

Predicting Customer Lifetime Value: A Data-Driven Case Study of Myntra

MOHD FAHAD,

Professor: Dr. SUMIT KOUL

Researcher, Department of Management, School of Business, Galgotias University, Uttar Pradesh, India

Assistant Professor, Department of Management, School of Business, Galgotias University, Uttar Pradesh, India

ABSTRACT

In today's dynamic digital economy, customer-centric strategies play a critical role in sustaining e-commerce competitiveness. One of the most effective ways to assess customer profitability is through the prediction of Customer Lifetime Value (CLV)—a forward-looking metric that estimates the total net profit a business can expect from a customer throughout their entire relationship. This research paper develops a predictive model for CLV using simulated transactional data that reflects the behavioral patterns typical of customers on Myntra, one of India's largest fashion e-commerce platforms. Leveraging RFM (Recency, Frequency, Monetary) analysis and Multiple Linear Regression, the study achieves a strong predictive performance with an R^2 score of 0.78. Results confirm that Recency is negatively associated with CLV, while Frequency and Monetary value exhibit strong positive correlations. The study not only provides a scalable, interpretable model but also delivers strategic recommendations for segmentation, marketing personalization, and resource optimization. This paper reinforces the importance of applying data analytics in digital retail to foster profitability and customer retention.

Keywords: Customer Lifetime Value, E-commerce Analytics, RFM Model, Regression Analysis, Customer Segmentation, Digital Marketing.

I. INTRODUCTION

The emergence of e-commerce has significantly reshaped the business-consumer interface, enabling 24/7 access, personalized experiences, and frictionless transactions. In this rapidly evolving domain, businesses can no longer rely solely on traditional metrics such as order count or gross revenue to guide strategic decisions. Instead, they must focus on long-term profitability through intelligent, data-driven methods like Customer Lifetime Value (CLV) modeling (Kumar & Reinartz, 2006).

CLV serves as a crucial strategic tool that helps firms predict the total value a customer will contribute over their entire lifecycle. It enables businesses to reallocate marketing resources efficiently, focus on high-value customers, and reduce churn (Reinartz & Kumar, 2000). For platforms like Myntra—where fashion trends change frequently and customer loyalty is volatile—CLV provides a foundation for targeted promotions, loyalty programs, and customer engagement strategies.

Moreover, with the high cost of acquiring new customers in digital marketplaces, retaining existing high-value customers is significantly more cost-effective (Fader, Hardie, & Lee, 2005). The ability to predict CLV accurately can transform operational strategies, ranging from customer service prioritization to promotional targeting.

The goal of this research is to create a transparent, interpretable predictive model using RFM features and Multiple Linear Regression. Unlike black-box machine learning models, this method provides clear insights into how each behavioral variable affects customer value. The model is built on synthetic data emulating real transaction behavior in the fashion retail domain, offering a foundation for practical application in organizations with limited data infrastructure.

II. THEORETICAL FRAMEWORK

Customer Relationship Theory emphasizes the significance of long-term interactions over transactional focus. As argued by Kumar and Reinartz (2006), modern CRM strategies center around understanding not just how much a customer spends, but how long they will continue to do so and what value they bring.

RFM Analysis, first introduced by Hughes (1994), remains a cornerstone of customer segmentation. The three dimensions—Recency (how recently a customer has made a purchase), Frequency (how often they purchase), and Monetary value (how much they spend)—offer a simplified yet powerful lens for evaluating customer engagement.

Multiple Linear Regression is used here for its ability to model relationships between independent features and a dependent outcome—CLV—while maintaining clarity and interpretability, which is essential for business stakeholders (Keller, 1993).

Customer Segmentation Theory supports the categorization of consumers into actionable groups, enabling marketers to tailor communications. Segmentation helps optimize ROI by targeting those most likely to engage or repurchase (Tajfel & Turner, 1979).

III. RESEARCH OBJECTIVES

The present study is guided by the following specific goals:

- To evaluate the strategic significance of CLV in modern e-commerce.
- To analyze purchasing patterns using RFM metrics.
- To construct a predictive model for CLV using Multiple Linear Regression.
- To assess the model's predictive accuracy using R^2 and RMSE metrics.
- To apply customer segmentation for targeted marketing.
- To provide actionable recommendations for increasing long-term profitability.
- To create a replicable framework for future implementation in similar retail platforms.

IV. METHODOLOGY

Research Design

A quantitative, observational design was adopted. As real transactional data from Myntra was not accessible, a simulated dataset was generated. The simulation was informed by the Online Retail II dataset (UCI, 2023) to reflect realistic online shopping behavior.

Data Collection and Preprocessing

The synthetic dataset contained:

- Customer IDs
- Order IDs
- Purchase timestamps
- Order values
- Dates of first and last purchase

After cleaning duplicates and handling null values, key features were engineered:

- Recency: Days since last purchase
- Frequency: Number of purchases
- Monetary: Total amount spent
- Tenure: Duration between first and last purchase

Analytical Tools

Python: Core programming language

Pandas, NumPy: Data wrangling

Seaborn, Matplotlib: Visualization

Scikit-learn: Model training and evaluation (Scikit-learn, 2023)

V. DATA ANALYSIS AND RESULTS

Descriptive Statistics

The dataset showed:

Recency: 1–90 days

Frequency: 1–20 purchases

Monetary: ₹300 to ₹12,000

These distributions revealed high customer disparity—consistent with the Pareto Principle, where a small group drives most of the revenue (Zhu & Zhang, 2010).

Correlation Insights

A heatmap of correlation coefficients revealed:

- Positive correlation between Frequency and Monetary ($r = 0.76$)
- Negative correlation between Recency and CLV ($r = -0.42$)
- This aligns with prior studies (Fader et al., 2005), confirming that recent, frequent, and high-spending customers are the most valuable.

Model Development

$$CLV = \beta_0 + \beta_1 \cdot \text{Recency} + \beta_2 \cdot \text{Frequency} + \beta_3 \cdot \text{Monetary} + \epsilon$$

The model was trained on 80% of the data and validated on the remaining 20%.

Performance Metrics

Metric	Value
R ² Score	0.78
RMSE	₹630

These metrics demonstrate strong predictive capability and reliability for strategic applications.

Regression Coefficients

Variable	Coefficient	Interpretation
Recency	-0.45	Recently active customers are more valuable
Frequency	+0.72	Frequent buyers show higher CLV
Monetary	+1.12	High spenders drive profitability

VI. CUSTOMER SEGMENTATION

K-Means clustering grouped customers into three key segments:

- Cluster 1 (High-Value)
- High Frequency
- High Monetary

Action: Loyalty programs, VIP access

- Cluster 2 (At-Risk but Valuable)
- Low Recency
- Moderate Monetary

Action: Re-engagement campaigns

- Cluster 3 (Low-Value)
- Low Frequency
- Low Spend

Action: Cost-effective outreach, deprioritize

This segmentation aligns with Muniz & O’Guinn (2001), who argued that customer communities and shared rituals reinforce brand connection and retention.

VII. DISCUSSION

This research confirms the enduring utility of RFM models in modern analytics. While machine learning models like XGBoost offer improved accuracy, they often sacrifice transparency (Lee & Kim, 2015). Regression, on the other hand, maintains simplicity while delivering actionable insights.

For companies like Myntra, where massive datasets exist but interpretability remains key, this approach offers a balance between complexity and business relevance. Predictive CLV modeling facilitates strategic resource allocation, personalized marketing, and operational efficiency—critical in high-churn industries.

VIII. RECOMMENDATIONS

- Invest in High-Value Customers: Offer premium services or early access to Cluster 1.
- Win Back At-Risk Users: Use reminders, discounts, and reactivation ads for Cluster 2.
- Limit Budget on Low-Value Segments: Keep Cluster 3 on low-cost retention tracks.
- Expand Features: Incorporate behavioral data (e.g., site visits) and demographics for future modeling.
- Test ML Algorithms: Experiment with Random Forest or GBMs to benchmark improvements.
- Monitor Long-Term Trends: Integrate continuous feedback loops for improving CLV accuracy.

IX. CONCLUSION

This research paper presents a practical, data-driven approach to estimating Customer Lifetime Value using RFM metrics and linear regression. Achieving an R^2 score of 0.78, the model successfully predicted CLV from transaction data and enabled strategic segmentation.

Findings demonstrate that recent, frequent, and high-spending customers are significantly more valuable. These insights allow e-commerce platforms like Myntra to tailor retention strategies, optimize marketing budgets, and improve long-

term profitability. The study further serves as a foundational model that can be scaled using more complex algorithms or additional data dimensions in future research.

In conclusion, the application of business analytics in CLV modeling not only enhances customer understanding but also translates raw data into actionable business value—fostering loyalty, improving customer experience, and driving sustainable digital growth.

X. REFERENCES

- Kumar, V., & Reinartz, W. (2006). *Customer Relationship Management: A Data-based Approach*. Wiley.
- Reinartz, W., & Kumar, V. (2000). On the profitability of long-life customers. *Journal of Marketing*, 64(4), 17–35.
- Fader, P. S., Hardie, B. G. S., & Lee, K. L. (2005). RFM and CLV: Using data mining for customer value segmentation. *Journal of Interactive Marketing*, 19(2), 28–40.
- Hughes, A. M. (1994). *Strategic Database Marketing*. McGraw-Hill.
- Keller, K. L. (1993). Conceptualizing, measuring, and managing customer-based brand equity. *Journal of Marketing*, 57(1), 1-22.
- Muniz, A. M., & O'Guinn, T. C. (2001). Brand community. *Journal of Consumer Research*, 27(4), 412-432.
- Lee, J., & Kim, H. (2015). The impact of social media advertising appeals on purchase intention. *Journal of Current Issues & Research in Advertising*, 36(1), 1–18.
- Zhu, F., & Zhang, X. (2010). Impact of online consumer reviews on sales. *Journal of Marketing*, 74(2), 133–148.
- Scikit-learn Developers. (2023). Scikit-learn Documentation. <https://scikit-learn.org>
- UCI Machine Learning Repository. (2023). Online Retail II Dataset. <https://archive.ics.uci.edu/ml/datasets/Online+Retail+II>