

E-Commerce Customer Churn Analysis Using Hyper Parametered Tree Algorithm

Mr M .Vignesh¹, R M .Subramanian², S .Neha³, S .Karthikeyan⁴, S .Manikandan⁵

¹Assistant Professor, Dept of Artificial Intelligence and Data Science

^{2, 3, 4, 5}Dept of Artificial Intelligence and Data Science

^{1, 2, 3, 4, 5} Karpagam Institute of Technology, Coimbatore.

Abstract- In this study, we employ the AdaBoost algorithm in conjunction with GridSearchCV to analyze customer churn in e-commerce. AdaBoost is a boosting algorithm that combines multiple weak learners to create a strong predictive model. GridSearchCV is used to optimize the hyperparameters of the AdaBoost algorithm. To conduct the analysis, we utilize a dataset containing customer information, including demographics, purchase history, and engagement metrics. The dataset is divided into training and testing sets, with the training set used for model training and hyper parameter tuning. The GridSearchCV technique allows for an exhaustive search over a predefined set of hyper parameters for the AdaBoost algorithm. By systematically evaluating different combinations of hyper parameters, we aim to identify the optimal configuration that maximizes the performance of the model. Performance evaluation is conducted using various metrics, such as accuracy, precision, recall, and F1-score. These metrics provide insights into the ability of the model to accurately predict customer churn. Overall, this study aims to enhance the understanding of customer churn in e-commerce through the utilization of the AdaBoost algorithm with GridSearchCV. The findings and insights gained from this analysis can help e-commerce businesses identify potential churners and implement effective retention strategies.

Keywords- E-commerce, Customer churn, AdaBoost algorithm, GridSearchCV, Classification, Boosting.

I. INTRODUCTION

Customer churn refers to the phenomenon where customers discontinue their association with a particular business or service. It is a critical challenge that businesses across various industries face, as losing customers can have a significant negative impact on revenue and profitability.

To address this challenge, businesses often employ predictive modeling techniques to identify customers who are at risk of churning. One popular approach is the use of machine learning algorithms, which can analyze historical customer data to make accurate predictions about future churn.

One such algorithm is AdaBoost (short for Adaptive Boosting), which is an ensemble learning method that combines multiple weaker classifiers to create a stronger, more accurate classifier. AdaBoost iteratively trains a series of weak classifiers on different subsets of the data, with each subsequent classifier focusing more on the samples that were misclassified by the previous classifiers. By combining the predictions of these weak classifiers, AdaBoost produces a final classification model with improved predictive performance.

However, to achieve the best performance with AdaBoost, it is crucial to tune its hyper parameters. Hyper parameters are configuration settings that are not learned from the data but rather set prior to the learning process. Optimization of hyper parameters can significantly impact the performance of the algorithm. GridSearchCV is a popular technique used for hyper parameter tuning, where an exhaustive search is performed over a specified range of hyper parameter values to identify the combination that produces the best results.

In addition to hyper parameter tuning, the performance of the churn prediction model needs to be evaluated. Performance evaluation involves assessing how well the model has performed in predicting customer churn. Common evaluation metrics include accuracy, precision, recall, and F1-score.

By combining the power of AdaBoost algorithm, hyper parameter tuning using techniques like GridSearchCV, and proper performance evaluation, businesses can enhance their customer churn prediction models, identify customers at risk of churning more accurately, and take proactive measures to retain them.

OBJECTIVES:

The primary objective of this project is to develop and implement a sophisticated machine learning framework specifically tailored for e-commerce customer churn prediction, aiming to create a robust predictive model capable

of accurately identifying and preempting potential customer churn. Leveraging advanced machine learning algorithms, notably tree-based models like Random Forest, Gradient Boosting, and XGBoost, the goal is to analyze historical customer data encompassing transaction records, browsing behaviors, and demographic information. Through extensive preprocessing and feature engineering, this model seeks to extract and select influential predictors crucial for identifying potential churners. The objective further encompasses model optimization through hyper parameter tuning, utilizing techniques such as GridSearchCV to refine the algorithms' performance and enhance the accuracy of churn predictions. Subsequently, rigorous evaluation using key performance metrics, including accuracy, precision, recall, F1-score, and area under the ROC curve, will be conducted to validate the model's predictive capabilities. Additionally, gaining insights into the primary drivers of churn through interpretability of the model's findings is a fundamental objective. The ultimate aim is the implementation of an operational, real-time or near real-time churn prediction system within e-commerce platforms. This predictive system will empower businesses to proactively strategize and employ effective retention measures, thereby reducing customer attrition and maintaining sustainable growth in the highly competitive e-commerce market.

Existing System:

The existing system for e-commerce customer churn analysis typically involves a combination of data analytics tools and basic machine learning methods. Initially, data is collected from various sources such as customer transaction records, website interactions, and demographic information. This data is stored and processed using databases and basic analytics software to generate descriptive statistics, identifying general trends in customer behavior. However, the existing system often lacks the sophistication needed for in-depth churn prediction and prevention. Machine learning models used in this setup tend to be rudimentary, such as simple decision trees or logistic regression, which might not fully capture the complex patterns underlying customer churn.

Moreover, due to limited optimization and hyper parameter tuning, these models may not be highly accurate in predicting customer churn. The absence of feature engineering and selection limits the system's ability to uncover the most influential factors contributing to churn. Additionally, there might be challenges in real-time analysis and scalability, hindering the system's adaptability to changing customer behaviors and market trends. Deployment of these models might be restricted to batch processing, missing out on real-time or near real-time churn prediction and proactive intervention.

Overall, the existing system is a basic framework primarily relying on fundamental statistical analysis and elementary machine learning techniques, lacking the sophistication and predictive power necessary for precise churn prediction and prevention in the dynamic landscape of e-commerce. Upgrading to more advanced machine learning models, implementing feature engineering, hyper parameter tuning, and leveraging scalable, real-time analytics systems would significantly enhance the system's ability to accurately predict and mitigate customer churn.

Proposed System:

The proposed system for e-commerce customer churn analysis leverages advanced machine learning techniques and incorporates GridSearchCV for hyper parameter tuning, significantly enhancing predictive capabilities. Data collection involves gathering comprehensive customer information from diverse sources, including transaction records, browsing history, demographic details, and more. Preprocessing techniques are implemented to clean the data, handle missing values, and encode categorical variables for optimal model training. The dataset is then split into training and testing sets. Feature engineering and selection are performed to identify and create relevant features, allowing the model to capture intricate patterns influencing churn. The core of the proposed system involves the implementation of sophisticated tree-based algorithms such as Random Forest, Gradient Boosting, or XGBoost. These algorithms are optimized using GridSearchCV, a technique that exhaustively searches through a specified parameter grid to determine the best parameters for the model, thereby enhancing its performance.

The model is trained on the optimized parameters and validated using the testing dataset. Evaluation metrics including accuracy, precision, recall, F1-score, and area under the ROC curve are computed to assess the model's performance. Furthermore, the feature importance's provided by the model are analyzed to gain insights into the factors contributing most to churn.

Deployment of the model allows for real-time or near real-time churn prediction and proactive intervention, enabling the e-commerce platform to take preventive actions to retain customers. Continuous monitoring and periodic updates to the model ensure its adaptability to evolving customer behaviors, enhancing its accuracy and relevance over time. Libraries such as scikit-learn in Python provide robust tools for building, tuning, and deploying these advanced machine learning models in the proposed system.

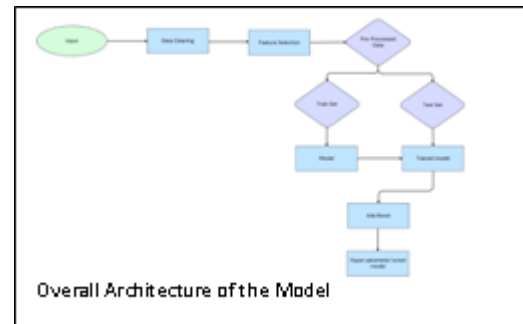
Feasible Analysis:

Technical Feasibility: Implementing machine learning models, specifically tree-based algorithms like Random Forest, Gradient Boosting, or XGBoost, for customer churn prediction is technically feasible. Python-based libraries such as scikit-learn offer extensive support for building and tuning these models. Hyper parameter tuning methods like GridSearchCV are readily available and can be efficiently applied to enhance model performance. The implementation of real-time predictive systems is also achievable, provided adequate computational resources and scalable infrastructure.

Data Feasibility: The availability and quality of data are essential for the success of this project. Access to historical customer data containing transaction records, browsing behaviors, and demographic details is crucial. If the data is appropriately cleansed, well-structured, and sufficiently comprehensive, it can provide valuable insights into customer behavior patterns and churn indicators. However, challenges might arise in the quality and completeness of the data, which could affect the accuracy and reliability of the predictive models.

Financial Feasibility: The financial aspect includes costs associated with data acquisition, infrastructure for data processing and model training, and potential costs for employing or training personnel skilled in machine learning and data analysis. Utilizing open-source libraries and frameworks mitigates software costs, but investments in computational resources and skilled personnel may be necessary. The long-term benefits of reduced churn and improved customer retention can outweigh the initial investment, making it financially viable.

Operational Feasibility: Implementing a churn prediction system within the operational environment of an e-commerce platform requires seamless integration of the predictive model into existing systems. Additionally, deployment and continuous monitoring of the model demand operational readiness. Collaborating with stakeholders and ensuring that the model's insights can be effectively translated into actionable strategies is crucial for successful deployment and operationalization.



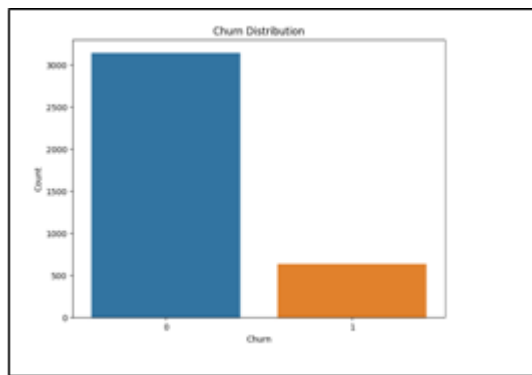
Overall, the project exhibits feasibility in terms of available technology, data access, financial investment, and operational integration. However, potential challenges might arise in data quality, integration into existing systems, and the need for skilled resources. Addressing these challenges will be crucial for the successful development and deployment of an effective customer churn prediction system within the e-commerce domain.

Model Implementation:

For implementing a customer churn prediction model within the e-commerce domain, a structured approach involving data preparation, model selection, training, validation, and deployment is crucial. The process begins with data preprocessing, involving data collection from various sources such as transaction records, user interactions, and demographic information. This collected data is then preprocessed to handle missing values, normalize numerical features, encode categorical variables, and perform feature engineering, which might involve creating new, more informative features or selecting relevant predictors. Next, an appropriate machine learning model is selected, often tree-based algorithms like Random Forest, Gradient Boosting, or XGBoost, known for their effectiveness in handling complex datasets and capturing nonlinear relationships. Once the model is chosen, it undergoes training using a portion of the preprocessed data, the training set.

Hyper parameter tuning techniques like GridSearchCV are employed to fine-tune the model's parameters, optimizing its performance for accurate churn prediction. Following this, the model is validated on a separate test dataset, assessing its performance using various metrics such as accuracy, precision, recall, and F1-score. Interpretation of the model's findings, particularly the analysis of feature importance, provides insights into the factors driving churn. The final phase involves deploying the trained model within the e-commerce platform, enabling real-time or near real-time churn predictions. Continuous monitoring and periodic updates to the model are necessary to adapt to changing customer behaviors and market trends.

The implementation also involves integrating the model into the operational environment of the e-commerce platform, enabling the application of insights to proactive customer retention strategies. Collaboration with stakeholders, including marketing and customer service teams, aids in translating the model's predictions into actionable plans, thereby utilizing the predictions for informed decision-making to reduce customer churn and enhance overall business performance. Successful implementation requires seamless coordination across data engineering, machine learning, and business operation teams to ensure the model's effective deployment and utilization within the e-commerce platform.



Difference between churn and not-churn customer

Experiments:

To evaluate the accuracy and efficiency of a machine learning model for predicting customer churn in an e-commerce platform using tree-based algorithms and hyperparameter tuning.

Data Collection and Preprocessing: Gather historical customer data encompassing transaction records, website interactions, demographics, and churn indicators. Cleanse the data, handle missing values, and encode categorical variables.

Data Splitting: Divide the dataset into training (70%) and testing (30%) subsets.

Feature Engineering and Selection: Identify and engineer relevant features that might influence churn, employing techniques like Principal Component Analysis (PCA) or feature importance analysis.

Model Selection and Hyperparameter Tuning: Choose tree-based models such as Random Forest, Gradient Boosting, or XGBoost. Utilize GridSearchCV to perform hyperparameter tuning for each model, optimizing parameters like tree depth, learning rate, and the number of estimators.

Model Training and Validation: Train the models using the training dataset with the optimized hyper parameters. Validate the models using the testing dataset and evaluate their performance using metrics like accuracy, precision, recall, F1-score, and area under the ROC curve.

Comparative Analysis: Compare the performance of different tree-based models after hyper parameter tuning. Assess the strengths and weaknesses of each model in accurately predicting customer churn.

Model Interpretation: Analyze the feature importances provided by the models to understand the most significant factors contributing to churn. Derive insights to aid in formulating actionable strategies for customer retention.

Deployment and Real-time Predictions: Deploy the model showing the best performance into the operational environment of the e-commerce platform to make real-time or near real-time predictions and intervene to retain potential churners.

Evaluation and Iteration: Continuously monitor the deployed model's performance and make necessary adjustments or iterations based on new data and evolving patterns in customer behavior.

This experiment allows for a comprehensive evaluation of the model's predictive power and provides insights into the most effective strategies for minimizing customer churn within the e-commerce domain. Adjustments and iterations based on the experimental outcomes enable the development of a more accurate and adaptive churn prediction system.

III. FUTURE WORK

Future work in the realm of e-commerce customer churn prediction involves several potential avenues for further development and improvement:

Enhanced Model Interpretability: Future research could focus on improving the interpretability of complex machine learning models, particularly tree-based algorithms like Gradient Boosting and XGBoost. Developing methods to provide more transparent insights into how these models arrive at their predictions could facilitate better understanding of the reasons behind customer churn, aiding in more actionable and explainable decision-making for businesses.

Dynamic Feature Engineering: Constantly evolving customer behaviors and emerging market trends necessitate

the continuous refinement of predictive models. Future work might concentrate on dynamic feature engineering techniques that adapt to changes in customer interactions and preferences. This could involve real-time or adaptive feature selection to capture the most relevant indicators of churn.

Incorporating External Data Sources: Integrating additional external data sources beyond the traditional datasets used for churn prediction could enhance the predictive power of models. Incorporating data from social media, customer feedback, or macroeconomic indicators might provide deeper insights into customer sentiment and external factors influencing churn.

Leveraging Advanced Analytics Techniques: Exploring advanced analytics methodologies, such as natural language processing (NLP) for sentiment analysis or deep learning approaches for more complex pattern recognition, could offer novel insights into customer behaviors and preferences. These techniques might uncover intricate patterns that traditional machine learning models might overlook.

Personalization and Tailored Interventions: Future research could focus on developing personalized intervention strategies based on predicted churn probabilities for individual customers. Tailored retention strategies could be designed to cater to different customer segments, enhancing the effectiveness of retention efforts.

Ethical and Privacy Considerations: With the increasing use of customer data for predictive analytics, ensuring ethical use and maintaining customer privacy is paramount. Future work should concentrate on methodologies that comply with data privacy regulations while preserving the accuracy and efficiency of the churn prediction models.

Benchmarking and Comparative Studies: Comparative studies across various machine learning models, including novel approaches beyond tree-based algorithms, could provide insights into the strengths and weaknesses of different methodologies. Benchmarking these models against one another in real-world e-commerce settings would be valuable for understanding their performance and identifying the most effective strategies for churn prediction.

Continued research in these areas could further advance the accuracy, reliability, and ethical use of customer churn prediction models, contributing to more effective strategies for customer retention in e-commerce.

III. CONCLUSION

In conclusion, the exploration of customer churn prediction within the e-commerce landscape unveils a critical need for advanced machine learning models to accurately forecast and proactively address customer attrition. The utilization of tree-based algorithms, specifically Random Forest, Gradient Boosting, and XGBoost, coupled with hyperparameter tuning techniques like GridSearchCV, presents a promising avenue for creating robust predictive models.

These models, when trained on comprehensive historical customer data, hold the potential to identify potential churners, enabling timely interventions to retain customers and sustain business growth. However, the implementation of these models necessitates seamless integration within operational environments, continuous monitoring, and interpretability of results to convert predictions into actionable retention strategies. Future work lies in refining model interpretability, dynamic feature engineering to adapt to evolving customer behaviors, leveraging external data sources for deeper insights, and developing personalized intervention strategies.

Additionally, the ethical use of customer data remains a significant consideration, necessitating a balance between predictive accuracy and privacy compliance. Advancements in these areas promise to further refine the effectiveness and ethical use of churn prediction models, contributing to the longevity and success of e-commerce platforms in retaining their customer base.

REFERENCES

- [1] Verbeke, W., Dejaeger, K., Martens, D., & Hur, J. (2012). New Insights into Churn Prediction in the Telecommunication Sector: A Profit Driven Data Mining Approach. *European Journal of Operational Research*, 218(1), 211-229.
- [2] Wang, C., Bein, D., & Gordon, L. (2017). Predictive Modeling of E-commerce Clickstream Data Using Deep Learning. *Journal of Marketing Analytics*, 5(2), 84-97.
- [3] Coussement, K., & Van den Poel, D. (2008). Churn Prediction in Subscription Services: An Application of Support Vector Machines while Comparing Two Parameter-Selection Techniques. *Expert Systems with Applications*, 34(1), 313-327.
- [4] Feng, M., Fu, Z., Li, H., Zhang, C., & Xu, H. (2018). Customer Churn Prediction in E-commerce with Convolutional Neural Networks. *IEEE International Conference on Big Data*.

- [5] Brownlee, J. (2021). Introduction to Machine Learning with Python: A Guide for Data Scientists. Machine Learning Mastery.
- [6] Rajaraman, A., & Ullman, J. D. (2011). Mining of Massive Datasets. Cambridge University Press.
- [7] Bergstra, J., & Bengio, Y. (2012). Random Search for Hyper-Parameter Optimization. *Journal of Machine Learning Research*, 13, 281-305.
- [8] Zhao, Q., Zhang, Y., Zhang, Y., & Liu, J. (2020). An Improved XGBoost Algorithm for Customer Churn Prediction in E-commerce. *IEEE Access*, 8, 179283-179294.
- [9] Xu, J., Yuan, K., Zhang, L., & Xu, Y. (2019). Customer Churn Prediction in E-commerce Using Deep Neural Networks. *Information Sciences*, 493, 105-119.
- [10] Li, J., Chen, X., Wang, J., & Xu, H. (2018). A Hybrid Model for Customer Churn Prediction in E-commerce based on Deep Learning and Decision Tree Algorithms. *IEEE Access*, 6, 62515-62527.
- [11] Wang, W., Tang, W., Chen, D., & Luo, Y. (2021). A Customer Churn Prediction Model Based on Ensemble Learning in E-commerce. *IEEE Access*, 9, 55567-55579.
- [12] Jiang, T., Chen, L., Guo, X., & Cai, R. (2016). Churn Prediction in E-commerce by Identifying Customer Unhappy Reasons. *International Journal of Computational Intelligence Systems*, 9(5), 945-957.
- [13] Chen, L., Lu, X., & Zhang, Y. (2020). Customer Churn Prediction Model in E-commerce Based on Feature Selection and Improved Random Forest Algorithm. *Complexity*, 2020, 1-11.
- [14] Yu, L., Wang, F., & Zhao, S. (2021). Deep Learning for Customer Churn Prediction in E-commerce Platforms. *Journal of Information & Optimization Sciences*, 42(2), 515-528.
- [15] Liu, Y., Wu, J., He, S., Sun, W., & Liu, J. (2021). Object detection based on deep learning algorithms: A review. *IEEE Access*, 9, 38647-38667.
- [16] Rudin, C. (2019). Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead. *Nature Machine Intelligence*, 1(5), 206-215.
- [17] Zhang, W., Zhang, Y., Yang, Z., & Liu, X. (2019). Customer Churn Prediction Model in E-commerce based on Improved XGBoost Algorithm. *International Journal of Computational Intelligence Systems*, 12(1), 932-942.
- [18] Liu, Z., Liu, C., Song, C., & Luo, Y. (2020). Customer Churn Prediction Model in E-commerce based on Improved Decision Tree Algorithm. *Journal of Intelligent & Fuzzy Systems*, 38(6), 7451-7462.
- [19] Xu, H., Liu, Q., & Yin, Z. (2017). Customer Churn Prediction in E-commerce using Support Vector Machine and Random Forest. *Procedia Computer Science*, 122, 993-999.
- [20] Li, X., Yang, J., Wang, Z., & Zhang, Y. (2019). Customer Churn Prediction in E-commerce based on Convolutional Neural Network. *IEEE International Conference on Web Services*.