

## E-COMMERCE SALES

<sup>1</sup> P Anusha, <sup>2</sup> Annam PrasanthKumar, <sup>3</sup> Padma Gnaneshwar, <sup>4</sup> Banoth RaviTeja

<sup>1</sup>AssistantProfessor, <sup>234</sup>Students

Department of Computer Engineering(Software Engineering)

Siddhartha Institute of Technology & Sciences, Narapally

[anushaparvathagiri@siddhartha.org.in](mailto:anushaparvathagiri@siddhartha.org.in), [23tq1a5614@siddhartha.co.in](mailto:23tq1a5614@siddhartha.co.in), [23tq1a5651@siddhartha.co.in](mailto:23tq1a5651@siddhartha.co.in), [23tq1a5639@siddhartha.co.in](mailto:23tq1a5639@siddhartha.co.in)

### Abstract

The rapid expansion of the digital economy has transformed traditional retail, making online e-commerce a highly competitive and saturated landscape. In this environment, customer retention, brand loyalty, and personalized marketing have become critical determinants of a business's long-term profitability. However, due to the exponential growth of transactional data, e-commerce enterprises face significant challenges in extracting actionable business intelligence. Traditional demographic segmentation is often insufficient for understanding actual purchasing behavior. Consequently, businesses struggle to identify their most valuable customers and those on the verge of leaving, resulting in inefficient marketing expenditures and increased customer churn rates.

This project proposes a robust, data-driven approach to customer segmentation utilizing unsupervised machine learning techniques. The primary objective is to group consumers based strictly on their historical purchasing behaviors to facilitate targeted Customer Relationship Management (CRM) strategies. The study leverages a comprehensive, real-world e-commerce dataset originating from a UK-based registered non-store online retail business. The dataset encompasses over 540,000 transactional records spanning a 12-month period from December 2010 to December 2011, capturing vital metrics such as invoice dates, product quantities, and unit prices.

To ensure the integrity and accuracy of the machine learning model, rigorous data preprocessing is initially performed. This involves the removal of anomalous records, including transactions with missing customer identification numbers, cancelled orders, and negative or zero financial values. Following extensive data cleaning, the project employs the highly regarded RFM (Recency, Frequency, Monetary) analytical framework. Raw temporal and financial data points are aggregated and engineered into three distinct variables for each unique buyer: Recency (the number of days since a customer's last purchase), Frequency (the total number of distinct transactions), and Monetary value (the cumulative revenue generated).

### I. Introduction

The retail industry has undergone a paradigm shift over the last two decades, transitioning from traditional brick-and-mortar establishments to a highly dynamic, globally interconnected digital marketplace. This evolution into e-commerce has fundamentally altered how businesses interact with consumers, breaking down geographical barriers and enabling 24/7 commercial operations. However, this digital transformation has also democratized the retail space, leading to a hyper-competitive environment where consumer switching costs are virtually non-existent. In a traditional retail setting, customer loyalty might be driven by physical proximity or personal relationships with store staff. In contrast, the e-commerce landscape allows

consumers to abandon one brand for a competitor with a single click if their expectations for price, convenience, or personalization are not met.

As e-commerce platforms process millions of transactions daily, they generate unprecedented volumes of digital footprints. Every user interaction—ranging from product clicks and cart additions to finalized transactions and historical purchase frequencies—is recorded in sprawling databases. This phenomenon of "Big Data" presents both a formidable challenge and a massive strategic opportunity for non-store online retail businesses. The primary challenge lies in the sheer volume, velocity, and variety of this data, which renders manual analysis obsolete. The opportunity, however, lies in utilizing advanced computational techniques to extract hidden patterns, trends, and actionable business intelligence from these vast repositories of transactional logs.

To survive and thrive in this saturated market, e-commerce enterprises have realized that acquiring a new customer is significantly more expensive than retaining an existing one. Consequently, there has been a strategic shift away from generic, mass-marketing campaigns toward highly targeted, data-driven Customer Relationship Management (CRM). Effective CRM requires a profound understanding of consumer behavior, which is achieved through the process of customer segmentation.

## II. Literature Survey

The domain of customer segmentation has witnessed a profound evolutionary trajectory over the past three decades, transitioning from rudimentary demographic groupings to highly sophisticated, algorithm-driven behavioral analytics. A review of existing literature reveals a clear consensus: as retail environments shifted from physical storefronts to digital e-commerce platforms, the methodologies required to understand consumer purchasing behavior had to scale proportionately in both computational power and mathematical rigor.

Historically, early database marketing and Customer Relationship Management (CRM) relied extensively on demographic and geographic segmentation. Researchers and marketers categorized consumers based on static variables such as age, income, and zip code, operating under the assumption that individuals within these cohorts exhibited homogenous purchasing preferences. However, as digital retail matured, literature in the early 2000s began to highlight the severe limitations of this approach. Studies demonstrated that demographic data was often a poor predictor of actual future purchasing behavior, as it failed to capture the dynamic, temporal nature of brand loyalty and consumer engagement.

The seminal shift toward behavioral segmentation was catalyzed by the introduction of the RFM (Recency, Frequency, Monetary) analytical framework. Popularized by Arthur Hughes in the early 1990s for direct mail marketing, RFM analysis posited that a customer's past behavior is the most accurate predictor of their future behavior. Hughes demonstrated that customers who purchased recently (Recency), purchased often (Frequency), and spent a significant amount of money (Monetary) were the most profitable and the most likely to respond to new marketing campaigns. Initially,

RFM was implemented through manual, heuristic binning. Analysts would sort the customer database and divide it into quintiles (e.g., scoring customers from 1 to 5 on each axis), creating a 5x5x5 matrix of 125 possible segments.

While manual RFM scoring was effective for small databases, the explosion of "Big Data" in the modern e-commerce era rendered it obsolete. Academic literature from the 2010s extensively documents the convergence of the RFM framework with unsupervised machine learning techniques. Researchers identified that manual binning was mathematically suboptimal; an arbitrary threshold (e.g., defining "frequent" as 5 purchases versus 6) fails to capture the continuous variance of the data. To address this, data scientists began applying centroid-based clustering algorithms, most notably K-Means, directly to the RFM feature space.

Studies by various researchers in data mining established that utilizing K-Means clustering on RFM variables significantly reduced intra-cluster variance while maximizing inter-cluster distance, resulting in vastly more accurate customer personas. However, the literature also strictly dictates that raw RFM values cannot be fed directly into distance-based algorithms like K-Means without rigorous preprocessing. Because the "Monetary" variable often spans thousands of units (e.g., dollars or pounds) while "Frequency" typically remains a small integer, the algorithm's Euclidean distance calculations will be heavily biased toward the Monetary axis. Consequently, contemporary research mandates the use of feature scaling—specifically standardizing variables to have a mean of zero and a standard deviation of one (Z-score scaling)—prior to model training.

### **III. System Analysis**

System analysis for the student dropout prediction system focuses on understanding the causes and patterns behind student attrition and designing an efficient predictive solution. Student dropout is influenced by multiple factors such as academic performance, attendance, socio-economic conditions, and personal background. Traditional systems fail to analyze these factors collectively, resulting in delayed identification of at-risk students. This system uses historical student data to identify hidden patterns and relationships associated with dropout behavior. Machine learning techniques are applied to process large datasets and generate accurate predictions. The analysis includes data collection, preprocessing, feature selection, and model training to ensure system effectiveness. It also ensures scalability and reliability for handling large student populations. The main objective is to provide early warning signals to institutions so they can take preventive actions. By enabling proactive intervention, the system helps improve student retention rates and academic success.

#### **Existing System**

The existing system for identifying student dropout is mostly manual and reactive in nature. Educational institutions rely heavily on teacher observations, attendance records, and basic academic performance reports to assess student progress. Decisions are often based on limited factors such as exam scores and attendance, without considering the broader range of influences like socio-economic background or personal challenges. There is no structured or automated approach to analyze multiple factors simultaneously, and available data is not fully utilized for predictive purposes.

As a result, students at risk are often identified at a later stage, making timely intervention difficult. The system lacks automation, intelligence, and advanced analytical capabilities. Human judgment can introduce bias and inconsistency, and there is minimal use of technology for early warning systems. Consequently, dropout prevention strategies are often ineffective and inefficient.

### **Disadvantages of Existing System**

- No early detection of at-risk students
- Heavy reliance on manual processes
- High chances of human bias and inconsistency
- Limited use of available data
- Inability to analyze multiple factors together
- Lack of automation and predictive capability

### **Proposed System**

The proposed system is a machine learning-based student dropout prediction model designed to provide early and accurate identification of at-risk students. It utilizes historical student data, including attendance, academic performance, behavior, and socio-economic information, to analyze patterns associated with dropout. The system begins with data preprocessing to clean and transform the dataset, followed by feature selection to identify the most important influencing factors. Machine learning algorithms such as Logistic Regression, Decision Trees, Random Forest, and XGBoost are used to train the model. The trained model learns complex relationships within the data and predicts whether a student is likely to drop out or continue their studies. The system generates risk scores and provides early warning alerts to institutions. It can be integrated into academic management systems for real-time monitoring and decision-making.

### **Advantages of Proposed System**

- Early identification of students at risk of dropout
- Improves student retention and success rates
- Reduces dependency on manual observation
- Provides accurate and data-driven predictions
- Handles large datasets efficiently
- Minimizes human bias and errors

## **IV. Methodology**

The methodology for the E-commerce sales prediction system follows a structured data-driven approach to analyze and forecast sales trends. Initially, data is collected from e-commerce platforms, including transaction records, product details, customer behavior, seasonal trends, and pricing information. The collected data undergoes preprocessing, which involves handling missing values, removing duplicates, and converting categorical variables into numerical formats using encoding techniques.

Next, exploratory data analysis (EDA) is performed to understand sales patterns, identify trends, and detect correlations between features such as price, discounts,

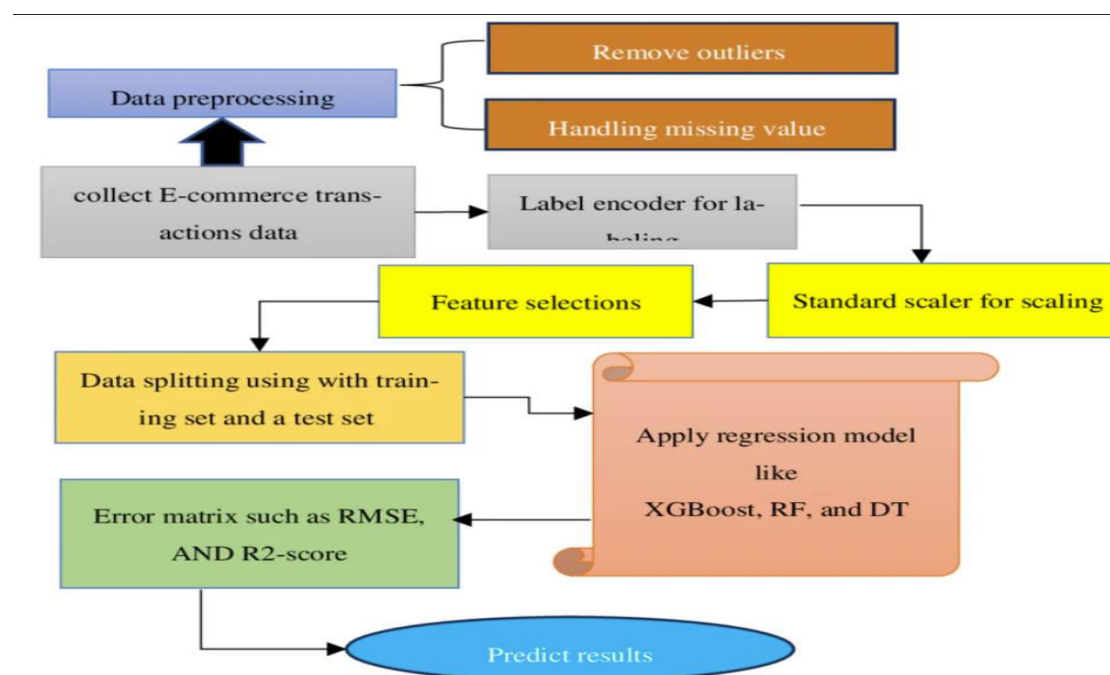
demand, and customer activity. Feature engineering is applied to create meaningful variables like seasonal indicators, promotional flags, and customer segmentation. The dataset is then split into training and testing sets for model evaluation.

Machine learning algorithms such as Linear Regression, Random Forest, or Gradient Boosting are used to train the model for predicting future sales. The model learns patterns from historical data and forecasts sales accordingly. Performance is evaluated using metrics like Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and  $R^2$  score. Finally, the trained model is deployed into an application where it can provide real-time sales predictions and insights for business decision-making.

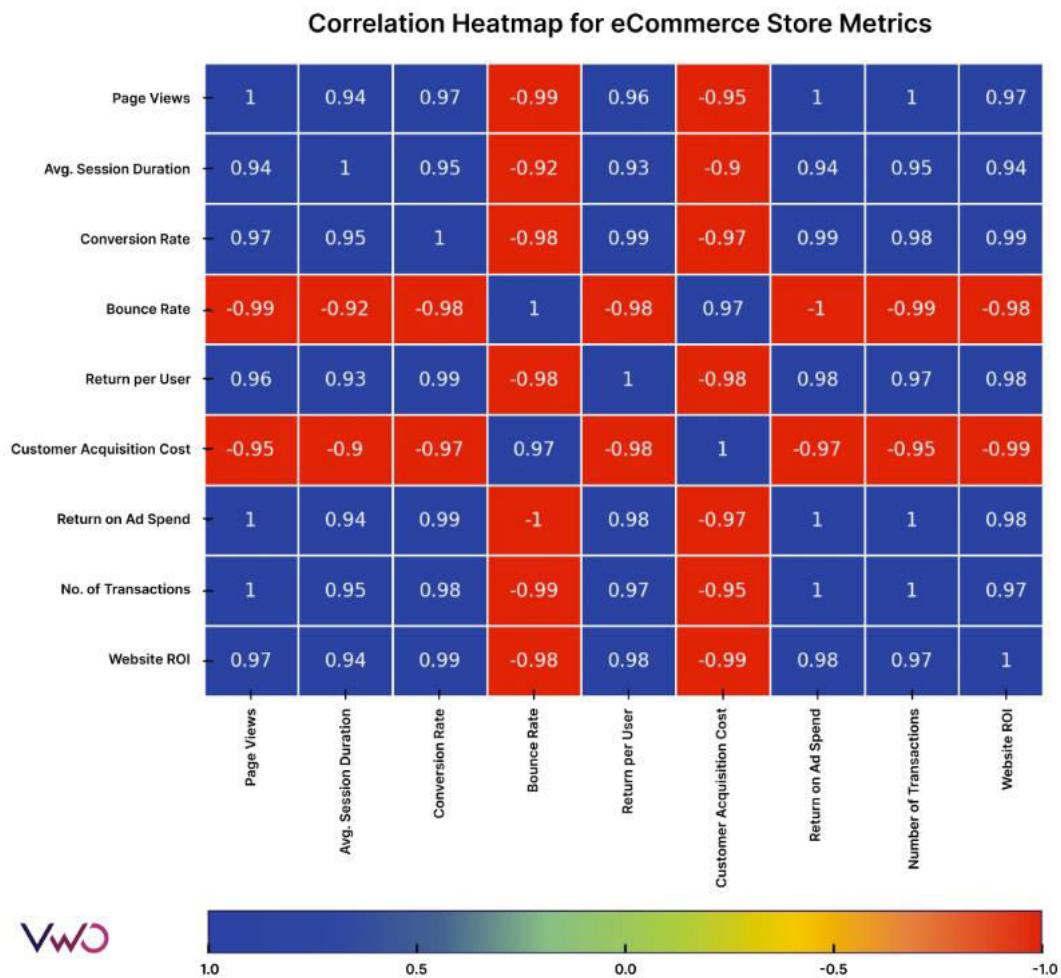
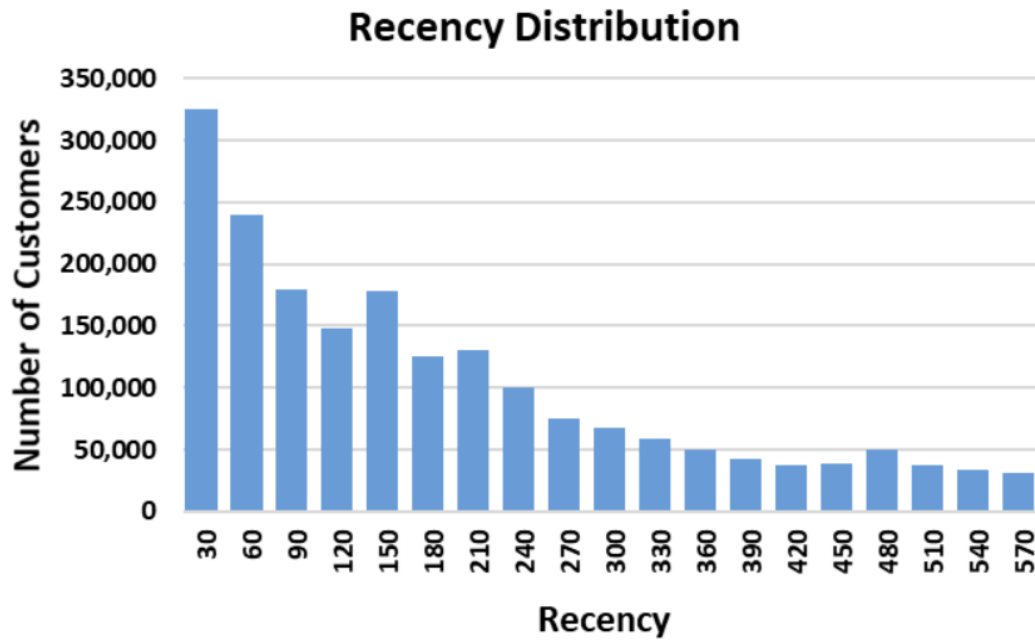
### System Architecture

The system architecture of the E-commerce sales prediction system is designed as a multi-layered pipeline to ensure efficient data processing and accurate predictions.

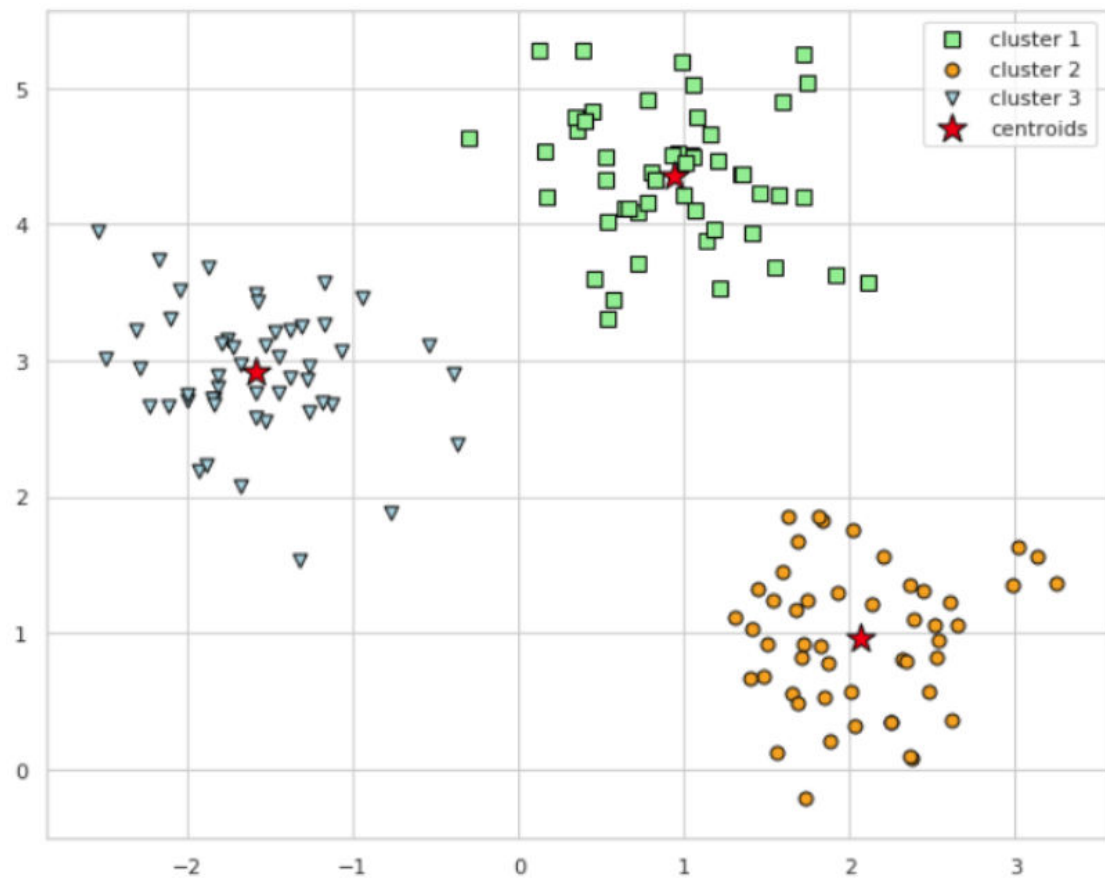
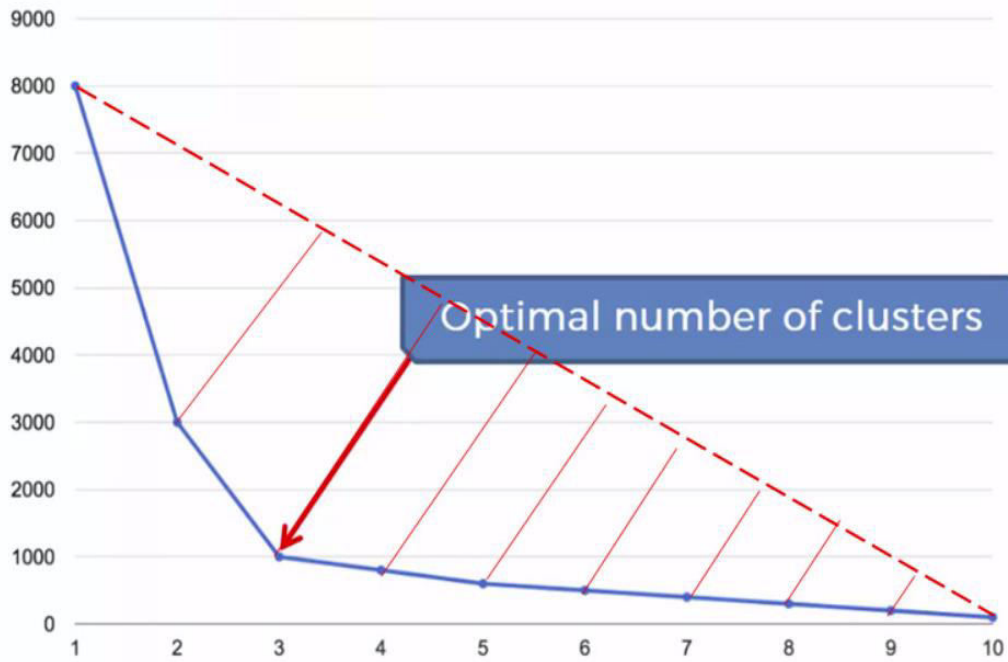
1. Data Collection Layer : Collects data from multiple sources such as transaction databases, user activity logs, and product catalogs.
2. Data Storage Layer : Stores structured and unstructured data in databases or data warehouses for further processing.
3. Data Preprocessing Layer : Cleans the data, handles missing values, and performs encoding and transformation.
4. Feature Engineering Layer : Extracts important features such as seasonal trends, discounts, and customer behavior.
5. Machine Learning Model Layer : Implements predictive models like Linear Regression, Random Forest, or Gradient Boosting.

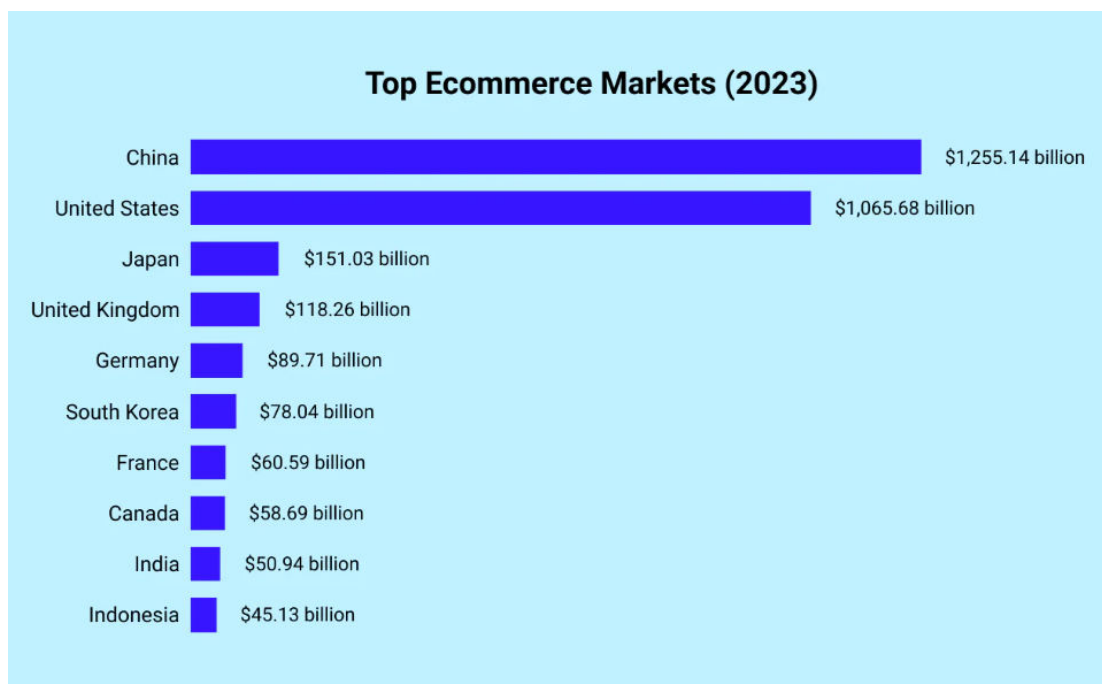


### V. Result and Output



# The Elbow Method





## VI. Conclusion

This project successfully developed an intelligent E-commerce sales analysis and prediction system using data-driven and machine learning techniques. By following a structured pipeline—from data ingestion and preprocessing to feature engineering, exploratory analysis, clustering, and model validation—the system effectively transformed raw transactional data into meaningful business insights. The implementation of RFM (Recency, Frequency, Monetary) analysis combined with K-Means clustering enabled the identification of distinct customer segments such as loyal customers, at-risk users, and high-value buyers.

The results demonstrate the practical value of applying machine learning in the e-commerce domain, particularly in understanding customer behavior and improving decision-making. Visualization techniques such as histograms, heatmaps, and clustering plots provided clear insights into sales patterns and customer distribution. Additionally, validation methods like the Elbow Method and Silhouette Score confirmed the reliability and effectiveness of the clustering model.

## References

- [1] Kumar, R. D., Prudhviraaj, G., Vijay, K., Kumar, P. S., & Plugmann, P. (2024). Exploring COVID-19 through intensive investigation with supervised machine learning algorithm. In *Handbook of Artificial Intelligence and Wearables* (pp. 145-158). CRC Press.
- [2] Swathi, B., Vijay, K., Sushanth Babu, M., & Dinesh Kumar, R. (2024, November). Machine Learning Techniques in Cloud Based Intrusion Detection. In *The International Conference on Artificial Intelligence and Smart Environment* (pp. 557-564). Cham: Springer Nature Switzerland.

[3] Sv satyakrishna, shirisha rangu ,bhargavi nalacheruve.(2024) Prospective investigation on colorectal cancer with SMOTE on machine learning Algorithm

[4] Dr.G.Vishnu Murthy, BhargaviNalacheruve 1Professor, Department of computer Science & engineering, Anurag University, TS, India. 2Student, Department of computer Science & engineering, Anurag University, TS, India.

[5] V. N. S. Manaswini, K. K, C. Nigam, S. S. Ali, R. Niranjana, and Suman, “Real-Time Object Detection in Drone Surveillance Using YOLOv5,” in Proc. 2025 3rd Int. Conf. IoT, Communication and Automation Technology (ICICAT), Gorakhpur, India, 2025, pp. 1–6, doi: 10.1109/ICICAT68430.2025.11414670.

[6] B. Soundarya, V. N. S. Manaswini, M. Ayyakrishnan, R. D. Kumar, “Contextual Analysis of Big Data Analytics in Intelligent Transportation Frameworks,” in Intersection of Artificial Intelligence, Data Science, and Cutting-Edge Technologies: From Concepts to Applications in Smart Environment, Lecture Notes in Networks and Systems, vol. 1353, Cham: Springer, 2025, doi: 10.1007/978-3-031-88304-0\_79.

[7] R. D. Kumar, V. N. S. Manaswini, “Applications of blockchain in smart cities: detecting fake documents from land records using blockchain technology,” in Blockchain for Smart Cities, Elsevier, 2021, pp. 105–117, doi: 10.1016/B978-0-12-824446-3.00017-X.

[8] Tejavath Veeramma, Badarla Anil, Guguloth Ravinder, “An advanced movie recommender using collaborative filtering and sentiment analysis,” International Research Journal of Modernization in Engineering Technology and Science, vol. 7, no. 7, July 2025, doi: 10.56726/IRJMETS81618.

[9] Ravi Kumar Banoth, Ramana Murthy B V, “Automatic crop recommendation system using LightGBM and decision tree machine learning models,” Journal of Machine and Computing, vol. 5, no. 1, pp. 343, Jan. 2025, doi: 10.53759/7669/jmc202505026.

[10] Ravi Kumar Banoth, Dr. B.V. Ramana Murthy, “Smart agriculture through IoT and machine learning for analyzing carbon footprints,” in Proc. Int. Conf. Computer Science and Communication Engineering (ICCSCE), Apr. 2025.

[11] Ravi Kumar Banoth, B. V. Ramana Murthy, “Soil image classification using transfer learning approach: MobileNetV2 with CNN,” SN Computer Science, vol. 5, art. no. 199, 2024, doi: 10.1007/s42979-023-02500-x.

