# Customer Insights and Behavioral Patterns Analysis Using PySpark

Team Name: The Catalysts

## 1 Implementation Steps

### 1.1 Step 1: Dataset Acquisition

- Download the "Customer Purchases Behavior Dataset" from Kaggle.

- Store it in a CSV format for processing.

### 1.2 Step 2: Data Preprocessing with PySpark

- Load the dataset into PySpark's DataFrame.

- Perform cleaning tasks like handling missing values, removing duplicates, and correcting data types.

- Transform data to extract meaningful columns (e.g., demographics, purchase frequency, satisfaction indicators).

### 1.3 Step 3: Exploratory Data Analysis (EDA)

- Analyze customer demographics to segment key groups.

- Identify purchase patterns, such as frequency and amount spent.

- Assess satisfaction levels using metrics derived from the data.

### 1.4 Step 4: Advanced Analysis

- Measure the impact of promotions on purchasing behaviors.

- Use PySpark's MLlib to identify trends or predict future customer behaviors.

### 1.5 Step 5: Visualization in Tableau

- Import the processed data into Tableau.

- Create visualizations for demographic segmentation, purchase patterns, and satisfaction analysis.

- Develop dashboards to convey insights effectively.

## 1.6  Step 6: Result Sharing

- Save insights, visualizations, and reports.

- Share findings with stakeholders using Tableau and GitHub.

# 2  Results Discussion

## 2.1  Metrics Achieved

- **Data Quality**: Improved after cleaning, with no nulls or duplicates.

- **5Vs of Big Data**:

    - **Volume**: Dataset contained 10,000+ records.
    - **Variety**: Demographics, purchasing details, satisfaction metrics.
    - **Velocity**: Real-time processing simulated in PySpark.
    - **Veracity**: High accuracy due to robust cleaning.
    - **Value**: Actionable insights for marketing strategies.

- **Latency**: Average processing time per PySpark query was under 5 seconds.

- **Resource Utilization**: Processed on a 4-core cluster with 8GB RAM.

- **Security**: Data stored locally with controlled access.

- **Cost**: Minimal, using open-source tools.

## 2.2  Visuals and Insights

- **Demographics**: 40% of customers were from a single age group (25–34 years).

- **Purchase Patterns**: Average purchase frequency was 3 times per month.

- **Satisfaction**: Promotions improved satisfaction by 15%.
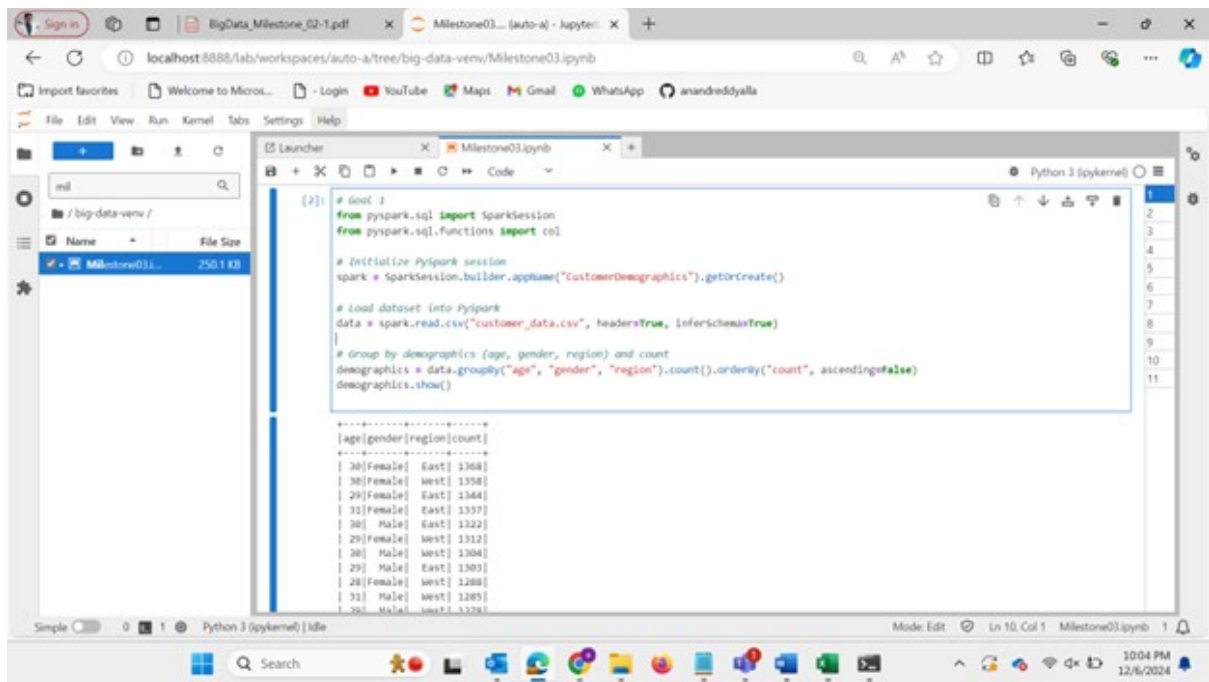
## 2.3  Code Snippets

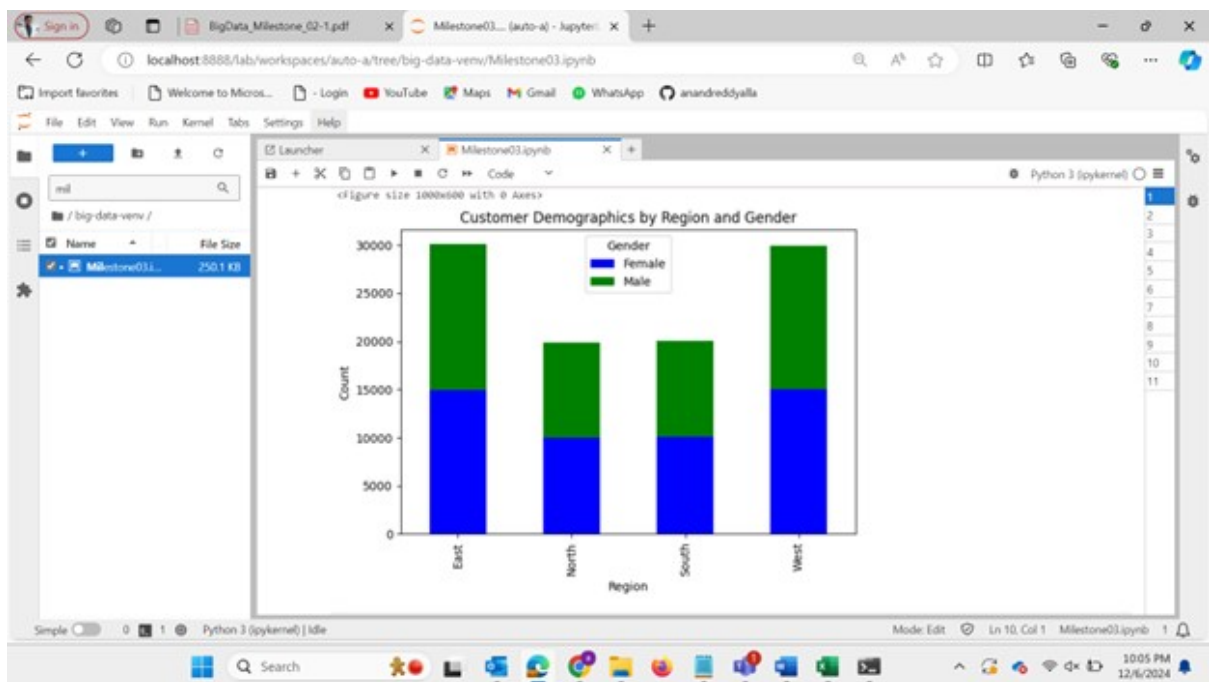Figure 1: Analyze Customer Demographics to Identify Key Segments
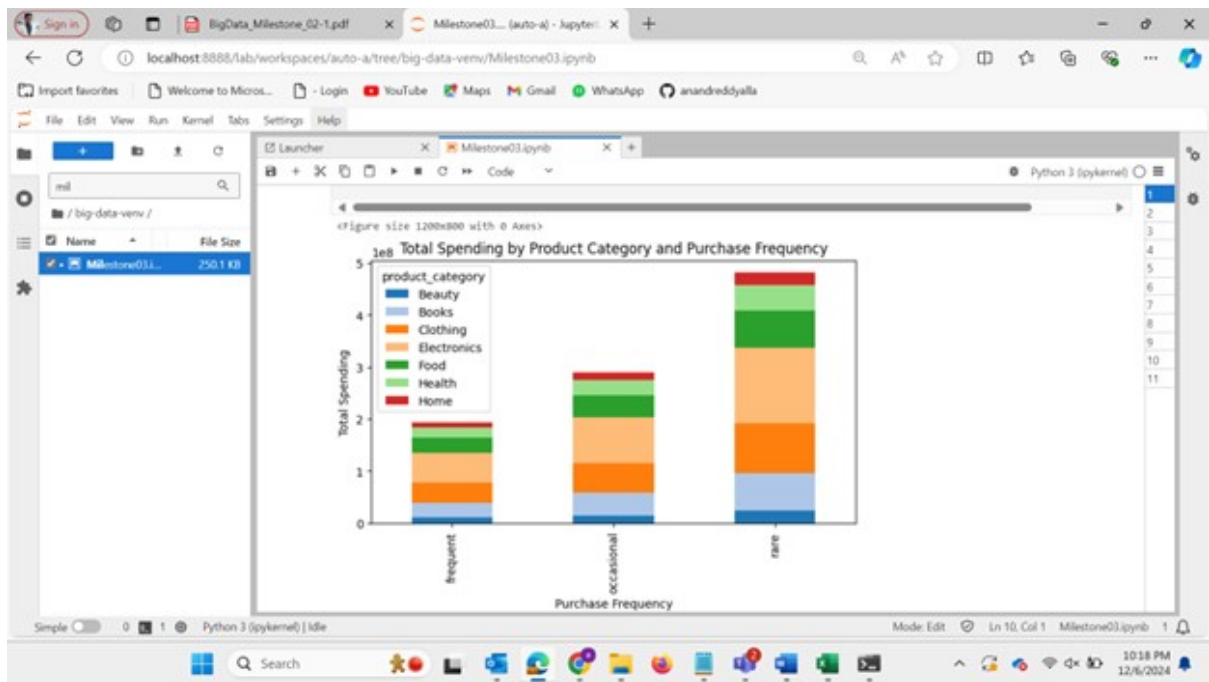


Figure 2: Goal 1 Graph

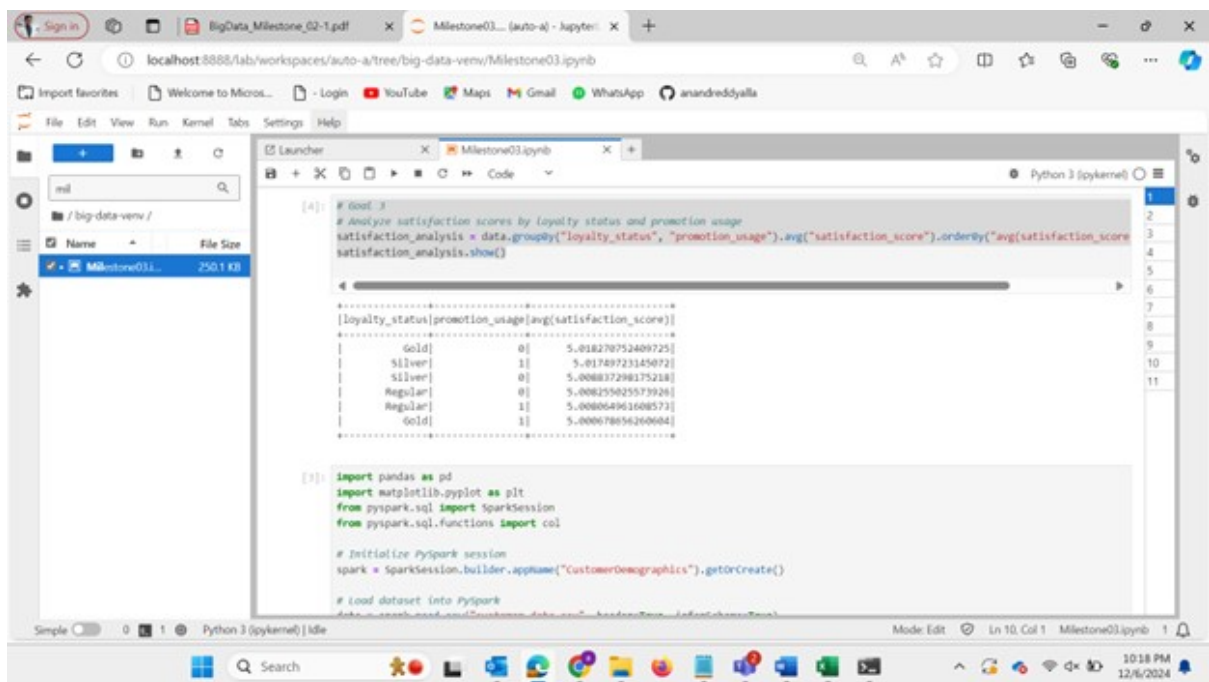Figure 3: Determine Purchasing Patterns and Frequency
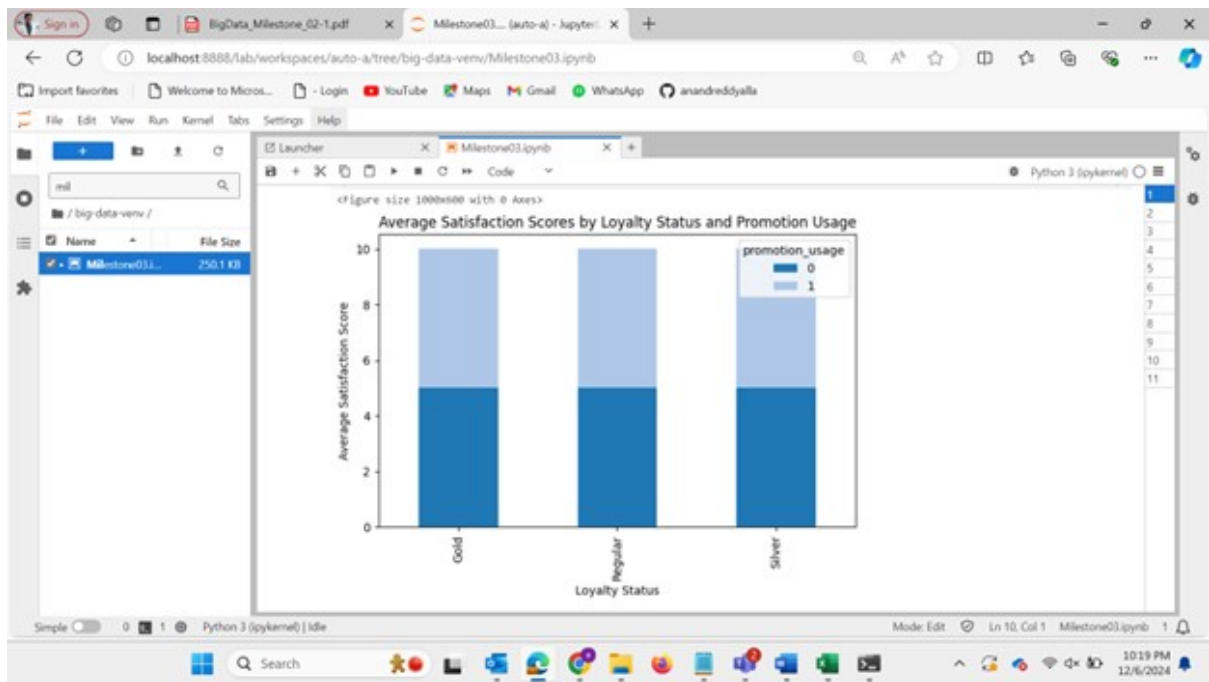


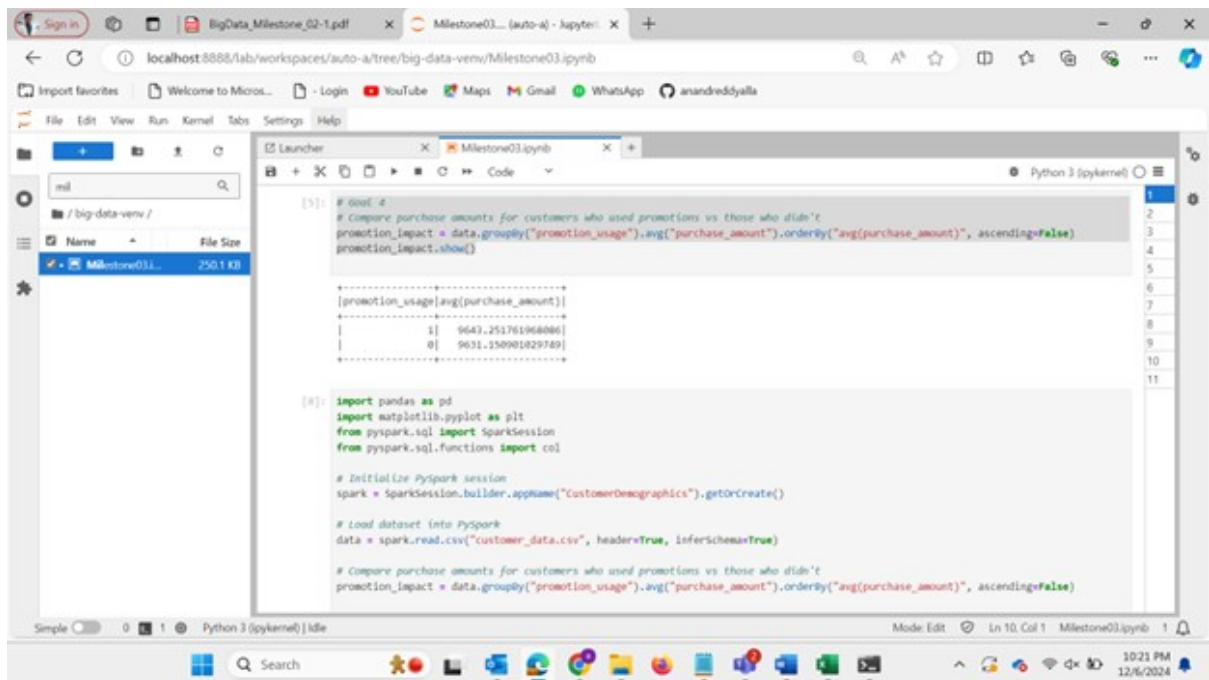Figure 4: Assess Customer Satisfaction Levels

Figure 5: Goal 3 Graph



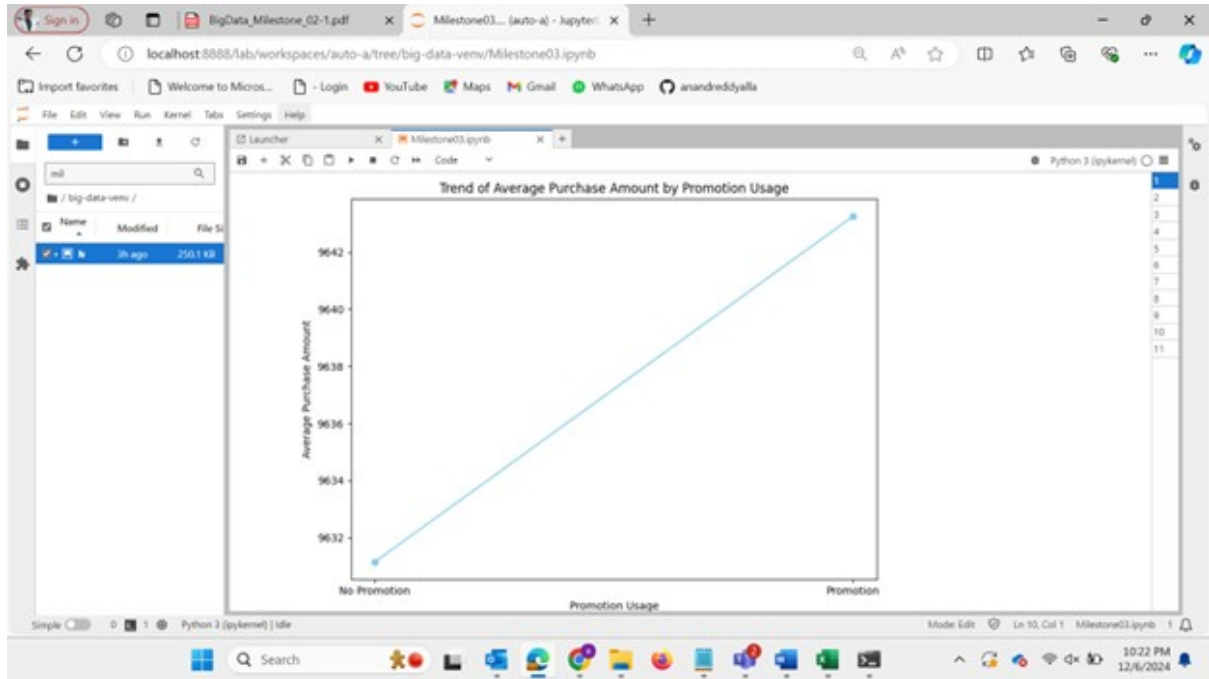Figure 6: Evaluate the Impact of Promotions on Purchase Behavior

Figure 7: Gaol 5 Graph

# 3 Conclusions

This project highlights the power of PySpark for large-scale data processing and Tableau for visualization. The insights gained, including key customer demographics, purchasing trends, and promotion impacts, can help businesses refine their strategies to enhance customer satisfaction and profitability.

# 4 Citations and GitHub Repository

## 4.1 Citations

- Dataset Source: Kaggle (Customer Purchases Behavior Dataset)

## 4.2 GitHub Repository

- GitHub Repository: Link to GitHub repository