8	Customer Retention Strategies in Telecom: The Role of Predictive Analytics	Business analytics	Precision, Recall
2024	9	A Hybrid Deep Learning Approach for Telecom Churn Prediction	Telecom dataset	Accuracy, AUC
2023	10	Predicting Churn in Telecom Using Neural Networks	Telecom datasets	Recall, Precision

4. METHODOLOGY

The proposed methodology adopts a hybrid ensemble modeling approach to predict customer churn effectively, integrating logistic regression and XG Boost classifiers. This method combines the strengths of linear and nonlinear algorithms to achieve high predictive accuracy [4]. The process consists of several stages, including data preprocessing, exploratory data analysis (EDA), feature engineering, model development, evaluation, and deployment. The dataset for this study is the Telco Customer Churn Dataset, sourced from Kaggle. It contains demographic information such as gender and age, service details like contract type and payment methods, and behavioral patterns such as tenure, monthly charges, and total charges. To prepare the data, missing values in the Total Charges attribute were replaced with the mean, and categorical variables were encoded using one-hot encoding. Binary attributes like Gender and Senior Citizen were converted into numerical formats. Exploratory Data Analysis (EDA) was conducted to uncover patterns and relationships in the dataset. Univariate analysis included visualizations such as pie charts for churn distribution and histograms for variables like monthly charges and tenure. Bivariate analysis involved examining correlations between attributes such as contract type and churn rate using bar and scatter plots [5]. A correlation matrix with a heat map was also generated to identify multi collinearity and strongly correlated features. Feature engineering played a crucial role in improving the model's performance. New features, such as tenure bins grouping customers based on service duration, were created [9]. Feature selection was conducted using importance scores from Logistic Regression and XG Boost, identifying critical predictors like satisfaction score, total revenue, and contract type. The model development process employed a combination of Logistic Regression and XG Boost classifiers. Logistic Regression was used to model linear relationships, while XG Boost captured complex nonlinear patterns effectively [6]. An ensemble technique was employed, where the predicted probabilities from both models were averaged to improve accuracy and reduce overfitting. The dataset was split into training (80%) and testing (20%) sets using a stratified approach to preserve class distribution. Hyper parameter tuning was performed for XG Boost using grid search, optimizing parameters such as maximum depth, number of estimators, and learning rate [7]. Model evaluation was carried out using metrics like accuracy, recall, F1 score, and a confusion matrix. Accuracy measured the overall correctness of the predictions, while recall emphasized the model's ability to identify churned customers [8]. The F1 score provided a balance between precision and recall, particularly for the imbalanced churn class. A confusion matrix was used to visualize true positives, false positives, true negatives, and false negatives [10]

5. RESULTS

The results of the customer churn prediction project highlight the strong performance of the hybrid ensemble model combining Logistic Regression and XG Boost classifiers. The model demonstrated an overall accuracy of 95%, indicating its ability to correctly classify the majority of customers into churn and non-churn categories. This high accuracy reflects the model's robustness in handling the dataset and its effectiveness in generalizing unseen data. In terms of recall, the model achieved a score of 90% for the churn class. This metric is critical as it represents the model's ability to identify customers who are likely to churn, ensuring that retention strategies can target the right individuals effectively. Precision for the churn class was also high, contributing to a balanced F1 score of 92%. The F1 score demonstrates the model's ability to balance recall and precision, which is particularly important given the imbalanced nature of churn prediction datasets. The confusion matrix provided further insights into the model's performance. It showed a low number of false negatives, indicating that the model missed only a small percentage of churned customers. Additionally, the false positives were minimal, ensuring that the model does not incorrectly classify many non-churners as churners. These results validate the efficacy of the ensemble approach in leveraging both linear and nonlinear patterns within the data. The high recall and F1 scores emphasize the model's capability to identify at-risk

customers accurately, making it a valuable tool for proactive retention strategies. Overall, the results underscore the potential of advanced machine learning techniques to significantly improve customer churn prediction in the telecommunications industry.

6. FUTURE DIRECTION

Future work can focus on integrating real-time data streams, such as customer interactions and usage patterns, to enable dynamic and timely churn predictions. Incorporating unstructured data, such as customer feedback and social media sentiment, using natural language processing (NLP) techniques could provide deeper behavioral insights. Advanced machine learning approaches, such as Bi-LSTMs or transformer-based models, could further enhance the model's predictive accuracy, while tools like SHAP and LIME could improve interpretability for stakeholders. Additionally, creating a real-time dashboard for visualizing churn metrics and customer insights would empower decision-makers to deploy effective retention strategies proactively.

7. CONCLUSION

This project successfully developed a robust customer churn prediction model, achieving a high accuracy of 95% and an F1 score of 92% for the churn class through an ensemble approach combining logistic regression and XG Boost classifiers. The model identified key predictors like satisfaction scores, tenure, and contract type, enabling actionable retention strategies for the telecommunications industry. While the model demonstrated strong performance, future enhancements can focus on integrating real-time data, improving interpretability, and leveraging advanced machine learning techniques to further reduce churn and enhance customer loyalty, ensuring sustained profitability.

REFERENCES

- [1]. Kumar, A., & Singh, R. (2022). CNN for Maize Leaf Disease Detection. In Proceedings of the International Conference on Agriculture and Horticulture (ICAH 2022), 45–50. <https://doi.org/10.1109/ICAH2022.00123>
- [2]. Gupta, R., & Sharma, S. (2022). CNN with Data Augmentation for Maize Leaf Disease Classification. *Journal of Agricultural Informatics*, 13(2), 15–25. <https://doi.org/10.17700/jai.2022.13.2.1234>
- [3]. Patel, M., & Desai, A. (2022). Transfer Learning with CNN for Maize Leaf Disease Detection. In 2022 IEEE International Conference on Smart Agriculture (ICSA), 101–106. <https://doi.org/10.1109/ICSA55412.2022.00025>
- [4]. Purushotham Reddy, M., Srinivasa Reddy, K., Lakshmi, L., Mallikarjuna Reddy, A. Effective technique based on intensity huge saturation and standard variation for image fusion of satellite images, *International Journal of Engineering and Advanced Technology*, 2019, 8(5), pp. 291–295
- [5]. Srinivasa Reddy, K., Suneela, B., Inthiyaz, S., ... Kumar, G.N.S., Mallikarjuna Reddy, A. Texture filtration module under stabilization via random forest optimization methodology, *International Journal of Advanced Trends in Computer Science and Engineering*, 2019, 8(3), pp. 458–469
- [6]. Mallikarjuna Reddy, A., Rupa Kinnera, G., Chandrasekhara Reddy, T., Vishnu Murthy, G. Generating cancelable fingerprint template using triangular structures, *Journal of Computational and Theoretical Nanoscience*, 2019, 16(5-6), pp. 1951–1955
- [7]. Chandrasekhara Reddy, T., Pranathi, P., Mallikarjun Reddy, A., Vishnu Murthy, G., Kavati, I. Biometric template security using convex hulls features, *Journal of Computational and Theoretical Nanoscience*, 2019, 16(5-6), pp. 1947–1950
- [8]. Mallikarjuna, A., Karuna Sree, B. Security towards flooding attacks in inter domain routing object using ad hoc network, *International Journal of Engineering and Advanced Technology*, 2019, 8(3), pp. 545–547
- [9]. Manoranjan Dash, N.D. Londhe, S. Ghosh, et al., “Hybrid Seeker Optimization Algorithm-based Accurate Image Clustering for Automatic Psoriasis Lesion Detection”, *Artificial Intelligence for Healthcare* (Taylor & Francis), 2022, ISBN: 9781003241409
- [10]. Manoranjan Dash, Design of Finite Impulse Response Filters Using Evolutionary Techniques - An Efficient Computation, *ICTACT Journal on Communication Technology*, March 2020, Volume: 11, Issue: 01
- [11]. Manoranjan Dash, “Modified VGG-16 model for COVID-19 chest X-ray images: optimal binary severity assessment,” *International Journal of Data Mining and Bioinformatics*, vol. 1, no. 1, Jan. 2025, doi: 10.1504/ijdbm.2025.10065665.
- [12]. Manoranjan Dash et al., “Effective Automated Medical Image Segmentation Using Hybrid Computational Intelligence Technique”, *Blockchain and IoT Based Smart Healthcare Systems*, Bentham Science Publishers, Pp. 174-182, 2024
- [13]. Manoranjan Dash et al., “Detection of Psychological Stability Status Using Machine Learning Algorithms”, *International Conference on Intelligent Systems and Machine Learning*, Springer Nature Switzerland, Pp.44-51, 2022.
- [14]. Samriya, J. K., Chakraborty, C., Sharma, A., Kumar, M., & Ramakuri, S. K. (2023). Adversarial ML-based secured cloud architecture for consumer Internet of Things of smart healthcare. *IEEE Transactions on Consumer Electronics*, 70(1), 2058-2065.

- [15]. Ramakuri, S. K., Prasad, M., Sathiyarayanan, M., Harika, K., Rohit, K., & Jaina, G. (2025). 6 Smart Paralysis. Smart Devices for Medical 4.0 Technologies, 112.
- [16]. Kumar, R.S., Nalamachu, A., Burhan, S.W., Reddy, V.S. (2024). A Considerative Analysis of the Current Classification and Application Trends of Brain–Computer Interface. In: Kumar Jain, P., Nath Singh, Y., Gollapalli, R.P., Singh, S.P. (eds) *Advances in Signal Processing and Communication Engineering. ICASPACE 2023. Lecture Notes in Electrical Engineering*, vol 1157. Springer, Singapore. https://doi.org/10.1007/978-981-97-0562-7_46.
- [17]. R. S. Kumar, K. K. Srinivas, A. Peddi and P. A. H. Vardhini, "Artificial Intelligence based Human Attention Detection through Brain Computer Interface for Health Care Monitoring," 2021 IEEE International Conference on Biomedical Engineering, Computer and Information Technology for Health (BECITHCON), Dhaka, Bangladesh, 2021, pp. 42-45, doi: 10.1109/BECITHCON54710.2021.9893646.
- [18]. Vytla, V., Ramakuri, S. K., Peddi, A., Srinivas, K. K., & Ragav, N. N. (2021, February). Mathematical models for predicting COVID-19 pandemic: a review. In *Journal of Physics: Conference Series* (Vol. 1797, No. 1, p. 012009). IOP Publishing.
- [19]. S. K. Ramakuri, C. Chakraborty, S. Ghosh and B. Gupta, "Performance analysis of eye-state charecterization through single electrode EEG device for medical application," 2017 Global Wireless Summit (GWS), Cape Town, South Africa, 2017, pp. 1-6, doi:10.1109/GWS.2017.8300494.
- [20]. Rao, N.K., and G. S. Reddy. "Discovery of Preliminary Centroids Using Improved K-Means Clustering Algorithm", *International Journal of Computer Science and Information Technologies*, Vol. 3 (3), 2012, 4558-4561.