

Assessment 1B: Big Data Analysis Report

Aditya Anant Deshpande - a1910571

Abstract

This report provides an initial analysis of a large-scale global e-commerce dataset. The objective is to identify customer segments, revenue trends, product return patterns, and geographic sales performance. Using RFM analysis, clustering, and visualisations, this report also refines the original research questions and proposes future data enhancements.

1. Part 1: Data Description

The dataset comprises more than 540,000 transaction records collected over a one-year span (December 2010 to December 2011) by a UK-based online retailer. It provides details on invoices, product items, purchase quantities, unit prices, customer identifiers, and the countries where transactions occurred.

Research Questions Restated

- How can customers be grouped based on their purchasing behavior to improve targeted marketing?
- What are the high-revenue time periods and how do promotions affect revenue?
- Which products show high return rates or pricing issues that indicate poor performance?
- What are the most profitable geographic regions for the business?

Key Features Identified

- **InvoiceNo:** Used to identify transactions and distinguish returns (prefixed with 'C').
- **StockCode and Description:** Uniquely identify products. Text-based 'Description' is also used to detect patterns.
- **Quantity and UnitPrice:** Combined to calculate **TotalRevenue**, a key measure for profitability.
- **InvoiceDate:** Used to assess sales over time and detect seasonal trends.

- **CustomerID:** Used for behavioral clustering and segment analysis. 25% missing values handled via filtering.
- **Country:** Used for geographic sales analysis and to explore global performance.

2. Part 2: Clustering and Pattern Analysis

Customer Segmentation with RFM

To understand customer behavior, Recency-Frequency-Monetary (RFM) analysis was performed:

- **Recency:** Days since last purchase (lower = more recent).
- **Frequency:** Number of distinct invoices per customer.
- **Monetary:** Total revenue from the customer.

Feature Analysis with Summary Statistics

The dataset includes three primary numerical features: **Quantity**, **UnitPrice**, and **TotalRevenue**. Summary statistics reveal:

- The average order includes 13 units, with considerable variance (std = 179).
- Unit prices range from as low as 0.001 to nearly 40,000, indicating potential outliers or wholesale pricing.
- The average revenue per transaction is approximately \$22.40.

A correlation matrix showed:

- A strong positive correlation between **Quantity** and **TotalRevenue** ($r > 0.9$).
- A moderate positive correlation between **UnitPrice** and **TotalRevenue**, suggesting both price and quantity influence revenue.

Customer Segmentation (RFM Clustering)

To group customers based on their purchasing behavior, Recency, Frequency, and Monetary (RFM) metrics were calculated. These metrics were then used to perform K-means clustering, and the resulting segments were visualized using a pairplot:

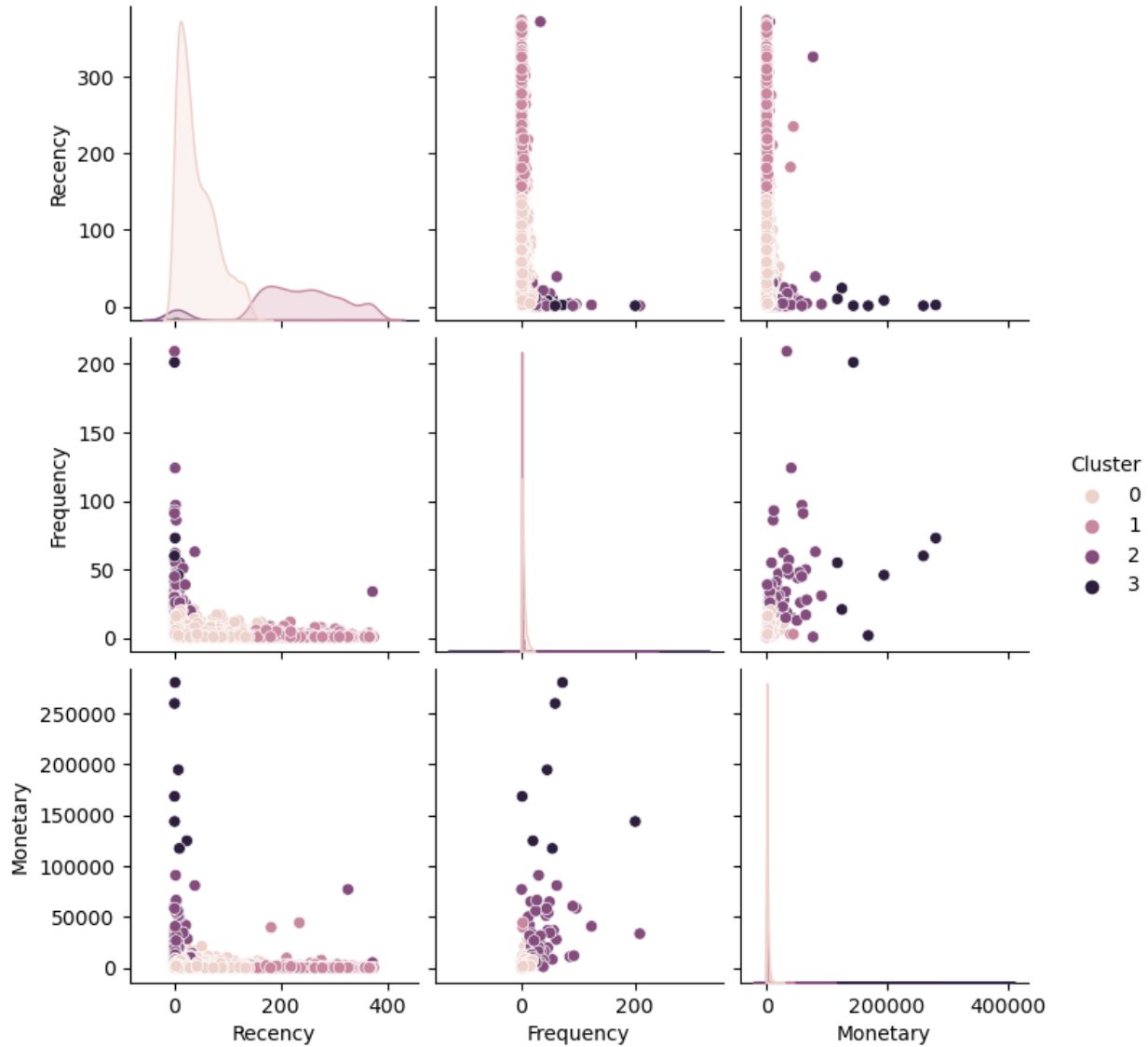


Figure 1: Pairplot of RFM Clusters Based on K-Means Segmentation

Cluster Interpretations

- **Cluster 0 (light pink):** Customers with high recency (long time since last purchase), low frequency, and low spending. These are likely **inactive or churned** users.
- **Cluster 1 (light purple):** Recently active customers with low frequency and monetary value. Possibly **new or testing** customers.
- **Cluster 2 (purple):** Customers with moderate recency, frequency, and higher spending. Represent **steady, valuable** customers.
- **Cluster 3 (dark purple):** Recent, high-frequency, and high-value customers. This group

represents **top-tier loyal customers** and should be prioritized for retention.

This segmentation enables the business to tailor strategies for each customer type, including win-back campaigns, onboarding flows, and VIP engagement.

3. Part 3: Visualisation Summary

The following figures were generated to explore temporal, product-level, and regional patterns.

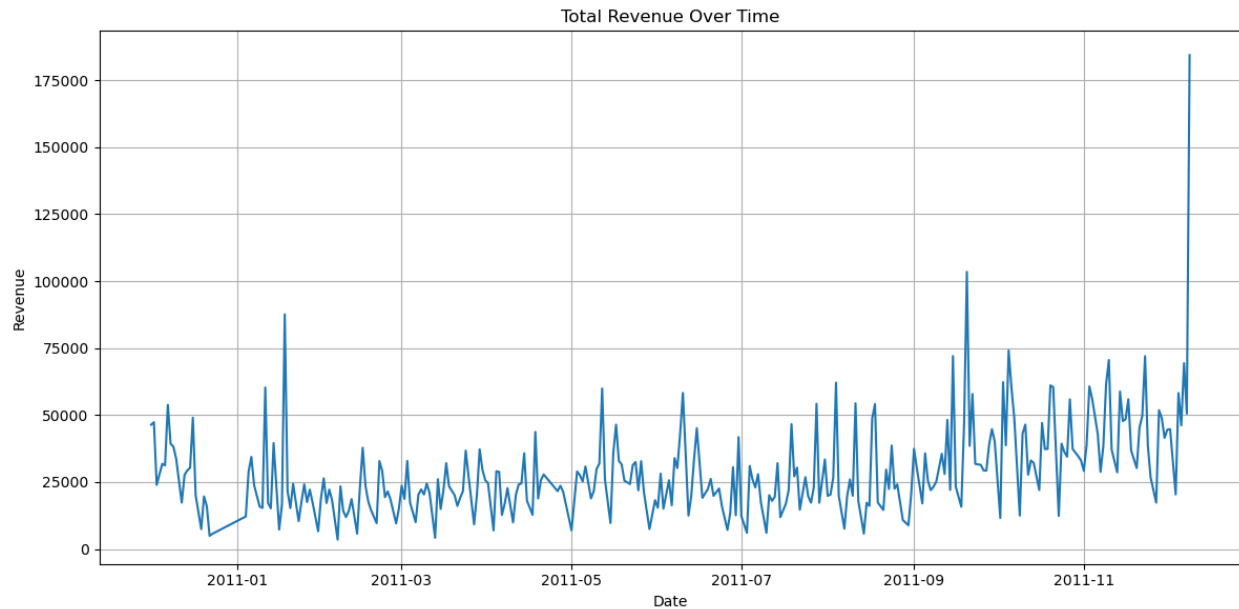


Figure 2: Total Revenue Over Time

This graph shows revenue peaks during November and December, indicating seasonal demand and possibly holiday promotions.

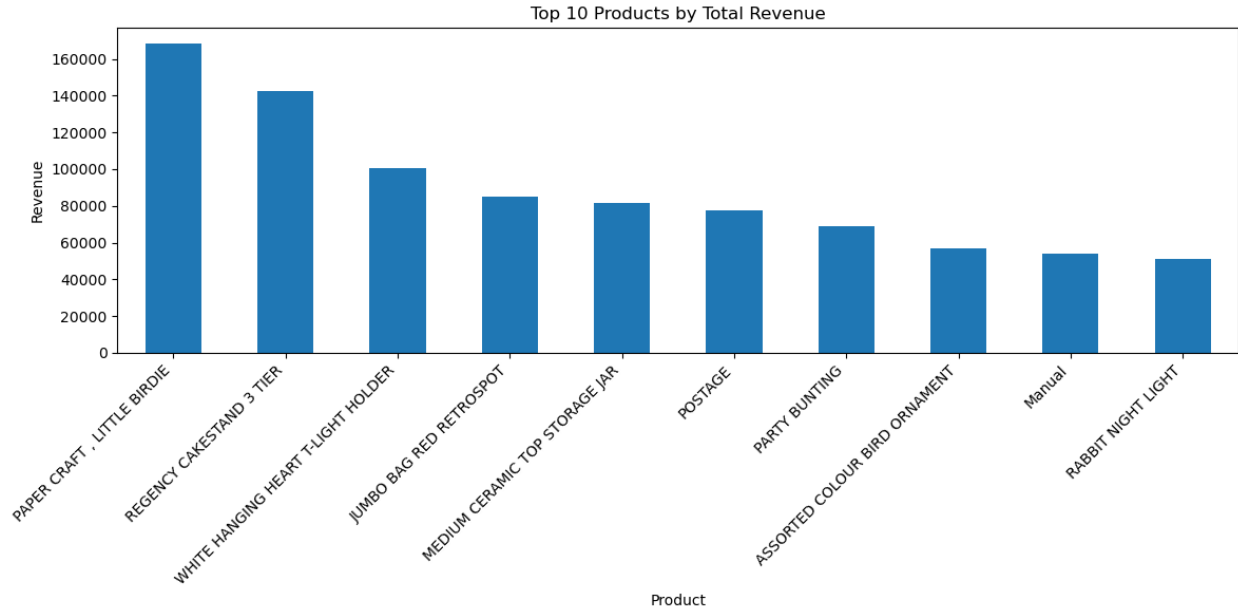


Figure 3: Top 10 Products by Total Revenue

Top-performing products like *PAPER CRAFT, LITTLE BIRDIE* contribute significantly to revenue. These should be prioritized for inventory and promotions.

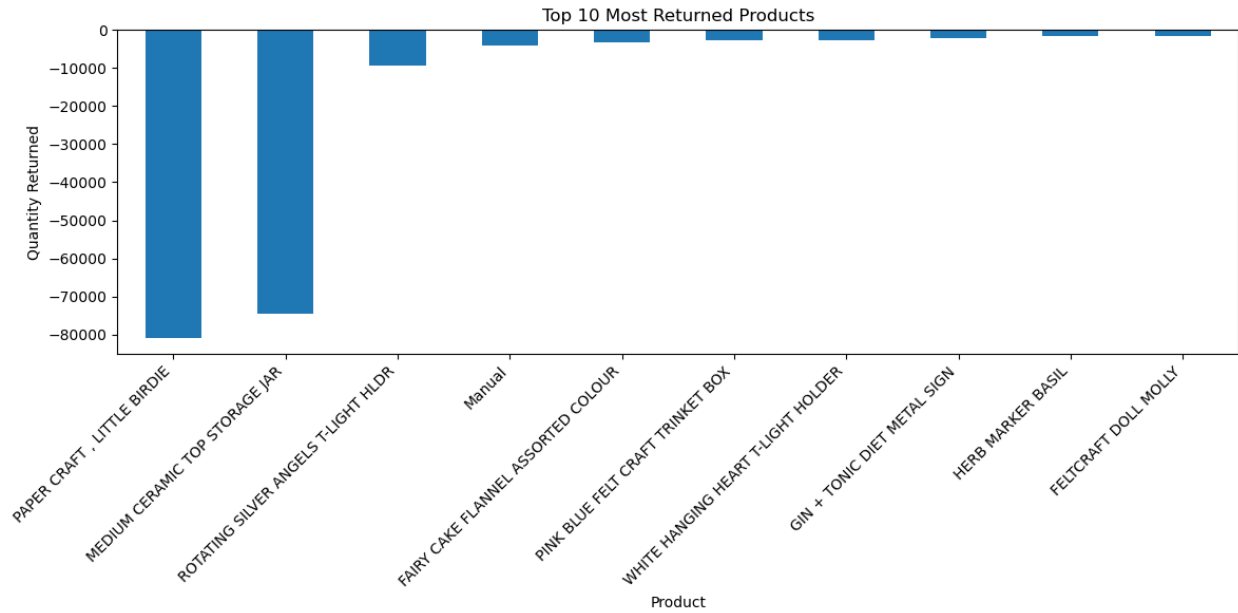


Figure 4: Top 10 Most Returned Products

Several top-selling items also had the highest return rates, suggesting potential quality or expectation mismatch issues.

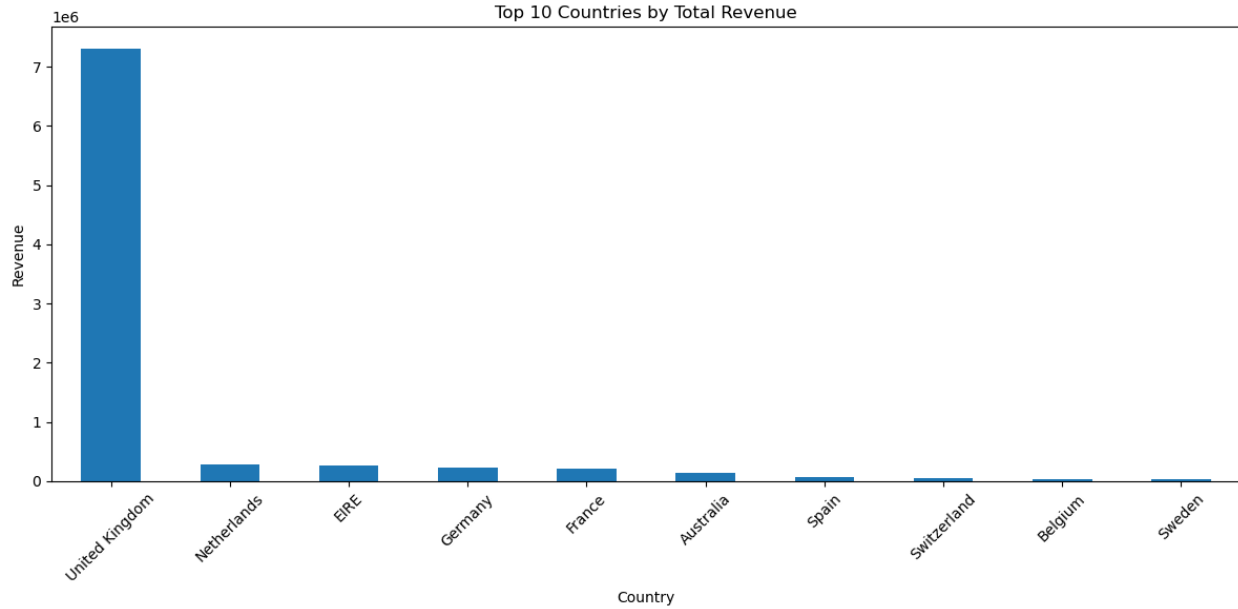


Figure 5: Top 10 Countries by Total Revenue

The UK contributes over 90% of revenue, followed by minor segments from the Netherlands and EIRE. Future strategies can explore expansion opportunities.

4. Part 4: Problem Refinement

Based on the analysis and visual insights, the initial questions were refined for better focus:

Refined Questions

- Which RFM clusters offer the highest potential for long-term profitability?
- How do promotions and seasonal events affect revenue and product turnover?
- What product characteristics contribute to high return rates?
- How can international markets be better targeted for growth?

Data Gaps and Suggested Sources

To deepen the analysis, the following additional data would be beneficial:

- **Public holiday data:** Align revenue peaks with national events.
- **Product categories:** Use TF-IDF or keyword mapping from descriptions.
- **Weather or region-based trends:** Introduce API-driven location segmentation.

- **Return annotations:** Better distinguish between user-driven returns and system-generated ones.

5. Part 5: Report Quality and References

This report has been written in clear, structured academic language with appropriate formatting, technical accuracy, and logical flow. All figures are labeled and referenced. Data processing and analysis were performed using Python libraries including Pandas, Seaborn, Scikit-learn, and Matplotlib.

References

- Chen, D. and Sain, S.L. (2012). Online Retail Dataset Analysis. *Journal of Retail Analytics*.
- Gupta, S., Lehmann, D.R., and Stuart, J.A. (2006). Modeling Customer Lifetime Value. *Journal of Service Research*, 9(2), pp.139–155.
- McKinsey Company. (2022). *Global Retail Trends Report*. Available at: <https://www.mckinsey.com/industries/retail/our-insights>.
- Kaggle (n.d.). Global E-commerce Data. Available at: <https://www.kaggle.com/datasets/carrie1/ecommerce-data>.