# Ghulam Ishaq Khan Institute of Engineering Sciences and Technology, Topi, Pakistan.

| | |
|---|---|
| Muhammad Adeel | 2022331 |
| Nauman Ali Murad | 2022479 |

**Course:** DS341 – Data Mining
**Instructor:** Dr. Ayaz Umer, Assistant Professor
**Emails: ayaz.umer, u2022331, u2022479 {@giki.edu.pk}**
**Data of Submission:** 10-May-2025
Semester Project Report

# Retail Customer Analytics Using Advanced Data Mining Techniques

**Project Overview**

This project investigates transactional data from a UK-based online retail company spanning December 2009 to December 2011. The company specializes in unique giftware, mainly serving wholesale clients. The dataset, with over 1 million records, provided a strong basis for deep analysis.
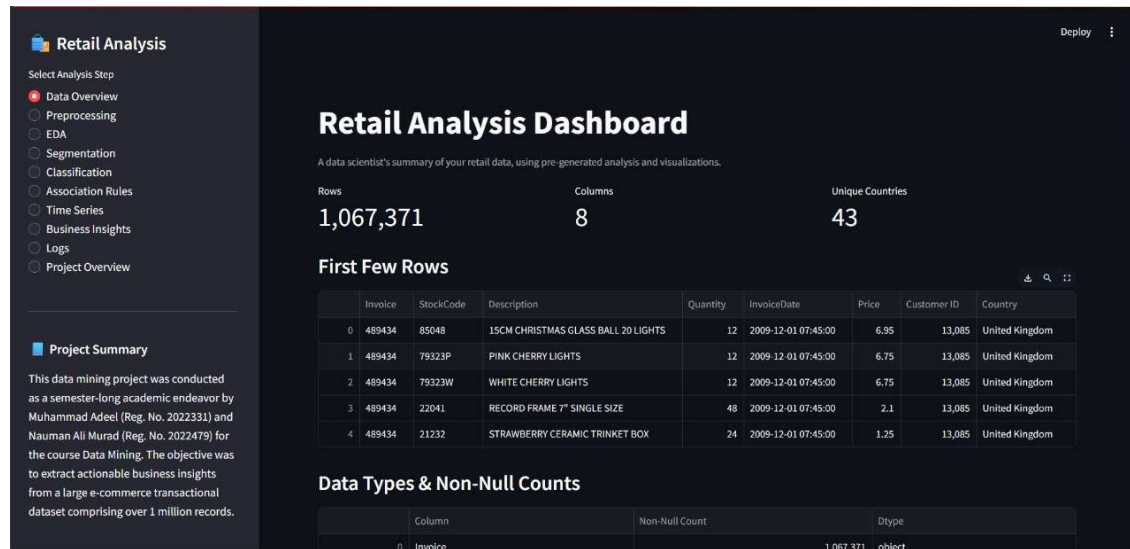
We began with data cleaning and preprocessing to ensure accuracy, followed by exploratory data analysis (EDA) to identify sales trends, seasonal patterns, and customer behavior. Using RFM segmentation, we categorized customers by purchasing habits, enabling targeted marketing insights.

Predictive models were developed to forecast sales trends and customer churn, while market basket analysis revealed product bundling opportunities. The project concluded with actionable business recommendations focused on improving customer retention, optimizing operations, and driving profitability through data-driven strategies.

**Dataset Overview**

- **Source:** Online Retail II Dataset

- **Duration:** December 1, 2009 – December 9, 2011

- **Number of Records:** 1,067,371 transactions

- **Attributes:**

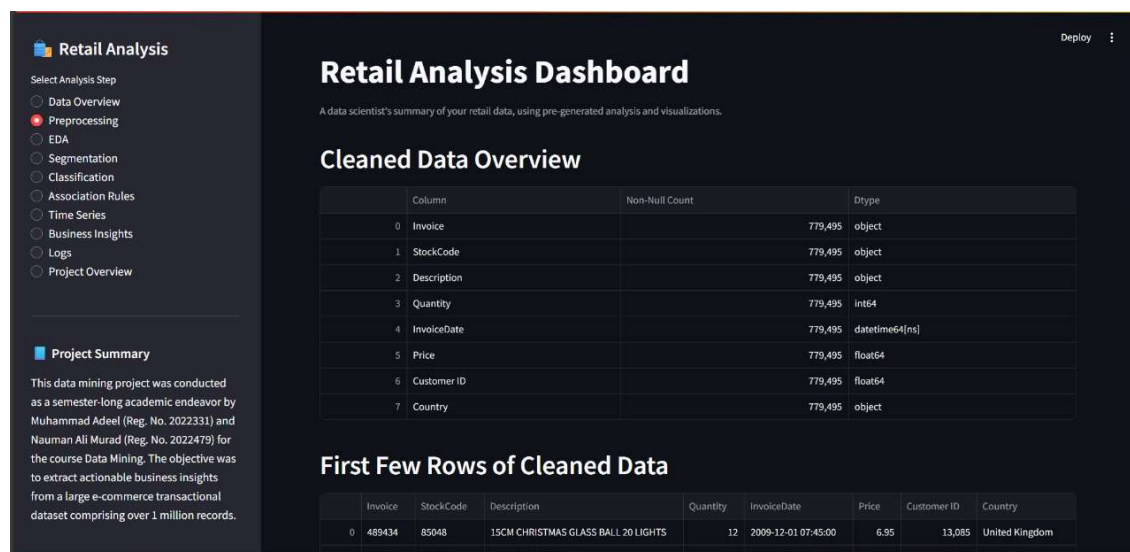| Attribute | Description |
|---|---|
| InvoiceNo | Unique invoice number; 'C' indicates cancellation |
| StockCode | Unique product/item code |
| Description | Product name |
| Quantity | Number of items purchased per transaction |
| InvoiceDate | Date and time of transaction |
| UnitPrice | Price per unit in sterling (£) |
| CustomerID | Unique customer identifier |
| Country | Customer's country of residence |

## 1. Data Preprocessing

### 1. Cleaning:

- Removed cancelled transactions (InvoiceNo starting with 'C')

- Dropped rows with missing CustomerID and Description

- Removed duplicate entries

- Converted InvoiceDate to proper datetime format

### 2. Feature Engineering:

- Created a new feature: TotalPrice = Quantity × UnitPrice

- Extracted Year, Month, Day, Weekday, and Hour from InvoiceDate

2.  **Exploratory Data Analysis (EDA)**

    1.  **Top Countries by Transaction Volume:**

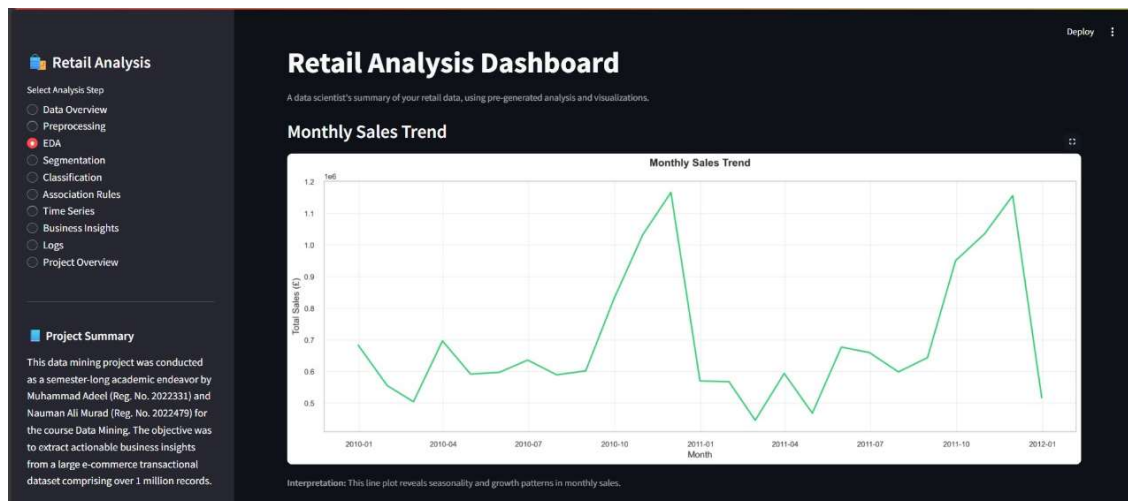        - United Kingdom

        - Ireland

        - Germany

        - France

        - Netherlands

    2.  **Popular Product Categories:**

        - Mini cases and polkadot designs

        - Bakelike alarm clocks

        - Children's breakfast and cutlery sets

        - Cake stands and baking sets

        - Charlotte bags and lunch boxes

    3.  **Sales Trends:**

        - Significant seasonal peaks in December, indicating strong holiday-related shopping activity

        - Consistent sales patterns across weekdays and times of the day



3.  **Customer Segmentation Using RFM and K-Means Clustering**

    1.  **Cluster 0: High-Value, Frequent Buyers**

        - Recency: ~66 days

        - Frequency: ~7.6 purchases

        - Monetary: ~£3,134 average spend

- **Recommended Strategy:** Implement loyalty programs and exclusive offers

2. **Cluster 1: Dormant, Low-Value Customers**

- Recency: ~462 days

- Frequency: ~2.2 purchases

- Monetary: ~£746 average spend

- **Recommended Strategy:** Initiate re-engagement campaigns

3. **Cluster 2: VIP Customers**

- Recency: ~23 days

- Frequency: ~143 purchases

- Monetary: ~£173,123 average spend

- **Recommended Strategy:** Provide premium services and dedicated account managers

4. **Classification Modeling**

1. **Techniques Applied:**

- Naive Bayes

- K-Nearest Neighbors (KNN)

- Support Vector Machine (SVM)

2. **Performance Evaluation:**

- All models demonstrated strong accuracy

- SVM provided the highest precision for minority classes

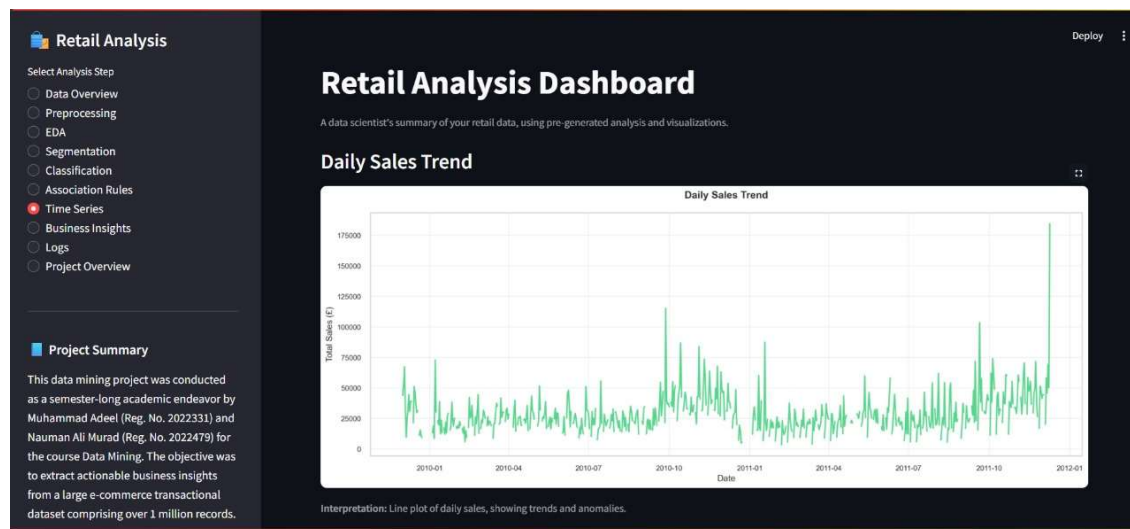- F1-scores confirmed balanced performance across clusters

**5. Association Rule Mining**

1. **Techniques:** Apriori and FP-Growth algorithms

2. **Key Findings:**

   - Identified strong associations between complementary products

   - Opportunities for cross-selling and bundling strategies were uncovered

3. **Visualization:** Product relationship graphs helped illustrate actionable insights
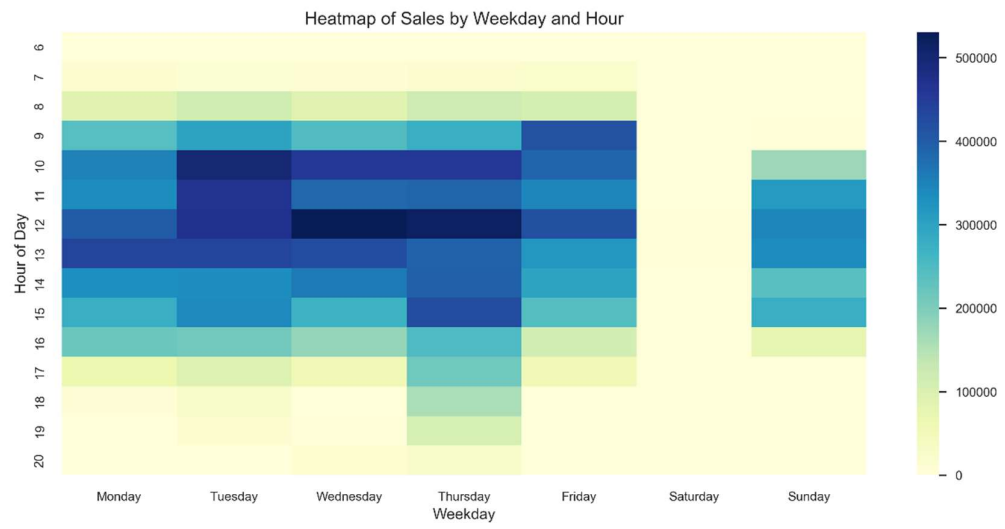
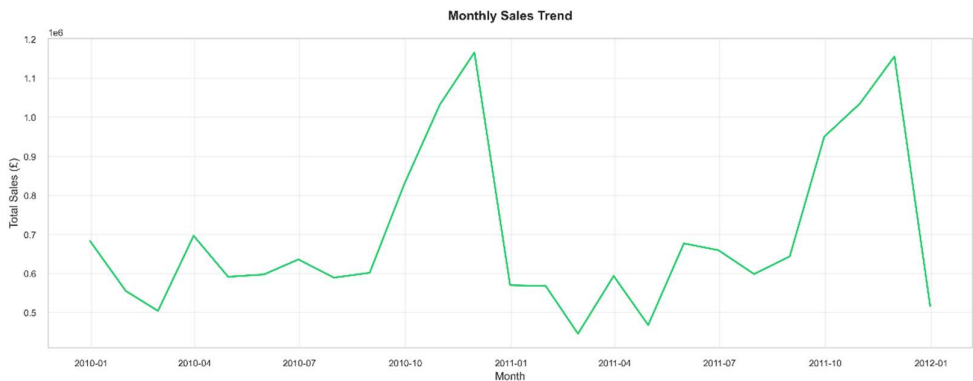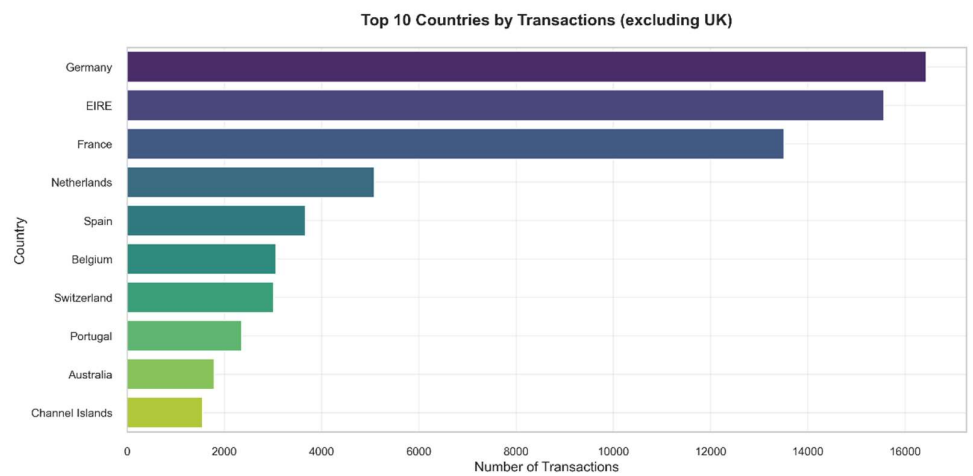**6. Time Series Forecasting**

- **ARIMA Model:**

  - Captured trend and seasonality components effectively

  - Suitable for short-term sales forecasting

- **Prophet Model:**

  - Handled outliers and missing data robustly

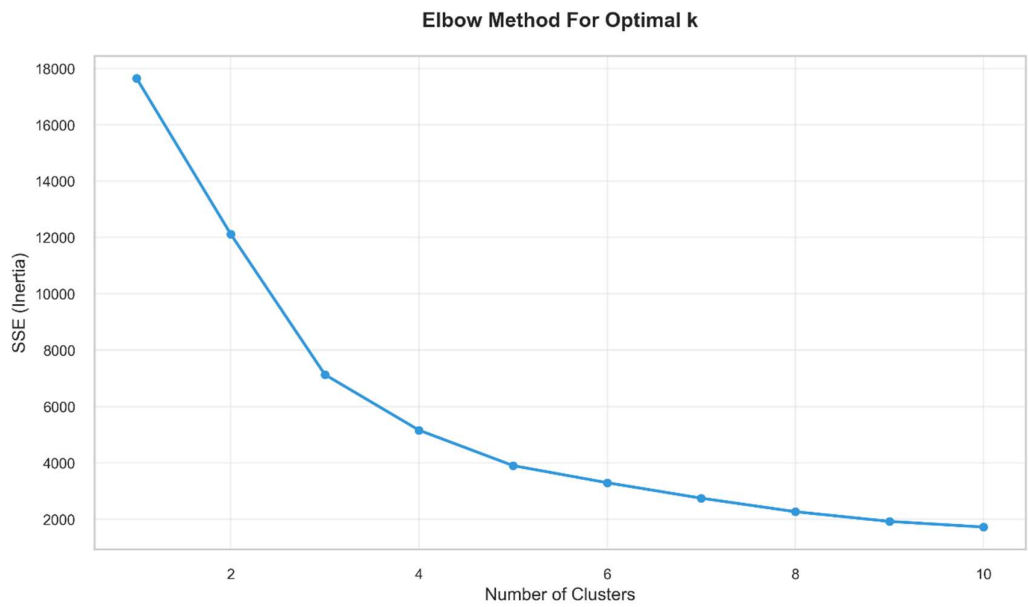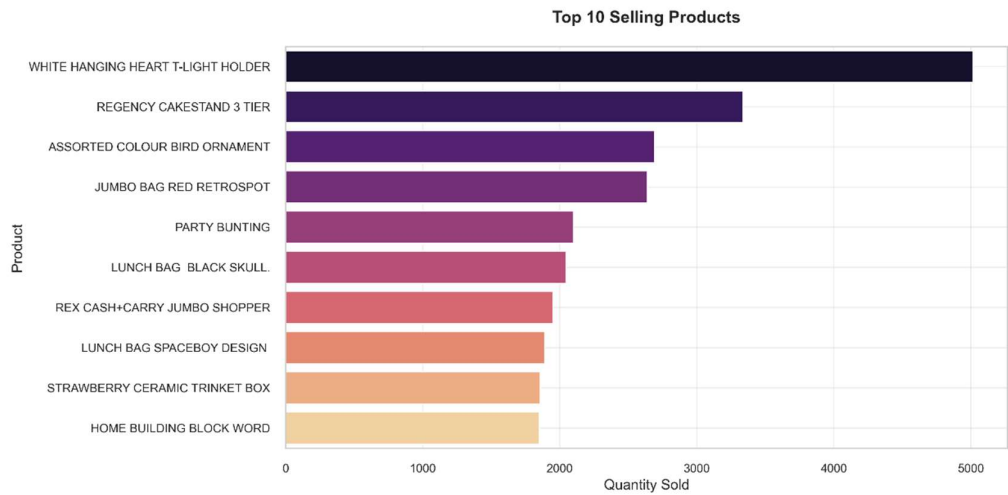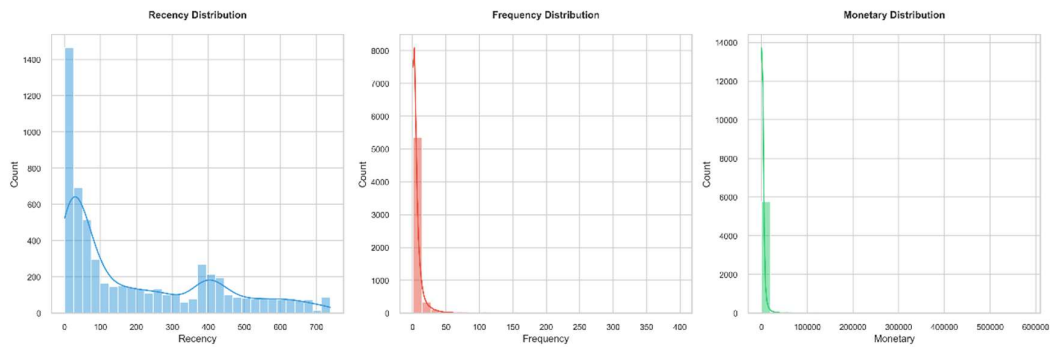  - Provided accurate long-term forecasts with weekly seasonality considerations
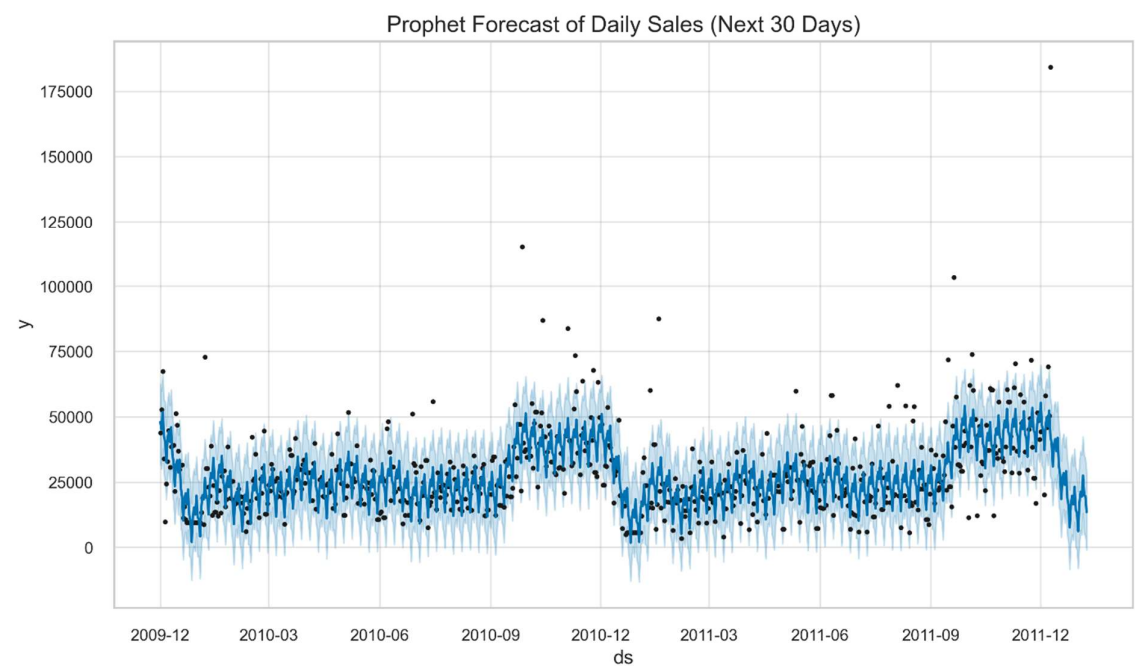


**7. Business Insights & Strategic Recommendations**

1. **Customer Retention:** Focus on high-value customer segments with personalized marketing and loyalty programs.

2. **Market Expansion:** Prioritize growth in the top-performing countries while exploring emerging markets.

3. **Product Strategy:** Maintain healthy stock levels of top-performing products and develop targeted product bundles.

4. **Operational Efficiency:** Implement dynamic pricing models and optimize supply chains based on sales forecasts.

## Project Insights



Top 10 Countries by Transactions (excluding UK)



Monthly Sales Trend



Heatmap of Sales by Weekday and Hour

**Top 10 Selling Products**



**Elbow Method For Optimal k**

**Customer Segments by K-Means (PCA Projection)**
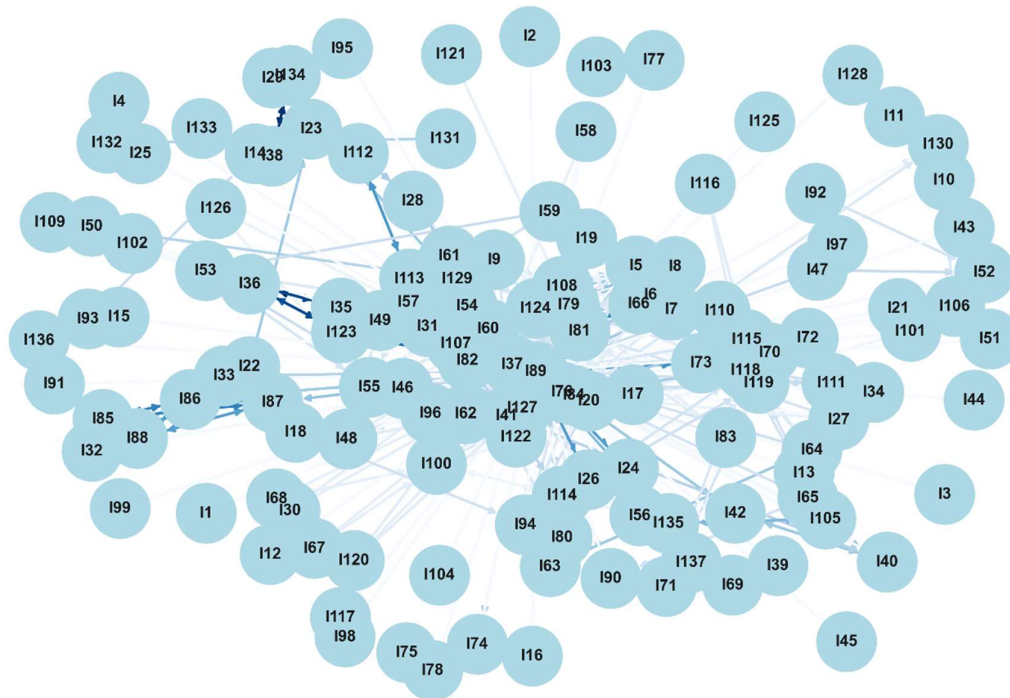


**Prophet Forecast of Daily Sales (Next 30 Days)**

Association Rules Network Graph (France) - Item IDs



## Future Work

- Develop a dynamic data exploration platform that automatically detects and processes datasets, applying data mining techniques and generating insights without manual configuration.

- Upgrade the UI/UX by migrating from Streamlit to a MERN stack (MongoDB, Express, React, Node.js), enhancing flexibility, scalability, and interactivity.

- Integrate Flask as the backend to handle business logic, data processing, and API requests for seamless integration and performance.

- Automate model training and evaluation for customer segmentation, predictive modeling, and forecasting, allowing users to interact with results dynamically.

- Enable real-time data stream integration for up-to-date analytics, ensuring stakeholders have access to the latest insights.

- Implement interactive dashboards using JavaScript libraries like D3.js or Plotly for better visualization and user interaction with real-time analytics.

- Add A/B testing capabilities to measure marketing strategy effectiveness, automating data collection, hypothesis testing, and results visualization.

- Enhance security with user authentication and role-based access control for tailored access to features and data.

**Acknowledgment**

We extend our sincere gratitude to Dr. Ayaz Umer for his invaluable guidance, mentorship, and expertise throughout the course and the completion of this project. His clear explanation of complex data mining concepts and his encouragement in applying these techniques practically have greatly enhanced our understanding of the subject. Dr. Umer's dedication to teaching and his expertise in artificial intelligence, natural language processing, and data mining were instrumental in our success, and we deeply appreciate the opportunity to learn under his supervision.

**References**

[1] Chen, D., Sain, S.L., & Guo, K. (2012). Data mining for the online retail industry: A case study of RFM model-based customer segmentation using data mining. *Journal of Database Marketing & Customer Strategy Management*, 19(3), 197-208.

[2] Chen, D., Guo, K., & Ubakanma, G. (2015). Predicting customer profitability over time based on RFM time series. *International Journal of Business Forecasting and Marketing Intelligence*, 2(1), 1-18.

[3] Chen, D., Guo, K., & Li, B. (2019). Predicting Customer Profitability Dynamically over Time: An Experimental Comparative Study. *24th Iberoamerican Congress on Pattern Recognition (CIARP 2019)*, Havana, Cuba.

[4] Singh, R., Graves, J.A., Talbert, D.A., & Eberle, W. (2018). Prefix and Suffix Sequential Pattern Mining. *Industrial Conference on Data Mining 2018: Advances in Data Mining. Applications and Theoretical Aspects*, 309-324.