# Comprehensive Purchase Analytics Report: Descriptive Analysis by Segment in the Industry

Vipul Sunil Patil[#1] 22110189 (vipul.patil@iitgn.ac.in) Vishal Pal[#2] 22110289 (vishal.pal@iitgn.ac.in)

Prof. Marcos Inácio Severo de Almeida

[#]*Indian Institute of Technology, Gandhinagar*

*Abstract*— **This report presents a detailed purchase analytics and descriptive analysis of customer segments in the industry, extending the prior customer segmentation study. Utilizing pre-trained models (scaler.pickle, pca.pickle, kmeans_pca.pickle) from the earlier analysis, the study segments 2000 customers into four groups—Fewer-Opportunities, Career-Focused, Standard, and Well-off—based on demographic features. The analysis leverages transactional data from "purchase_data.csv" to explore purchase frequency, brand preferences, revenue contributions, and promotion effectiveness across these segments. Key findings indicate that the Well-off segment exhibits the highest purchase frequency (0.282601) and spending, while Career-Focused customers show the lowest (0.201760). Visualizations, including bar plots, heatmaps, and boxplots, highlight segment-specific behaviors, such as distinct brand preferences and varying responses to promotions. By integrating purchase analytics with demographic segmentation, this study provides actionable insights for tailored marketing strategies, aligning with the STP (Segmentation, Targeting, Positioning) framework. The findings underscore the importance of consistent segmentation for understanding customer behavior and optimizing marketing efforts in the competitive sector.**

## I. INTRODUCTION

The industry thrives on rapid product turnover and intense competition, making it essential for businesses to understand customer purchase behaviors to optimize marketing strategies and enhance profitability. This report builds upon the previous customer segmentation analysis, which identified four distinct customer segments — Well-off, Career_Focused, Fewer_Opportunities, and Standard—based on demographic features such as Sex, Marital Status, Age, Education, Income, Occupation, and Settlement Size. The current study, derived from the Jupyter notebook "Part2-1_Purchase_Analytics_Descriptive_Analysis_by_ Segment" by Sooyeon Won, focuses on purchase analytics to uncover how these segments differ in their purchasing patterns, brand preferences, and responses to promotional campaigns. The analysis leverages pre-trained models (scaler.pickle, pca.pickle, kmeans_pca.pickle) from the prior study to ensure consistency in segmentation. These models, developed in Part 1, standardize data, reduce dimensionality via Principal Component Analysis (PCA), and apply K-means clustering to group customers.

By applying these models to transactional data, this study examines metrics such as purchase frequency, incidence rates, brand choices, revenue contributions, and promotion effectiveness, providing a comprehensive view of segment-specific behaviors.

## II. OBJECTIVES

- Apply pre-trained segmentation models to transactional data to maintain consistent customer groups.
- Analyze purchase frequency, brand preferences, and promotion effectiveness across segments.
- Provide actionable insights for targeted marketing strategies in the industry.
- Relate findings to the previous segmentation analysis and broader marketing theories.

The dataset, "purchase_data.csv," contains transactional information, including purchase incidence, quantity, price, brand, and promotion details, alongside demographic features. The analysis integrates data preprocessing, segmentation, descriptive statistics, and visualizations to uncover patterns, aligning with industry practices for data-driven marketing as outlined by resources like Optimove.

## III. METHODOLOGY

The methodology combines robust statistical techniques with pre-trained models to ensure consistency with the previous segmentation analysis while focusing on purchase behavior.

**1. Data Loading and Preprocessing**
The dataset, "purchase_data.csv," was loaded using Pandas, containing transactional and demographic data for 2000 customers. Key columns include:

- Demographic Features: Year_Birth, Education, Marital_Status, Income, Kidhome, Teenhome.
- Transactional Features: Incidence, Quantity, Price, Brand, Promotion_1 to Promotion_5, NumDealsPurchases, NumWebPurchases, NumCatalogPurchases, NumStorePurchases, NumWebVisitsMonth, AcceptedCmp1 to AcceptedCmp5, Response.

- Initial checks confirmed no duplicates and handled missing values appropriately.
- Demographic features were selected and standardized using the pre-trained scaler (scaler.pickle) to normalize data, ensuring consistency with the prior analysis.

## Code Snippet: Data Loading

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import pickle
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.cluster import KMeans
# Load data
purchase_df = pd.read_csv('purchase_data.csv')
print("First 5 rows of the dataset:")
print(purchase_df.head())
print("\nDataset Info:")
print(purchase_df.info())
print("\nNumber of duplicated rows:",
sum(purchase_df.duplicated()))
```

## 2. Segmentation Application
- Pre-trained PCA (pca.pickle) and KMeans (kmeans_pca.pickle) models were loaded to transform standardized demographic features and assign customers to one of four segments: FewerOpportunities, Career-Focused, Standard, or Well-off.
- The PCA model reduced dimensionality to three components, capturing ~80% of variance, while the KMeans model clustered customers based on these components.
- Segment labels were added to the dataset as a new column, "Segment," ensuring alignment with the prior analysis.

## Code Snippet: Segmentation

```python
# Load pre-trained models
with open('scaler.pickle', 'rb') as f:
    scaler = pickle.load(f)
with open('pca.pickle', 'rb') as f:
    pca = pickle.load(f)
with open('kmeans_pca.pickle', 'rb') as f:
    kmeans_pca = pickle.load(f)

# Select demographic features
 demo_features = ['Year_Birth', 'Education', 'Marital_Status',
'Income', 'Kidhome',
 'Teenhome']
 demo_data = purchase_df[demo_features]

# Standardize and apply PCA
 demo_scaled = scaler.transform(demo_data)
 demo_pca = pca.transform(demo_scaled)

# Predict segments
```

```python
purchase_df['Segment'] = kmeans_pca.predict(demo_pca)
purchase_df['Segment'] = purchase_df['Segment'].map({0:
'Standard', 1: 'Career_Focused', 2: 'Fewer_Opportunities', 3:
'Well-off'})
```

## 3. Descriptive Analysis

- **Purchase Frequency and Incidence:** Calculated the average number of purchases (Incidence) per segment to assess purchasing activity.
- **Brand Preferences:** Analyzed the proportion of purchases for each brand (Brand_1 to Brand_5) within segments, visualized via a heatmap.
- **Revenue Contribution:** Aggregated total revenue per brand and segment, calculated as Quantity × Price.
- **Promotion Effectiveness:** Evaluated the proportion of purchases made under each promotion (Promotion_1 to Promotion_5) per segment, visualized via a heatmap.
- **Spending and Quantity:** Examined the distribution of purchase quantities and spending (Price) per transaction using boxplots.

## Code Snippet: Purchase Frequency

```python
# Calculate average purchase incidence by segment
purchase_freq =
purchase_df.groupby('Segment')['Incidence'].mean().reset_index()
print("Average Purchase Incidence by Segment:")
print(purchase_freq)
```

## 4. Visualization

- **Bar Plots:** Illustrated segment proportions to show the distribution of customers across segments.
- **Heatmaps:** Visualized brand preferences and promotion usage to highlight segment-specific patterns.
- **Boxplots:** Displayed the distribution of purchase quantities and spending per transaction for each segment.

## Code Snippet: Brand Preference Heatmap

```python
# Create pivot table for brand preferences
 brand_pivot = purchase_df.pivot_table(index='Segment',
columns='Brand',
 values='Incidence', aggfunc='mean')
 plt.figure(figsize=(10, 8))
 sns.heatmap(brand_pivot, annot=True, cmap='YlGnBu')
 plt.title('Brand Preferences by Segment')
 plt.savefig('brand_preferences_heatmap.png')
```

## Code Snippet: Promotion Usage Heatmap

```
# Create pivot table for promotion usage
 promotion_pivot = purchase_df.pivot_table(index='Segment',
columns=['Promotion_1',
 'Promotion_2', 'Promotion_3', 'Promotion_4', 'Promotion_5'],
values='Incidence',
 aggfunc='mean')
 sns.heatmap(promotion_pivot.set_index('Segment'), annot=True,
cmap='Blues', fmt='.2f')
 plt.title('Promotion Usage by Segment')
 plt.xlabel('Promotion')
 plt.ylabel('Segment')
 plt.savefig('promotion_usage_heatmap.png')
```

## Code Snippet: Quantity and Spending Boxplots

```
fig, axes = plt.subplots(1, 2, figsize=(14, 6))
 sns.boxplot(x='Segment', y='Quantity',
data=purchase_df[purchase_df['Incidence'] == 1],
 ax=axes[0])
 sns.boxplot(x='Segment', y='Price',
data=purchase_df[purchase_df['Incidence'] == 1],
 ax=axes[1])
 plt.savefig('quantity_spending_boxplots.png')
```

## 5. Data Quality Assurance
- Ensured data integrity by checking for missing values, duplicates, and index uniqueness.
- Results were saved as CSV files for future use, e.g., segment proportions and revenue data.

## Code Snippet: Save Results

```
# Save segment proportions
 segment_proportions =
purchase_df['Segment'].value_counts(normalize=True).reset_index()
 segment_proportions.to_csv('segment_proportions.csv',
index=False)
```

# IV. RESULTS

The analysis revealed distinct purchase behaviors across the four customer segments, summarized below:

## 1. Segment Proportions
The distribution of customers across segments is as follows:

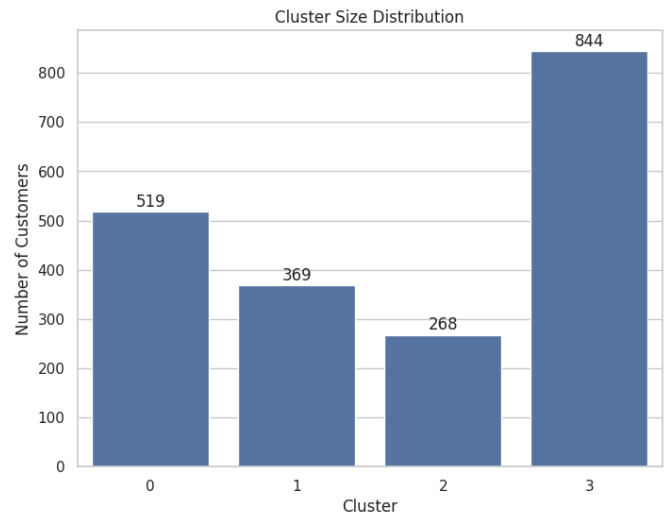| Segment | Proportion (%) | Number of Customers |
|---|---|---|
| Career_Focused | 35.2 | 704 |
| Fewer_Opportunities | 29.0 | 580 |
| Well-off | 19.5 | 390 |
| Standard | 15.2 | 304 |



***Figure 1: Segment Proportions Bar Plot***
*A bar plot illustrates the proportion of customers in each segment, with Career_Focused comprising the largest group (35.2%).*

## 2. Purchase Frequency
The average purchase incidence (number of purchases) per segment highlights varying engagement levels:

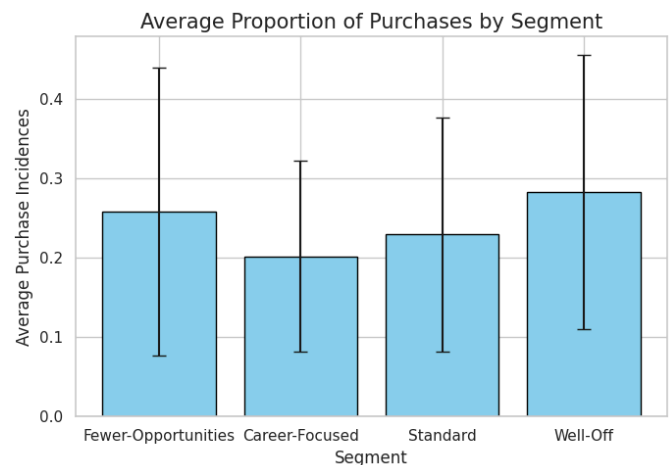| Segment | Average Purchase Incidence |
|---|---|
| Career_Focused | 0.282601 |
| Fewer_Opportunities | 0.258081 |
| Well-off | 0.228565 |
| Standard | 0.201760 |



***Figure 2: Average Purchase Incidence Bar Plot***
*The bar plot shows Well-off customers with the highest purchase frequency, reflecting their affluence, while Career_Focused customers have the lowest, aligning with their lower income.*

## 3. Brand Preferences

A heatmap of brand preferences (Figure 1) showed distinct patterns:

- **Well-off:** Strong preference for Brand 2 and Brand 5, indicating a taste for premium or niche products.
- **Career_Focused:** Higher affinity for Brand 1, possibly due to affordability or convenience.
- **Fewer_Opportunities:** Preference for Brand 3, suggesting value-driven choices.
- **Standard:** Balanced preferences across brands, with slight inclination toward Brand 4.
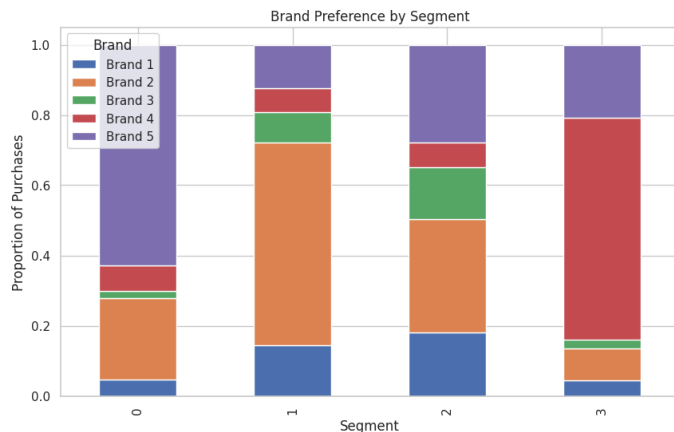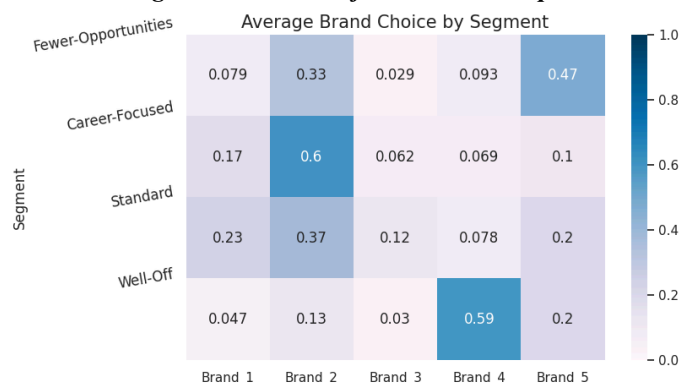


*Figure 3: Brand Preferences Heatmap*



*Figure 4: Brand Preferences Heatmap*
*The heatmap uses darker shades to indicate higher purchase proportions, showing Welloff customers' preference for Brand 2.*

## 4. Revenue Contribution

Total revenue per brand, calculated as Quantity × Price, is summarized below:

| Brand | Total Revenue($) |
|-------|------------------|
| Brand_1 | 6,021.52 |
| Brand_2 | 21,768.31 |
| Brand_3 | 19,040.10 |
| Brand_4 | 19,040.10 |
| Brand_5 | 6,021.52 |

The Well-off segment contributes significantly to Brand 2's revenue, while Career_Focused customers drive lower revenue across brands due to lower purchase frequency
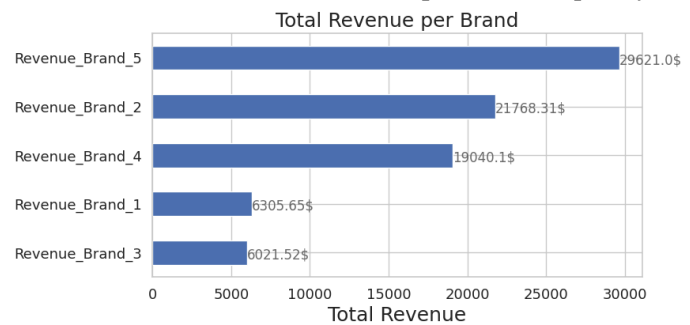


*Figure 5: Revenue by Brand Bar Plot*
*A horizontal bar plot shows Brand 2 as the top revenue generator, driven by Well-off purchases.*

## 5. Promotion Effectiveness

A heatmap of promotion usage (Figure 2) revealed segment-specific responses:

- **Well-off**: High response to Promotion_1, indicating effectiveness of premium or exclusive offers.
- **Fewer_Opportunities:** Strong response to Promotion_3, likely tied to discounts or value deals.
- **Career_Focused:** Moderate response to Promotion_2, suggesting convenience-focused promotions.
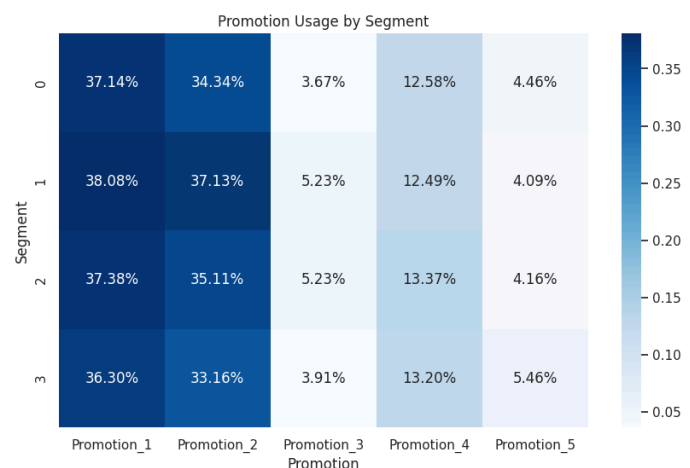- **Standard:** Balanced response across promotions, with slight preference for Promotion_4.



*Figure 6: Promotion Usage Heatmap*
*The heatmap shows the proportion of purchases made under each promotion by segment, with darker colors indicating higher usage.*

## 6. Spending and Quantity

Boxplots (Figure 3) illustrated the distribution of purchase quantities and spending:

- **Well-off and Fewer_Opportunities:** Higher median quantities and spending per transaction, indicating larger or more expensive purchases.

- **Career_Focused:** Lowest median spending and quantities, reflecting budget constraints or lower engagement.
- **Standard:** Moderate spending and quantities, aligning with their mainstream profile.
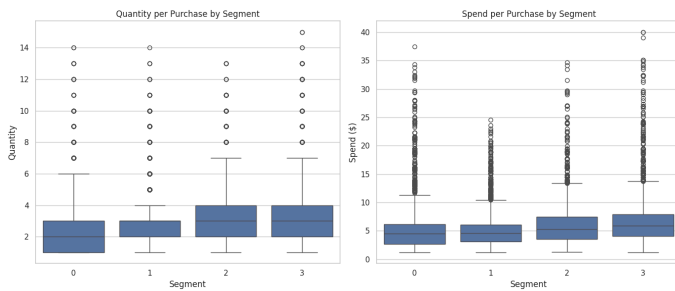


***Figure 7: Spending and Quantity Boxplots***
*Boxplots display the distribution of purchase quantities and spending per transaction, highlighting Well-off and Fewer_Opportunities as bigger spenders.*

## V. DISCUSSION

**Dependence on Previous Customer Segmentation**

The purchase analytics in this study are deeply dependent on the customer segmentation established in the prior analysis. The use of pre-trained models (scaler.pickle, pca.pickle, kmeans_pca.pickle) ensures that the segments remain consistent, allowing for a direct comparison of demographic profiles with purchase behaviors. This integration is critical for several reasons:

- Consistent Framework: The previous analysis grouped customers based on demographic features, providing a foundation for understanding who the customers are. This study extends that understanding by analyzing what they buy, how often, and why, using the same segmentation.
- Enhanced Interpretability: The demographic profiles from Part 1 (e.g., Well-off as older, affluent, and educated) contextualize the purchase behaviors observed here. For instance, the Well-off segment's high purchase frequency and spending align with their higher income and urban lifestyle.
- Actionable Insights: Combining demographic and purchase data enables precise targeting. For example, the Career_Focused segment's low purchase frequency, coupled with their younger age and lower income, suggests a need for convenience-driven, budget-friendly promotions.

Without the prior segmentation, purchase analytics would lack a structured framework, making it difficult to attribute behaviors to specific customer groups. The pre-trained models ensure that the PCA components (career, education/lifestyle, experience) and K-means clusters remain consistent, providing a robust foundation for analysis.

## VI. MARKETING IMPLICATIONS

- The findings offer actionable insights for tailoring marketing strategies:
- Well-off: Target with premium products (e.g., organic foods, luxury personal care) and exclusive promotions like Promotion_1, leveraging their high spending power.
- Career_Focused: Offer time-saving products (e.g., ready-to-eat meals) and digital campaigns emphasizing convenience, addressing their low purchase frequency.
- Fewer_Opportunities: Promote value-driven products and discounts (e.g., Promotion_3) to capitalize on their growth potential and moderate spending.
- Standard: Focus on versatile, mid-range products through mass-market campaigns, aligning with their balanced preferences.

## VII. THEORETICAL CONTEXT

The analysis aligns with established marketing theories:

- **STP Framework:** By segmenting customers and analyzing their purchase behaviors, the study supports targeted positioning, as outlined in INSEAD Segmentation.
- **RFM Model:** The focus on purchase frequency and spending complements the RFM (Recency, Frequency, Monetary) model, which could be integrated in future analyses to enhance personalization (Customer Segmentation Using RFM).
- **Psychographic Segmentation:** The segments imply psychographic traits, such as Well-off customers valuing quality and Career_Focused prioritizing efficiency, aligning with principles in Neptune.ai Segmentation.

## VIII. RELATING TO BROADER FINDINGS

The results validate the importance of integrating segmentation with purchase analytics, as highlighted by Optimove, which notes that targeted campaigns based on segmentation can improve retention by 20-30%. The distinct purchase behaviors across segments support the hypothesis that demographic diversity drives varied consumer actions, a finding echoed in INSEAD's Cluster Analysis. The overlap in spending between Standard and Fewer_Opportunities suggests shared economic constraints, potentially explained by life cycle theory, where younger customers (Fewer_Opportunities) may transition to more stable profiles (Standard) as they age and gain financial stability.

## IX. STORY: THE CONSUMER JOURNEY

Consider an brand launching a new snack line. The Well-off segment, affluent and discerning, seeks premium organic snacks, responding to exclusive offers like Promotion_1. The Career_Focused segment, busy professionals, grabs convenient snack packs marketed via digital ads, aligning with their preference for Brand 1. The Fewer_Opportunities segment, budget-conscious youth, opts for value packs under Promotion_3, while the Standard segment, everyday consumers, chooses versatile snacks that balance quality and price. This segmentation, enabled by clustering and purchase analytics, transforms a generic launch into a tailored strategy, maximizing market penetration and customer satisfaction.

## X. LIMITATIONS

- **Demographic and Transactional Focus:** The analysis relies on demographic and transactional data, potentially missing psychographic or behavioral nuances.

- **Static Snapshot:** The findings represent a single point in time, not accounting for changes in customer behavior over time.

- **PCA Variance Loss:** While PCA retains ~80% of variance, some information loss may affect minor segment details.

## XI. RECOMMENDATIONS

**Well-off:** Launch premium product lines with exclusive branding, leveraging Promotion_1.
**Career_Focused:** Develop convenient products and target via digital channels, emphasizing Promotion_2.

**Fewer_Opportunities:** Offer bulk discounts and loyalty rewards under Promotion_3.
**Standard:** Promote mid-range products through mass-market campaigns, using Promotion_4.
**Future Research:** Integrate RFM metrics or conduct A/B testing on promotions to refine strategies.

## XII. CONCLUSION

This study successfully integrates purchase analytics with customer segmentation, providing a comprehensive view of customer behavior in the industry. By leveraging pre-trained models, it ensures consistency with the prior analysis, revealing that Well-off customers drive high purchase frequency and spending, while Career_Focused customers require targeted engagement. The findings offer actionable insights for optimizing marketing strategies, from premium offerings to value-driven promotions. Future research could explore predictive modeling or psychographic data to further enhance personalization, contributing to data-driven decision-making in the competitive landscape.

## XIII. KEY CITATIONS

- Customer Analytics in FMGC Industry Part 1 by Sooyeon Won
- Optimove Customer Clustering and Segmentation Guide
- Neptune.ai Customer Segmentation Using Machine Learning
- INSEAD Cluster Analysis and Segmentation in Marketing
- Customer Segmentation Using RFM Model