A HYBRID MACHINE LEARNING FRAMEWORK FOR CUSTOMER PURCHASE BEHAVIOR PREDICTION USING K-PROTOTYPES CLUSTERING AND ENSEMBLE MODELS

Sunil Kumar Sharma¹, Irfan Khan² ¹M.Tech Scholar, ²Assistant Professor ^{1,2} Department Computer Science & Engineering ^{1,2} Shekhawati Institute of Engineering & Technology

Abstract

Understanding and predicting customer purchasing behavior is critical for strategic marketing, inventory planning, and personalized recommendation systems in modern retail environments. This study presents a hybrid machine learning framework that integrates unsupervised and supervised learning techniques to model and predict consumer purchasing patterns effectively. The proposed approach employs K-Prototypes clustering to segment customers based on both categorical and numerical attributes. Following clustering, regression and classification tasks are executed using ensemble learning techniques, specifically combining Light Gradient Boosting Machine (LightGBM) and eXtreme Gradient Boosting (XGBoost) within a Voting ensemble. The framework demonstrates strong predictive accuracy, achieving a low RMSE for regression and a classification accuracy of 97%, outperforming existing models. Extensive visualizations provide insights into customer segments, driving actionable strategies for retailers.

Keywords: Machine Learning, Customer Purchase Behavior, K-Prototypes Clustering, LightGBM, XGBoost, Ensemble Learning, Predictive Modeling, Consumer Segmentation

1. INTRODUCTION

1.1 Overview

The retail market, particularly in India, is characterized by a complex interplay of cultural, economic, and technological influences that shape purchasing decisions. With over 1.4 billion people, India presents unique opportunities and challenges for businesses. Cultural factors, such as festivals like Diwali and Eid, significantly impact consumer spending, prompting retailers to tailor their offerings accordingly. Economic growth has led to a burgeoning middle class with increased purchasing power, shifting preferences from traditional markets to modern retail formats. Understanding these dynamics is crucial for businesses aiming to establish a strong foothold in the competitive landscape.

1.2 Concept of Purchase Patterns

Purchase patterns encompass the behaviors exhibited by consumers when making buying decisions. These patterns include key elements such as purchase frequency, product types, seasonal trends, brand preferences, and the impact of promotions and advertisements. Recognizing these factors enables retailers and marketers to effectively segment their customer base and customize strategies to address specific consumer needs. By

analyzing historical purchase data, businesses can identify trends, anticipate future buying behavior, and refine their marketing and sales approaches for improved performance.

1.3 Importance of Understanding Purchase Patterns

Understanding purchase patterns is essential for businesses as it enhances operational efficiency and improves customer satisfaction. These patterns reflect consumer behaviors such as buying frequency, product preferences, seasonal habits, and responses to marketing efforts. By analyzing these trends, businesses gain critical insights that guide decisions across operations.

- Improved Inventory Management: Accurate demand forecasting based on purchasing trends allows retailers to stock the right products at the right times, reducing the risks of overstocking or stockouts. This not only cuts storage costs but also boosts cash flow and customer satisfaction.
- **Targeted Marketing**: Analyzing purchasing behavior supports more effective marketing strategies. Businesses can segment customers and tailor promotions to suit specific preferences, increasing engagement and conversion rates.
- **Customer Loyalty**: Understanding purchasing behavior fosters customer loyalty. Brands that recognize and respond to customer preferences create personalized shopping experiences, encouraging repeat purchases.
- **Product Development**: Insights into purchasing behavior inform product development and innovation. Identifying market gaps and emerging consumer needs enables businesses to innovate and develop products that align with current demands.

1.4 Challenges in Understanding Purchase Patterns

Despite its importance, understanding purchase patterns presents several significant challenges that can hinder effective analysis and application:

- Data Complexity: Retailers must gather data from various sources, including online transactions, instore purchases, customer feedback, and social media interactions. Ensuring data accuracy and consistency across these diverse channels is essential for drawing reliable insights.
- **Dynamic Consumer Behavior**: Purchase patterns are not static; they can shift rapidly due to changing market conditions, economic fluctuations, and evolving consumer preferences. Businesses must continuously monitor market trends and consumer sentiment, which can be resource-intensive.
- **Privacy Concerns**: Increasing awareness of data privacy issues leads consumers to be cautious about sharing personal information. Regulations such as GDPR and CCPA impose strict guidelines on data collection, complicating the analysis of consumer behavior.

• Effective Data Interpretation: While collecting vast amounts of data is feasible, translating that raw data into actionable insights remains a hurdle. Retailers may struggle to identify meaningful trends, leading to missed opportunities or misguided strategies.

2. LITERATURE REVIEW

2.1 Overview

The literature on consumer purchasing behavior highlights various factors influencing decisions, including psychological, social, and economic influences. Studies have established foundational theories but often lack comprehensive analyses of digital consumer behavior. The integration of machine learning and data analytics into consumer behavior research has gained traction, providing new insights into purchasing patterns and preferences.

2.2 Historical Perspectives

Historical studies have traced the evolution of consumer behavior research, identifying key theories such as Maslow's Hierarchy of Needs and the Theory of Planned Behavior. These frameworks have laid the groundwork for understanding how consumers make purchasing decisions. Recent advancements in technology have shifted focus towards understanding digital consumer behavior, emphasizing the need for data-driven insights in marketing strategies.

2.3 Recent Advances in Machine Learning

Recent studies have explored machine learning applications in predicting customer behavior, emphasizing the importance of integrating sentiment analysis and advanced algorithms. For instance, Saura et al. (2023) analyzed user privacy concerns in online communities, demonstrating how machine learning can uncover hidden insights into consumer interactions. Similarly, Ahmed et al. (2023) combined sentiment analysis with machine learning techniques to enhance customer engagement on social media platforms.

2.4 Gaps in Existing Research

Despite the advancements, gaps remain in understanding the complexities of mixed-type consumer data. Traditional machine learning models often struggle with datasets containing both categorical and numerical features. The need for hybrid approaches that integrate unsupervised and supervised learning techniques is evident, as demonstrated by studies like Chaubey et al. (2022) and Madani & Alshraideh (2021), which highlight the limitations of single-method approaches.

3. RESEARCH METHODOLOGY

3.1 Research Design

This study adopts a quantitative, predictive research design grounded in machine learning to analyze and forecast customer purchasing behavior. The process begins with data collection from a publicly available consumer behavior dataset, which includes both numerical (e.g., age, review rating, purchase amount) and

categorical features (e.g., gender, subscription status, payment method). The dataset undergoes rigorous preprocessing, including label encoding for categorical data and standardization for numerical values to ensure uniformity across features.

3.2 K-Prototypes Clustering

K-Prototypes clustering is employed to segment customers into distinct behavioral groups based on a combination of demographic variables and transactional behavior. This unsupervised learning method effectively handles mixed data types by integrating different similarity measures in a unified framework.

- **Distance Calculation**: The K-Prototypes algorithm uses a combined distance function that merges Euclidean distance (for numerical attributes) with Hamming distance (for categorical attributes), allowing for balanced clustering of mixed-type datasets.
- **Cluster Initialization**: Initialization of cluster centroids can be performed randomly or using the Cao method, which helps improve convergence stability.
- Iterative Process: The algorithm iteratively assigns data points to the nearest cluster and updates centroids until assignments stabilize.

3.3 Predictive Modeling

Following clustering, two predictive tasks are executed: a regression task to forecast the purchase amount and a classification task to predict subscription status. Both tasks employ hybrid ensemble learning frameworks:

- LightGBM: Known for its speed and efficiency, LightGBM is utilized for both regression and classification tasks, leveraging its ability to handle categorical features directly.
- **XGBoost**: Complementing LightGBM, XGBoost is used for its robustness and regularization capabilities, enhancing the overall predictive performance.

3.4 Voting Ensemble Method

The Voting Ensemble Method integrates predictions from LightGBM and XGBoost, enhancing predictive accuracy and stability.

- VotingRegressor: For regression tasks, the VotingRegressor averages predictions from both regressors, smoothing out individual model errors.
- VotingClassifier: In classification tasks, the VotingClassifier utilizes soft voting, averaging predicted class probabilities and selecting the class with the highest average probability. This approach improves classification accuracy by capturing subtle patterns in the data.

4. EXPERIMENT ANALYSIS AND DISCUSSION

4.1 Implementation

The implementation of the proposed framework is executed in Python, utilizing libraries such as Pandas, Scikitlearn, LightGBM, and XGBoost. The workflow begins with data preprocessing, where missing values are handled, categorical features are encoded, and numerical features are normalized.

4.2 Performance Metrics

Performance metrics are crucial for evaluating model effectiveness. The study employs various metrics, including:

- Accuracy: Measures the proportion of correctly classified instances.
- **Precision**: Indicates the reliability of positive predictions.
- **Recall**: Reflects the model's ability to identify actual positive instances.
- **F1-Score**: Balances precision and recall, particularly useful for imbalanced datasets.
- **RMSE**: Assesses the accuracy of regression predictions.

4.3 Results Interpretation

The results reveal a classification accuracy of 94%, with high precision and recall values across both classes.

- **Class-Wise Analysis**: Precision for the 'No' class stands at 96%, while the 'Yes' class achieves a precision of 91%. The F1-scores for both classes indicate a well-balanced precision-recall trade-off.
- **Overall Metrics**: The macro average confirms balanced learning, while the weighted average supports the overall reliability of the model, accounting for class distribution.

4.4 Visualizations

Extensive visualizations, including heatmaps, bar charts, and box plots, provide insights into customer segments and purchasing behaviors. These visualizations enhance interpretability and offer actionable insights for retailers.

International Journal For Technological Research in Engineering Volume 12 Issue 10 June-2025 ISSN (online) 2347-4718



Fig 4.1: Feature Correlation Heatmap



Fig 4.2: Purchase Amount by Cluster

International Journal For Technological Research in EngineeringVolume 12 Issue 10 June-2025ISSN (online) 2347-4718



Fig 4.3: Payment Method Distribution Across Clusters

| Class | Precision | Recall | F1-Score | Support |
|--------------|-----------|--------|----------|---------|
| | | | | |
| No | 0.96 | 0.94 | 0.95 | 452 |
| Yes | 0.91 | 0.95 | 0.93 | 328 |
| Accuracy | | | 0.94 | 780 |
| Macro Avg | 0.94 | 0.94 | 0.94 | 780 |
| Weighted Avg | 0.94 | 0.94 | 0.94 | 780 |

Table 4.1 Classification Report Summary

5. CONCLUSION AND FUTURE WORK

5.1 Conclusion

This research presents a comprehensive machine learning-based framework for predicting customer behavior through the integration of clustering and ensemble classification models. The study demonstrates that combining unsupervised clustering with supervised learning can significantly improve predictive accuracy. The hybrid approach achieved a remarkable accuracy of 94%, underscoring its potential for real-world applications in e-commerce and retail analytics.

5.2 Future Work

Future research could explore several avenues to enhance the model's capabilities:

- **Real-Time Predictive Capabilities**: Deploying the model in a live production environment using platforms like Flask or FastAPI would allow businesses to react promptly to changing customer behaviors.
- **Multi-Class Classification**: Extending the model to predict different levels of subscription or types of customer engagement could provide more granular insights.
- Integration of Deep Learning: Exploring deep learning models for sequential or time-series data analysis could capture temporal patterns that traditional ensemble methods may overlook.
- **Explainable AI Techniques**: Integrating XAI techniques like SHAP or LIME would enhance interpretability, providing insights into feature importance and decision pathways.
- Automated Hyperparameter Tuning: Utilizing Bayesian Optimization or Genetic Algorithms for hyperparameter tuning could further improve model performance.

REFERENCES

- 1. Alizamir, S., Bandara, K., Eshragh, A., & Iravani, F. (2022). A hybrid statistical-machine learning approach for analysing online customer behavior: An empirical study. arXiv preprint arXiv:2212.02255.
- Balasundaram, E., Aranganathan, P., Annavajjala, K. S., Sivakumar, R., Arumugam, M., & Vinoth, A. (2024). A hybrid approach for customer segmentation and loyalty prediction in e-commerce. Prabandhan: Indian Journal of Management, 17(10).
- 3. Deniz, E., & Çökekoğlu Bülbül, S. (2024). Predicting customer purchase behavior using machine learning models. Information Technology in Economics and Business, 1(1).
- 4. Gustriansyah, R., Alie, J., & Suhandi, N. (2024). A hybrid machine learning model for market clustering. Engineering, Technology & Applied Science Research, 14(6), 18824–18828.
- Chaubey, B., Bisen, D., Arjaria, M., & Yadav, R. (2022). Predicting customer purchasing behavior using traditional and ensemble machine learning techniques. Journal of Retailing and Consumer Services, 66, 102878.
- Madani, S., & Alshraideh, H. (2021). Application of classical machine learning algorithms to predict consumer purchasing decisions in online food delivery. Journal of Retailing and Consumer Services, 59, 102365.
- Saura, J. R., Palacios-Marqués, D., & Ribeiro-Soriano, D. (2023). Privacy concerns in social media UGC communities: Understanding user behavior sentiments in complex networks. Information Systems and e-Business Management, 1-21.

- Ahmed, C., ElKorany, A., & ElSayed, E. (2023). Prediction of customer's perception in social networks by integrating sentiment analysis and machine learning. Journal of Intelligent Information Systems, 60(3), 829-851.
- 9. Zhang, Q., Zhang, Z., Yang, M., & Zhu, L. (2021). Exploring coevolution of emotional contagion and behavior for microblog sentiment analysis: a deep learning architecture. Complexity, 2021(1), 6630811.
- 10. Tadesse, M. M., Lin, H., Xu, B., & Yang, L. (2018). Personality predictions based on user behavior on the Facebook social media platform. IEEE Access, 6, 61959-61969.
- 11. Guimaraes, R. G., Rosa, R. L., & Bressan, G. (2017). Age groups classification in social network using deep learning. IEEE Access, 5, 10805-10816.
- 12. İş, H., & Tuncer, T. (2021). A Profile Analysis of User Interaction in Social Media Using Deep Learning. Traitement du signal, 38(1).
- Fan, H., Du, W., Dahou, A., Yousri, D., Elaziz, M. A., & Al-qaness, M. A. (2021). Social media toxicity classification using deep learning: real-world application UK Brexit. Electronics, 10(11), 1332.
- 14. Hayat, M. K., Daud, A., Alshdadi, A. A., & Banjar, A. (2019). Towards deep learning prospects: insights for social media analytics. IEEE Access, 7, 36958-36979.
- 15. Mohbey, K. K. (2020). Multi-class approach for user behavior prediction using deep learning framework on Twitter election dataset. Journal of Data, Information and Management, 2(1), 1-14.