

# Retail Stock Market Behavior: Customer Segmentation and Sales Pattern Analysis Using Data Mining Techniques

1<sup>st</sup> Tejmul Movin

*Department of CS and AI  
Rishihood University*

Sonipat, India

sansar.t23csai@nst.rishihood.edu.in

2<sup>nd</sup> A Jithendranath

*Department of CS and AI  
Rishihood University*

Sonipat, India

akula.j23csai@nst.rishihood.edu.in

3<sup>rd</sup> M Sree Sai Nath

*Department of CS and AI  
Rishihood University*

Sonipat, India

meesala.s23csai@nst.rishihood.edu.in

**Abstract**—This paper presents a comprehensive data-mining analysis of the UCI Online Retail Dataset to uncover patterns in customer behavior, product performance, and market trends. The project validates 15 testable hypotheses across six analytical dimensions: product performance and revenue dynamics, temporal patterns and seasonality, customer segmentation and behavior, association rule mining, predictive modeling and interpretability, and geographic and market expansion. We apply market basket analysis using the Apriori algorithm, customer segmentation using RFM modeling and K-Means clustering, and predictive models including basic time-series models, XGBoost, and Random Forest. Key findings demonstrate that higher prices correlate with lower purchase quantities, seasonal peaks occur during Q4, customers can be segmented into distinct behavioral groups, and specific product pairs are frequently purchased together. The analysis provides actionable insights for inventory management, pricing strategies, customer targeting, and cross-selling opportunities in retail operations.

**Index Terms**—retail analytics, data mining, customer segmentation, RFM analysis, association rule mining, market basket analysis, predictive modeling

## I. INTRODUCTION

Retail businesses generate large volumes of transactional data daily. Analyzing purchasing patterns hidden in this data can help retailers optimize inventory, pricing, and customer engagement strategies. Understanding customer behavior, product relationships, and sales dynamics is crucial for business decision-making in the competitive retail sector.

This project analyzes the UCI Online Retail Dataset spanning December 2010 to December 2011, containing transactional records from a UK-based online retail store. The analysis explores six core research dimensions to answer fundamental business questions: What drives sales? When do customers buy? Who are the key customer groups? What products are bought together? How can future behavior be predicted? Where do growth opportunities exist?

The paper is organized as follows: Section II presents the literature review, Section III describes the methodology and analytical framework, Section IV outlines the data preprocessing approach, Section V presents the research questions, Section VI presents the hypotheses and testing approach,

Section VII discusses the implementation plan, and Section VIII concludes with expected outcomes and implications.

## II. LITERATURE REVIEW

### A. Market Basket Analysis and Association Rule Mining

Market Basket Analysis (MBA) is the foundational technique for identifying product relationships. The Apriori algorithm is the classic approach for discovering frequent itemsets and generating association rules, utilizing metrics such as Support, Confidence, and Lift. Lift measures how much more likely item Y is purchased given item X, relative to their independent probabilities.

The FP-Growth algorithm provides a more memory and time-efficient alternative to Apriori. These techniques are essential for validating hypotheses H9 and H10, which focus on identifying frequently co-purchased products and assessing cross-selling opportunities.

### B. Customer Segmentation and RFM Modeling

The Recency, Frequency, and Monetary (RFM) model is the widely accepted standard for quantifying customer value and behavior. Recency measures time since last purchase, Frequency counts total purchases, and Monetary represents total spending. K-Means Clustering is the dominant technique for partitioning customers based on normalized RFM scores.

Methods such as the Elbow Method and Silhouette Score serve as standard practices for objectively determining the optimal number of customer segments. This approach directly supports hypothesis H7, which aims to identify distinct customer groups with different purchasing behaviors.

### C. Predictive Modeling and Time-Series Analysis

Traditional basic time-series models serve as robust baselines for sales prediction. Ensemble Machine Learning models including XGBoost and Random Forest often demonstrate superior performance in forecasting tasks. For model interpretability, feature importance techniques are preferred for explaining predictions from tree-based models in business-understandable terms.

#### D. Retail Benchmarks and Empirical Principles

The Pareto Principle, commonly known as the 80/20 rule, is frequently validated in retail datasets, suggesting approximately 20% of products account for 80% of total revenue. The holiday shopping effect and other temporal cycles are well-documented phenomena in retail seasonality. These empirical observations guide expectations for hypotheses H3 and H4.

### III. METHODOLOGY

#### A. Dataset Overview

The UCI Online Retail Dataset contains transactional records from December 2010 to December 2011. The dataset includes the following fields: InvoiceNo (unique transaction identifier), StockCode (product code), Description (product name), Quantity (units sold per transaction), InvoiceDate (date and time of transaction), UnitPrice (price per unit), CustomerID (unique customer identifier), and Country (customer country of residence).

#### B. Data Preprocessing

Preprocessing tasks are structured as follows.

1) *Handling Missing Values*: Rows without customer IDs, representing approximately 25% of records, are excluded from customer-level analyses but retained for aggregated product and sales insights. Rows with missing descriptions are dropped after verifying minimal impact on dataset size.

2) *Outlier and Anomaly Treatment*: Negative quantities indicating returns or cancellations are identified and processed separately for net sales calculations. These records are excluded from Association Rule Mining to focus on actual purchases. Zero or implausible UnitPrice values indicating data entry errors are removed prior to analysis.

3) *Feature Engineering*: The following engineered features are created: TotalSales = Quantity  $\times$  UnitPrice, temporal features including Year, Month, DayOfWeek, and HourOfDay extracted from InvoiceDate, and RFM metrics including Recency, Frequency, and Monetary computed per CustomerID.

### IV. RESEARCH QUESTIONS

The analysis is structured around answering core research questions across six dimensions.

#### A. Dimension A: Product Performance and Revenue Dynamics

- How does pricing influence the volume of units sold?
- Can we accurately predict high-value orders based on simple product metrics?
- Does the Pareto Principle (80/20 rule) hold true for product revenue distribution?

#### B. Dimension B: Temporal Patterns and Seasonality

- What is the impact of the Christmas season on overall sales volume?
- Are there significant differences in customer purchasing behavior between weekdays and weekends?
- Do purchase activities cluster around specific times of the day?

#### C. Dimension C: Customer Segmentation and Behavior

- Can customers be reliably segmented into distinct groups based on their purchasing characteristics (RFM)?
- Do purchasing behaviors, such as average order value and product preference, vary significantly across different countries?

#### D. Dimension D: Association Rule Mining

- Which specific product combinations frequently appear together in the same transaction?
- What is the potential revenue uplift from implementing cross-selling recommendation strategies?

#### E. Dimension E: Predictive Modeling and Interpretability

- Which transactional features (e.g., price, quantity, category) are the strongest predictors of order value?
- How does the performance of ensemble machine learning models compare to basic time-series models for sales forecasting?
- Can a customer's past purchasing activity accurately predict their future spending behavior?

#### F. Dimension F: Geographic and Market Expansion

- Is the total sales volume concentrated in a small number of key countries?
- Can countries be classified into distinct market types (e.g., Premium, Emerging) based on their purchasing metrics?

### V. HYPOTHESES AND TESTING FRAMEWORK

#### A. Dimension A: Product Performance and Revenue Dynamics

1) *H1: Price-Quantity Relationship*: Hypothesis: Higher-priced products sell in smaller quantities compared to lower-priced ones. Testing Approach: Compare unit price with quantity sold per product, visualize the relationship between price and quantity, and analyze variation across product categories. Expected Outcome: Higher prices correlate with fewer units sold.

2) *H2: High-Value Order Prediction*: Hypothesis: Product price, quantity, and category can predict whether an order exceeds £500. Testing Approach: Label orders above £500 as "big orders," train a machine learning model to classify order size, and measure accuracy while identifying key predictive factors. Expected Outcome: Model predicts high-value orders with  $\geq 75\%$  accuracy.

3) *H3: Pareto Principle in Retail*: Hypothesis: A small fraction of products account for most revenue. Testing Approach: Rank products by total sales, compute cumulative revenue contribution, and test Pareto distribution validity. Expected Outcome: Top 20% of products generate approximately 80% of total revenue.

## B. Dimension B: Temporal Patterns and Seasonality

1) *H4: Christmas Season Sales Impact:* Hypothesis: Sales rise during Q4 (October–December) compared to other quarters. Testing Approach: Aggregate monthly revenue, visualize trends across the year, and identify seasonal peaks and low periods. Expected Outcome: Peak activity during the holiday season with drops post-January.

2) *H5: Weekend versus Weekday Behavior:* Hypothesis: Customer purchasing behavior varies between weekends and weekdays. Testing Approach: Classify orders as weekday or weekend and compare average order values and frequency. Expected Outcome: Distinct spending or frequency patterns across days.

3) *H6: Daytime Purchase Clustering:* Hypothesis: Purchases cluster during daytime hours. Testing Approach: Extract transaction hours from timestamps, aggregate order counts and values per hour, and identify peak shopping hours. Expected Outcome: Increased activity during working or lunch hours.

## C. Dimension C: Customer Segmentation and Behavior

1) *H7: Customer Grouping by Behavioral Traits:* Hypothesis: Customers can be grouped into types by behavioral traits (frequency, spending, recency). Testing Approach: Compute RFM scores, perform K-Means clustering to identify segments, and characterize each cluster based on purchasing behavior. Expected Outcome: Three to four clear customer groups with distinct spending habits.

2) *H8: Geographic Purchasing Differences:* Hypothesis: Purchasing behavior differs significantly across countries. Testing Approach: Analyze top five countries by sales volume, compare average order value, frequency, and product preferences. Expected Outcome: Regional variations in product preferences and order values.

## D. Dimension D: Association Rule Mining

1) *H9: Frequent Product Co-occurrence:* Hypothesis: Certain product pairs co-occur in the same transaction frequently. Testing Approach: Generate association rules using market basket analysis, calculate support, confidence, and lift for each pair. Expected Outcome: Identification of 10–20 high-support product combinations.

2) *H10: Cross-Selling Effectiveness:* Hypothesis: Recommending complementary products increases transaction value. Testing Approach: Simulate adding top product pairs as recommendations and estimate revenue uplift based on cross-sell frequency. Expected Outcome: Recommendation strategies can raise average order value by 10–15%.

## E. Dimension E: Predictive Modeling and Interpretability

1) *H11: Revenue Influencing Factors:* Hypothesis: Price and quantity are the strongest predictors of order value. Testing Approach: Train regression and tree-based models for order value and use feature importance to rank predictors. Expected Outcome: Price and quantity dominate, followed by category and timing.

2) *H12: Machine Learning Model Performance:* Hypothesis: Machine learning models outperform basic time-series baselines in predicting future sales. Testing Approach: Build basic time-series, basic regression, and ensemble ML models, and compare prediction errors on unseen data. Expected Outcome: ML models yield lower error rates.

3) *H13: Customer Spending Prediction:* Hypothesis: Past six months' activity predicts future spending behavior. Testing Approach: Train model on first-half spending data and evaluate predictions against second-half results. Expected Outcome: Reliable identification of future high-value customers.

## F. Dimension F: Geographic and Market Expansion

1) *H14: Country-Level Revenue Contribution:* Hypothesis: A few countries account for most of the total sales volume. Testing Approach: Calculate and rank country-level revenue and analyze top contributors and monthly growth trends. Expected Outcome: Top five countries generate over 70% of revenue.

2) *H15: Market Type Classification:* Hypothesis: Countries can be classified as "Premium," "Emerging," or "Niche" based on purchasing behavior. Testing Approach: Aggregate metrics per country including average order value, product preference, and seasonality, then apply clustering. Expected Outcome: Discovery of three to four distinct market types with unique characteristics.

## VI. IMPLEMENTATION PLAN

The project execution follows a three-phase structure with defined leadership roles and deliverables.

### A. Phase 1: Planning, Documentation, and Data Understanding

Phase 1 establishes the project foundation through problem definition, dataset understanding, and preprocessing planning. Tejmul Movin serves as Phase 1 lead, responsible for coordinating documentation and research scope definition. Key deliverables include dataset description and rationale, data preprocessing plan, research objectives, hypothesis statements, methodology framework documentation, and literature review synthesis.

### B. Phase 2: Exploratory Data Analysis and Visualization

Phase 2 focuses on data exploration, cleaning, and pattern identification. A Jithendranath leads Phase 2, designing visualization dashboards and interpreting customer and product trends. Expected outputs include EDA notebooks, insights reports, correlation analysis, cleaned transaction datasets, and trend identification across all six analytical dimensions.

### C. Phase 3: Predictive Modeling, Evaluation, and Presentation

Phase 3 implements segmentation, association mining, and predictive models. M Sree Sai Nath leads Phase 3, responsible for model training, evaluation, and results interpretation. Phase 3 deliverables encompass customer segmentation results, association rule mining outputs, predictive model implementations

including classification and forecasting tasks, model performance evaluation, and final integrated analysis report.

#### D. Modeling Techniques

1) *Customer Segmentation:* K-Means clustering is applied to normalized RFM scores. The optimal number of clusters is determined using the Elbow Method and Silhouette Score metrics.

2) *Market Basket Analysis:* The Apriori algorithm is employed on cleaned transactions excluding returns. Association rules are prioritized based on high Lift and Confidence values for recommendation system development.

3) *Predictive Models:* For big order classification (H2), XGBoost Classifier is used with F1-Score, and Confusion Matrix evaluation metrics. For sales forecasting (H12), a basic time-series model serves as baseline with Gradient Boosting Regressor as the primary model, evaluated using RMSE. For feature influence analysis (H11), Random Forest Regressor with a feature importance ranking method is employed.

### VII. EXPECTED OUTCOMES

The analysis is expected to validate all 15 hypotheses through systematic testing. Product performance analysis should confirm inverse price-quantity relationships and validate the Pareto principle for revenue concentration. Temporal analysis should reveal clear seasonal patterns with Q4 peaks and distinct weekday/weekend variations. Customer segmentation should produce three to four actionable customer groups with different engagement patterns. Association analysis should identify multiple high-confidence product pairs suitable for bundling and cross-selling strategies. Predictive models should achieve or exceed accuracy targets and clearly identify the relative importance of different features. Geographic analysis should demonstrate revenue concentration in a small number of countries and enable market classification for targeted expansion strategies.

### VIII. CONCLUSION

This project applies comprehensive data-mining techniques to understand retail market behavior through systematic hypothesis testing. By analyzing product performance, temporal patterns, customer behavior, product associations, predictive capabilities, and geographic expansion opportunities, the project provides a structured framework for retail analytics.

The integration of RFM-based customer segmentation with market basket analysis creates opportunities for targeted customer engagement and revenue optimization. The predictive models, supported by feature importance analysis, bridge model performance with business interpretability.

The work demonstrates the applicability of established data-mining techniques to real-world retail datasets and provides actionable insights for inventory management, pricing strategies, customer targeting, and strategic market expansion.

### REFERENCES

- [1] Agrawal, R., and Srikant, R., "Fast algorithms for mining association rules," in *Proc. 20th Int. Conf. Very Large Data Bases*, 1994, pp. 487–499.
- [2] Han, J., Pei, J., and Yin, Y., "Mining frequent patterns without candidate generation," *ACM SIGMOD Rec.*, vol. 29, no. 2, pp. 1–12, 2000.
- [3] Dixon, C., and Wilkinson, N., *Marketing Analytics: Data Mining Techniques for Better Marketing Decisions*. Butterworth-Heinemann, 2011.