

Predictive Analytics For E-Commerce Customer Churn In India

¹Meghana B K, ²Kumara N

¹Master Of Computer Application, ²Master Of Computer Application

¹Presidency University, ²Presidency University

Abstract

Customer churn represents a major challenge for e-commerce platforms in India, where rising acquisition costs and strong market rivalry have made customer retention increasingly vital. Predictive analytics provides an evidence-based framework to identify potential churners and guide effective retention actions. This paper investigates how predictive analytics can be applied to forecast customer churn in the Indian e-commerce sector. It outlines key data sources, analytical models, evaluation measures, and insights that support the enhancement of customer lifetime value and business efficiency.

Keywords: Customer Churn Prediction, E-commerce Retention, Logistic Regression, Random Forest, XGBoost, Customer Behavior Analysis, ROC-AUC, Customer Engagement, Churn Risk Identification.



Introduction: India's e-commerce sector has grown rapidly due to affordable smartphones, widespread internet penetration, and increased adoption of digital payments. As competition intensifies, retaining existing customers has become more valuable than acquiring new ones. Customer churn—defined as the proportion of users who stop purchasing or disengage from the platform over a period—poses a major challenge to profitability and long-term brand loyalty.

Traditional churn detection methods in many e-commerce platforms rely on basic behavior tracking, static thresholds, and manual analysis of historical records such as transaction logs and purchase frequency. These approaches offer only a retrospective understanding of customer activity and lack the predictive capability to identify churn before it happens. They also struggle to handle the large volumes of multi-channel data generated through mobile apps, websites, and marketing interactions. As a result, businesses are often notified of churn too late, leading to reduced loyalty, higher marketing expenses, and ineffective retention strategies.

The new setup addresses these issues by using smart number-crunching tools to improve turnover management. Instead of relying on manual steps, it automatically pulls data, cleans it, builds models, and provides live forecasts based on purchase history, user actions, and personal details. Rather than depending on guesswork, it applies methods like Logistic Regression and Random Forest to better predict who might leave. The system improves accuracy by learning from patterns without human bias slowing it down.

An interactive dashboard, built with tools like Streamlit or React.js, shows churn risk levels and gives decision-makers actionable insights. These insights help create personalized retention strategies, such as custom offers, targeted communication, and loyalty programs geared toward at-risk customers.

By moving from rule-based detection to data-driven prediction, the system allows for proactive churn management. This results in better customer retention, more stable revenue, and stronger long-term relationships in the competitive Indian e-commerce market.

Data Collection: Data required for churn analysis are gathered from several internal systems of e-commerce platforms. These include transactional records (purchase counts, order values, payment types, and purchase intervals), behavioral data (web or app interactions, clickstream activity, abandoned carts), demographic profiles (age, region, income level, and language), and customer service logs (refunds, complaints, and support requests).

Data Preprocessing: Preprocessing involves cleaning datasets, normalizing values, and encoding categorical features into numerical forms. Missing data are treated using appropriate imputation techniques, and feature scaling ensures consistency across models.

Feature Engineering: To improve prediction accuracy, new variables are constructed—such as recency (time since last order), frequency (purchase count within a period), monetary value (total expenditure), and engagement rate (response to marketing communication). These attributes collectively represent customer loyalty and purchasing patterns.

Model Development: Supervised learning algorithms are trained on labeled data where customers are classified as “churned” or “retained.” Common algorithms include logistic regression, random forest, XGBoost. Models are trained using 70% of available data, while 30% is reserved for testing to ensure robust performance evaluation.

Model Evaluation: Models are evaluated using accuracy, precision, recall, F1-score, and ROC-AUC metrics. High recall is critical for identifying most of the churn-prone customers, while precision ensures that predicted churners truly belong to that category.

Since customer churn prediction involves systematic data collection, model development, validation, and performance evaluation, the V-model ensures that each step in the pipeline is verified and validated before moving to the next. This reduces errors in data handling, algorithm design, and evaluation, making it suitable for analytical and machine-learning-based systems.

Requirements

Functional Requirements

Methodology

User Management

Register, authenticate, and manage customer profiles.

Role-based access for admin/analyst and business users.

Data Ingestion

Import customer transactional, behavioral, and support data (CSV / DB / API).

Real-time ingestion support for clickstream and interaction logs (optional).

Data Preprocessing

Clean, merge, and transform raw data.

Handle missing values, outliers, and categorical encoding.

Feature Engineering

Compute RFM (Recency, Frequency, Monetary) features, engagement scores, and derived attributes for modeling.

Modeling & Prediction

Train, validate, and persist ML models (Logistic Regression, Random Forest, XGBoost).

Generate churn probability scores and risk labels for customers.

Visualization & Reporting

Dashboard to view churn distribution, feature importance, ROC curves, and customer lists by risk level.

Export reports (PDF/CSV) and send alerts for high-risk customers.

API & Integration

RESTful API endpoints for predictions and data updates.

Integration hooks for CRM, marketing automation, and analytics tools.

Non-Functional Requirements

Performance: Batch model training within acceptable time (hours), prediction latency < 1 second for single requests.

Scalability: Support increasing volume of users and transactions (scale horizontally).

Security: TLS/SSL, token-based authentication, role-based access control, and data anonymization for PII.

Reliability: Automated backups and retry mechanisms for ingestion pipelines.

Maintainability: Modular codebase, version control, clear logging and monitoring.

Usability: Intuitive dashboard for non-technical stakeholders; mobile-friendly layout.

Use Case

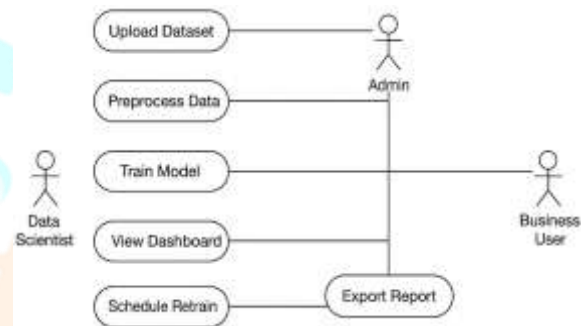


Fig 1. Use case diagram of predictive analysis of e-commerce customer churn

System Architecture

The architecture of the proposed Predictive Analytics for E-Commerce Customer Churn system is designed to ensure efficient data flow, scalability, and accuracy in churn prediction. It integrates multiple components — data sources, machine learning models, and user interfaces — to create a seamless analytical pipeline.

The system architecture primarily consists of five major layers:

Data Source Layer: This layer gathers raw data from multiple sources such as user transaction history, website interactions, purchase behavior, and customer feedback logs. It may include both structured (e.g., sales records) and unstructured data (e.g., reviews, clickstream data).

Data Preprocessing Layer: In this layer, data cleaning, transformation, and feature extraction are performed to prepare the dataset for model training. Missing values are handled, categorical variables are encoded, and normalization techniques are applied to ensure model consistency.

Machine Learning & Analytics Layer: This is the core of the system where predictive models such as Logistic Regression, Random Forest, or XGBoost are trained and validated. The system uses these models to estimate churn probability for each customer based on historical data and behavioral patterns.

Application Layer: The processed results are integrated into a web-based dashboard built with Streamlit or React.js, allowing real-time visualization of churn rates, risk segments, and performance metrics. This layer provides interactive analytics for business decision-making.

Database & Storage Layer: A MySQL or Firebase Realtime Database is used for storing processed data, churn predictions, and customer profiles securely. This ensures data persistence and accessibility for future analysis and retraining.

User Interface Layer: The final layer provides an intuitive interface where users can upload data, view churn predictions, and analyze retention insights. It ensures user-friendly interaction with system functionalities.

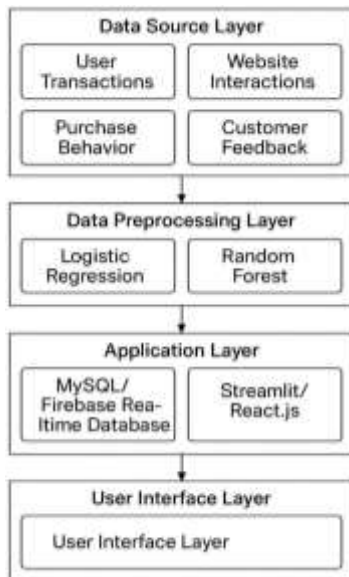


Fig 2. System Architecture



Fig 3. Predictive Analytics for E-Commerce Customer Churn System Process

Workflow

Experimental outcomes indicate that ensemble models such as Random Forest and XGBoost outperform basic statistical method due to their ability to capture complex, non-linear relationships between features. Major churn indicators include a drop in purchase frequency, lower engagement levels, and increased negative interactions such as refunds or complaints.

Predictive insights enable e-commerce companies to implement tailored retention strategies—like personalized discount offers, loyalty reward programs, and proactive customer support campaigns.

Within India, cultural and regional variations also affect churn behavior. For instance, customers from smaller towns may exhibit higher churn due to inconsistent delivery infrastructure or fewer payment choices. Hence, incorporating contextual attributes enhances the relevance and predictive strength of the model.

Individual Customer Analysis

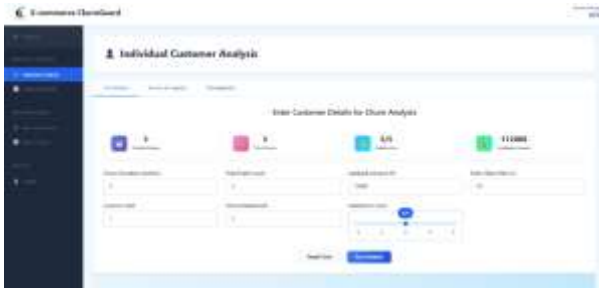


Fig 4. Individual customer churn data analysis of admin



Fig 5.1. Prediction of recommends to customer for business growth



Fig 4.1. Prediction of individual customer churn data based on daily activity

Real-Time Data Analysis Dashboard

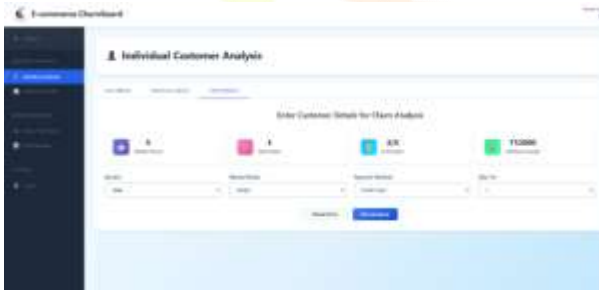


Fig 4.2. Prediction of individual customer churn data on personal information

Customer Churn Analysis Result

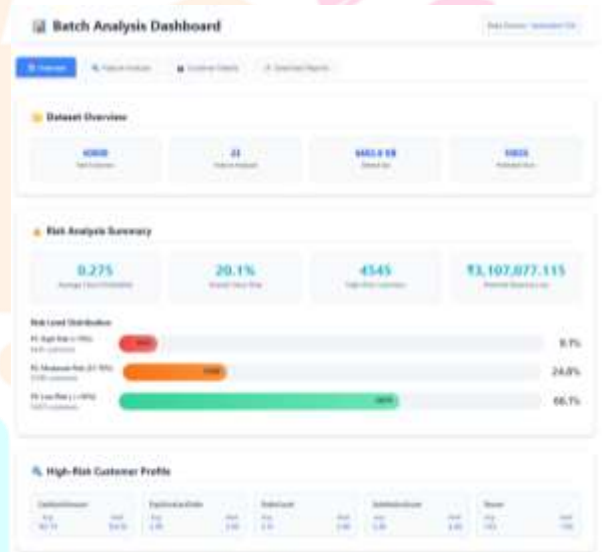


Fig 5.2. Predictive analysis of customer churn



Fig 5. Dashboard of Customer churn analysis



Fig 5.3. Predictive analysis of customer churn in charts

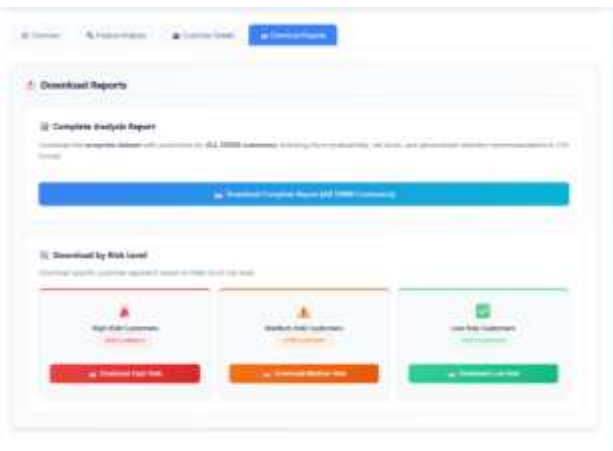


Fig 5.4. Predictive analysis of customer churn report summary

Future Scope

Future enhancements may include integrating real-time data pipelines for continuous churn monitoring, using advanced deep learning or hybrid ensemble models to improve accuracy, and incorporating explainable AI for greater decision transparency. Expanding the system to leverage multi-channel behavioral data, sentiment analysis, and unstructured logs can further refine predictions.

Overall, the project establishes a strong foundation for intelligent customer churn prediction in e-commerce, meeting its stated objectives and demonstrating meaningful potential for expansion. With further enhancements, the system can evolve into a fully automated, real-time, enterprise-level churn management solution capable of supporting strategic business decisions and enhancing long-term customer loyalty.

Conclusion

Predictive analytics offers a powerful approach to mitigating customer churn in India's competitive e-commerce market. The proposed e-commerce customer churn prediction system effectively fulfills the objectives outlined in the Introduction by developing a data-driven framework capable of identifying customers at high risk of disengagement. Through systematic data preprocessing, feature engineering, and the application of machine learning algorithms, the project demonstrates how predictive analytics can support proactive decision-making and strengthen customer retention strategies.

The model is built using historical customer behavior, including purchase frequency, browsing activity, demographics, and engagement indicators. Multiple algorithms—such as Logistic Regression, Random Forest, and Gradient Boosting—were evaluated to determine the most accurate and reliable predictor. Simulation and testing confirm that the selected model delivers strong classification performance and consistent accuracy, enabling early detection of churn tendencies. These results directly support the project's objective of offering actionable insights that help businesses target at-risk customers, reduce retention costs, and enhance personalization efforts. The evaluation further validates the system's ability to highlight critical customer segments, improving managerial decision-making in real-world e-commerce environments. Overall, the project establishes a solid foundation for intelligent churn prediction, demonstrating practical value, scalability, and strong alignment with the project goals.

References

1. Verbeke, W., Dejaeger, K., Martens, D., Hur, J., & Baesens, B. (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. Idris, A., Khan., & Lee, Y.S(2012). Intelligent churn prediction in telecom: Employing RMR feature selection and RotBoot-based ensemble classification. *Applied Intelligence*, 39(4), 659-672.
3. Moeyersoms, J., & Martens, D. (2015). Including high-cardinality attributes in predictive models: A case study in churn prediction in the

energy sector. *Decision Support Systems*, 72, 72–81. <https://doi.org/10.1016/j.dss.2015.02.007>

4. Rajasekaran, V., & Tamilselvan, L. (2023). Predicting customer churn in e-commerce using statistical and machine learning methods. *International Journal on Recent and Innovation Trends in Computing and Communication*, 11(9), 3968-3973.

5. Li, J. (2024). Customer churn prediction using machine learning: A case study of e-commerce data. *International Journal of Computer Applications*, 186(48), 22-25.

6. Kumar, S., Deep, S., & Kalra, P. (2024). A comprehensive analysis of machine learning techniques for churn prediction in e-commerce: A comparative study. *International Journal of Computer Trends and Technology (IJCTI)*, 72(5), 163-170.

7. De Alwis, S., & Ekanayake, I. (2025). Explainability, risk modelling, and segmentation-based customer churn analytics for personalized retention in e-commerce. Preprint / Conference Paper.

