

Understanding Customer Behavior through Transactional Data Analysis

Introduction:

Understanding customer behavior through transactional data analysis is a crucial aspect of modern business strategy. By leveraging transactional data, businesses can gain valuable insights into customer preferences, buying habits, and trends, allowing them to tailor their products and services to meet customer needs effectively.

Objectives:

The objective of analyzing transactional data is to gain a comprehensive understanding of customer behavior. This includes studying buying habits, social trends, and background factors that influence purchasing decisions. By doing so, businesses can create more enticing products and service offers, improve customer experiences, and ultimately boost profitability. Additionally, the objective is to utilize data analytics to gain complete visibility into the buyer journey, analyze browsing time, transaction data, and purchase histories, and create specific customer segments for targeted marketing and personalized communication.

Data Set Description:

Hackathon_Ideal_Data: Contains brand level data for 10 stores for the last 3 months.

Hackathon_Working_Data: Includes data for selected stores with missing and/or incomplete information.

Hackathon_Mapping_File: Provides column names and descriptions for better understanding of the dataset.

Hackathon_Validation_Data: Data stores and product groups for which total value prediction is required.

Sample Submission: Represents the format for uploading output, including columns and values required.

Methodology:

- **Data Preprocessing:** Clean the dataset, handle missing values, and transform variables if necessary.
- **Exploratory Data Analysis (EDA):** Explore the dataset through summary statistics, visualizations, and correlation analysis to understand customer behavior.
- **Customer Segmentation:** Utilize clustering algorithms to identify distinct customer segments based on their purchasing behavior.
- **Predictive Modeling:** Develop predictive models using regression or time series analysis to forecast customer spending or predict total sales value for validation data.
- **Interpretation and Insights:** Interpret the results of analysis and modeling to derive actionable insights for business stakeholders.

Importing necessary libraries:

```
✓ [1] #Importing necessary libraries
1s import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from mlxtend.frequent_patterns import apriori
from mlxtend.frequent_patterns import association_rules
```

Importing datasets:

```
✓ 0s #Importing datasets
ideal_data=pd.read_csv("/content/Hackathon_Ideal_Data.csv")
mapping_file=pd.read_csv("/content/Hackathon_Mapping_File.csv")
validation_data=pd.read_csv("/content/Hackathon_Validation_Data.csv")
working_data=pd.read_csv("/content/Hackathon_Working_Data.csv")
sample_submission=pd.read_csv("/content/Sample_Submission.csv")
print(ideal_data.head())
print(mapping_file.head())
print(working_data.head())
print(sample_submission.head())
```

Display the dataset:

```
0  MONTH STORECODE QTY VALUE GRP SGRP \
1  M1 P1 25 83 HAIR CONDITIONERS HAIR CONDITIONERS
2  M1 P1 6 22 HAIR CONDITIONERS HAIR CONDITIONERS
3  M1 P1 4 15 HAIR CONDITIONERS HAIR CONDITIONERS
4  M1 P1 15 60 HAIR CONDITIONERS HAIR CONDITIONERS
5  M1 P2 0 0 HAIR CONDITIONERS HAIR CONDITIONERS

SSGRP CMP MBRD \
0 HAIR CONDITIONERS HINDUSTAN UNILEVER LIMITED DOVE
1 HAIR CONDITIONERS HINDUSTAN UNILEVER LIMITED DOVE
2 HAIR CONDITIONERS HINDUSTAN UNILEVER LIMITED DOVE
3 HAIR CONDITIONERS L'OREAL INDIA GARNIER
4 HAIR CONDITIONERS HINDUSTAN UNILEVER LIMITED CLINIC PLUS

BRD
0 DOVE HAIR FALL RESCUE
1 DOVE INTENSE REPAIR
2 DOVE OXYGEN MOISTURE
3 FRUCTIS
4 CLINIC PLUS

File Name Column Name Column Description
0 Hackathon_Ideal_Data MONTH Month ID (M1, M2, M3)
1 NaN STORECODE STORE CODE (P1, P2, ..., P10)
2 NaN QTY Sales Unit
3 NaN VALUE Sales Value
4 NaN GRP Category

MONTH STORECODE DAY BILL_ID BILL_AMT QTY VALUE PRICE \
0 M1 N1 4 T375 225.0 1.0 225.0 225.0
1 M1 N1 4 T379 95.0 1.0 95.0 95.0
2 M1 N1 4 T381 10.0 1.0 10.0 10.0
3 M1 N1 4 T382 108.0 1.0 108.0 108.0
4 M1 N1 4 T384 19.0 1.0 19.0 19.0

GRP SGRP SSGRP \
0 BUTTER MARGR (4/94) BUTTER BUTTER SALTED
1 CONFECTIONERY ECLAIRS CONFECTIONERY - ECLAIRS CONFECTIONERY - ECLAIRS
2 CHOCOLATE CHOCOLATE PANNEED CHOCOLATE PANNEED
3 PACKAGED TEA MAIN PACKS MAIN PACKS
4 ALL IODISED SALT POWDERED SALT POWDERED SALT

CMP MBRD BRD
0 G C M M F AMUL AMUL
1 PARLE PRODS MELODY MELODY CHOCOLATY
2 MONDELEZ INTERNATIONAL CADBURY SHOTS CADBURY SHOTS
3 GUJ TEA PROCESSORS WAGH BAKRI WAGH BAKRI INSTANT
4 TATA CHEM TATA TATA SALT

ID TOTALVALUE
0 1112535 0
1 1112539 1
2 1112543 2
3 1112547 3
4 1112551 4

/usr/local/lib/python3.10/dist-packages/ipykernel/ipkernel.py:283: DeprecationWarning: `should_run_as`
and `should_run_async` are deprecated
```

Handling Missing Values:

To identify missing values and removing the missing values:

```
0s #Handling missing data
print(ideal_data.isnull().sum())
is_any_missing_data=ideal_data.isna().any().any()
print("Are there any missing values?")
print(is_any_missing_data)
ideal_data.dropna(axis=0, inplace=True)
ideal_data.dropna(axis=1, inplace=True)
```

```
MONTH 0
STORECODE 0
QTY 0
VALUE 0
GRP 0
SGRP 0
SSGRP 0
CMP 0
MBRD 0
BRD 0
dtype: int64
Are there any missing values?
False
```

Removing duplicates values:

```
0s ▶ ideal_data.drop_duplicates(inplace=True)
    ideal_data.duplicated().sum()
    0

0s ▶ working_data.drop_duplicates(inplace=True)
    working_data.duplicated().sum()
    0
```

DATA PROCESSING AND FEATURE ENGINEERING

```
0s ▶ from sklearn.preprocessing import StandardScaler, OneHotEncoder
    from sklearn.impute import SimpleImputer
    imputer = SimpleImputer(strategy="mean")
    ideal_data_filled = imputer.fit_transform(ideal_data.select_dtypes(include=['int', 'float']))
    ideal_data_filled = pd.DataFrame(ideal_data_filled, columns=ideal_data.select_dtypes(include=['int', 'float']).columns)
```

Data cleaning, normalization and standardization:

```
0s [17] # Data cleaning, normalization, and standardization
    scaler = StandardScaler()
    working_data_scaled = scaler.fit_transform(working_data.select_dtypes(include=['int', 'float']))
    working_data_scaled = pd.DataFrame(working_data_scaled, columns=working_data.select_dtypes(include=['int', 'float']))

0s ▶ encoder = OneHotEncoder()
    ideal_data_encoded = pd.get_dummies(ideal_data.select_dtypes(include=['object']))

    # Concatenate numerical and encoded categorical features
    ideal_data_preprocessed = pd.concat([ideal_data_filled, ideal_data_encoded], axis=1)

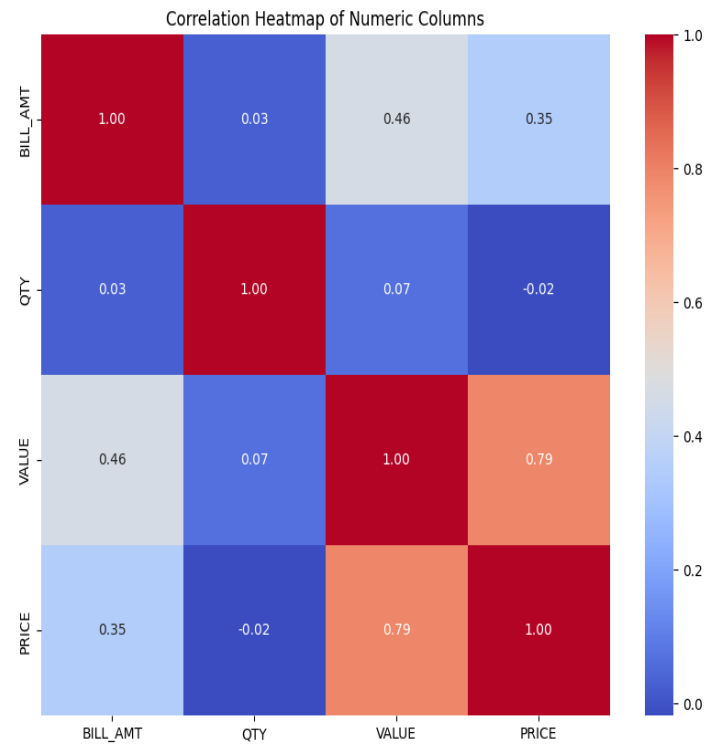
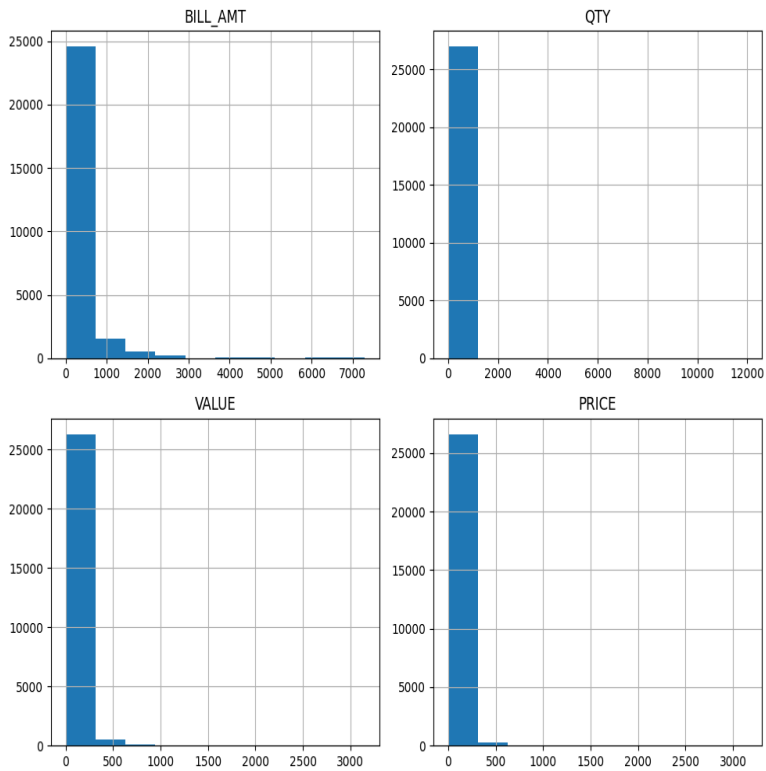
    # Print preprocessed data
    print("Ideal Data after preprocessing:")
    print(ideal_data_preprocessed.head())
    print("\nWorking Data after scaling:")
    print(working_data_scaled.head())

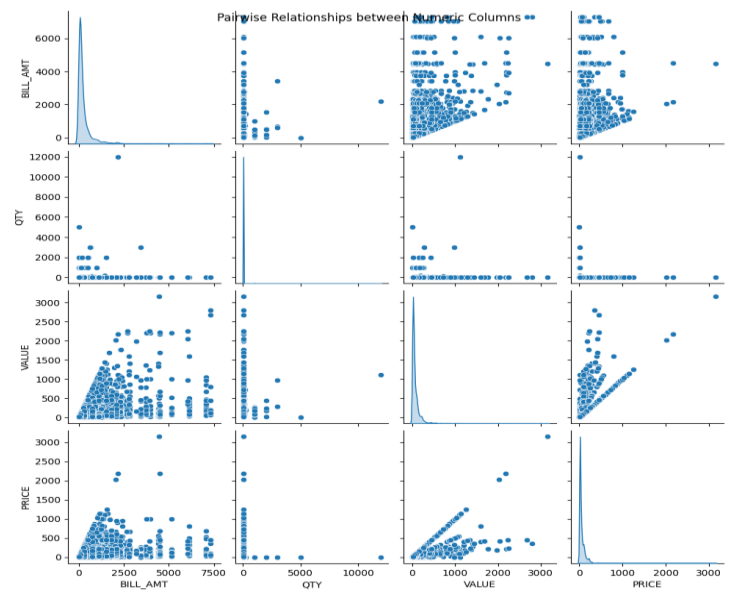
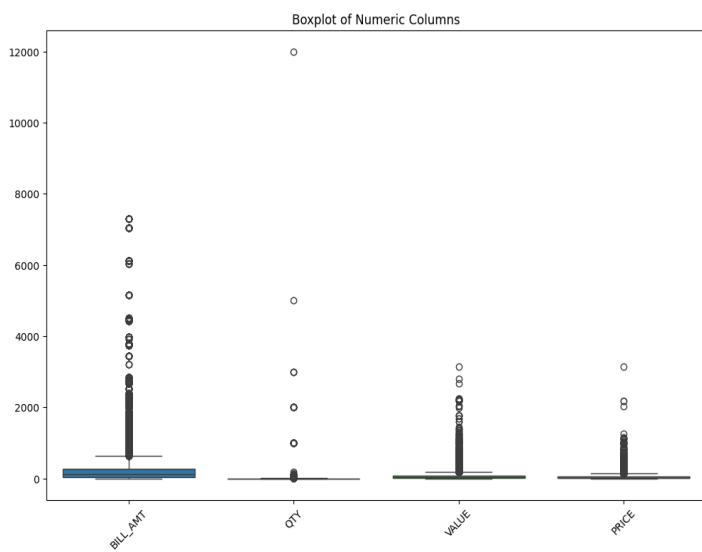
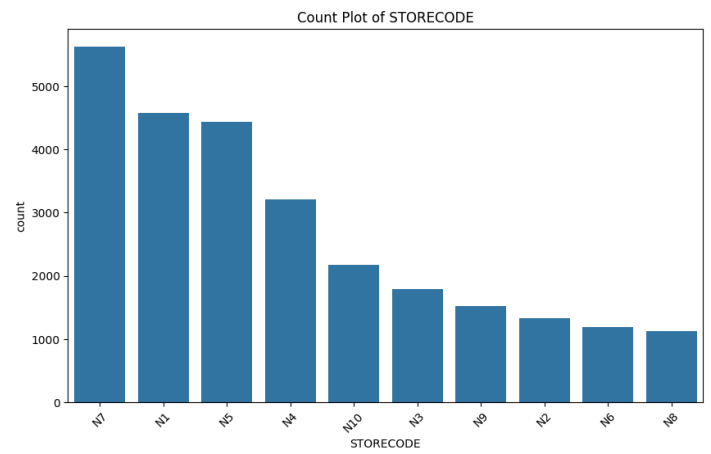
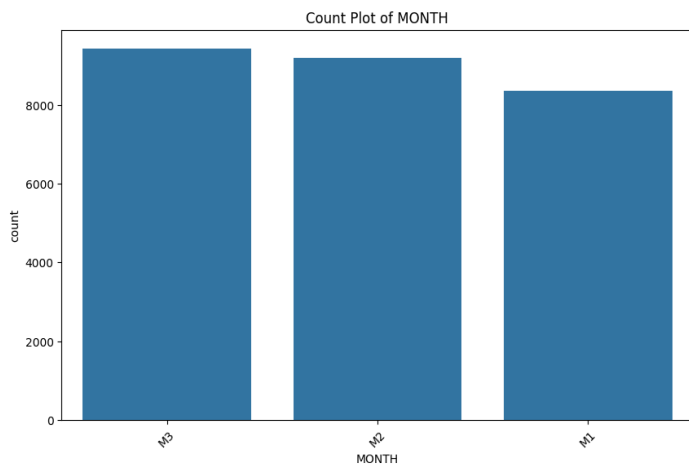
    ✓ 0s completed at 7:04 PM
```


EXPLORATORY DATA ANALYSIS (EDA)

```
# Plot boxplots to identify outliers
plt.figure(figsize=(12, 8))
sns.boxplot(data=working_data[numeric_cols])
plt.title('Boxplot of Numeric Columns')
plt.xticks(rotation=45)
plt.show()

# Plot scatter plots for pairwise relationships between numeric columns
sns.pairplot(working_data[numeric_cols], diag_kind='kde')
plt.suptitle('Pairwise Relationships between Numeric Columns')
plt.show()
```



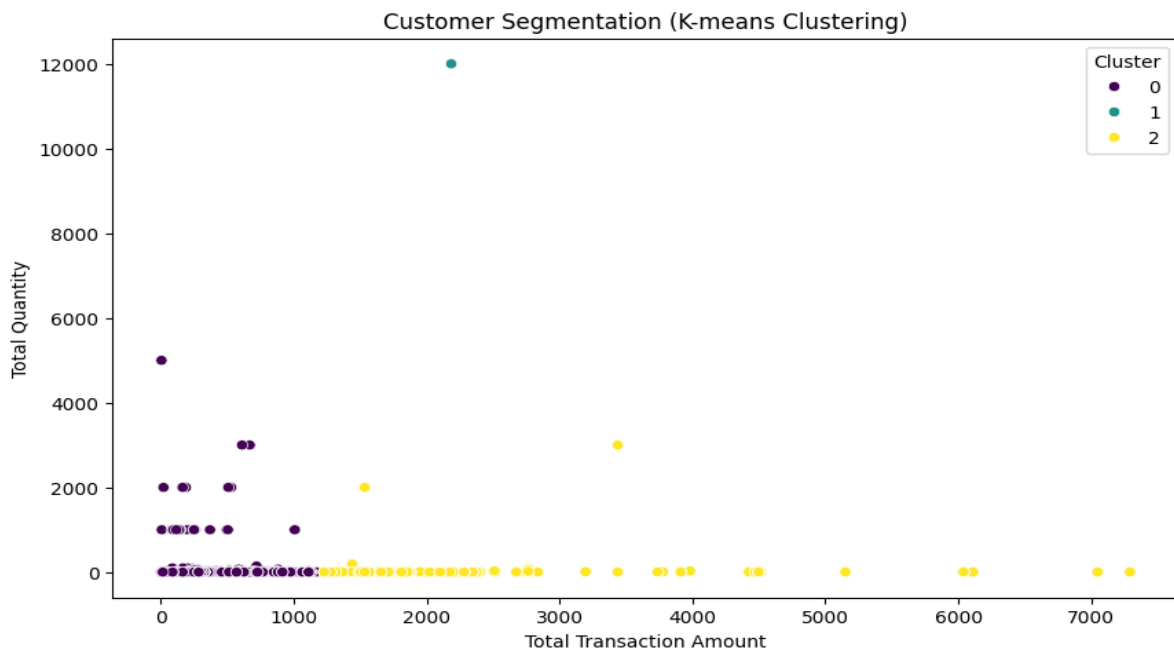


CUSTOMER SEGMENTATION WITH K-MEANS CLUSTERING

```
# Apply K-means clustering
kmeans = KMeans(n_clusters=n_clusters, init='k-means++', random_state=42)
working_data['Cluster'] = kmeans.fit_predict(X_scaled)

# Visualize clusters
plt.figure(figsize=(10, 6))
sns.scatterplot(data=working_data, x='BILL_AMT', y='QTY', hue='Cluster', palette='viridis')
plt.title('Customer Segmentation (K-means Clustering)')
plt.xlabel('Total Transaction Amount')
plt.ylabel('Total Quantity')
plt.legend(title='Cluster')
plt.show()
```

OUTPUT:



PREDICTIVE MODELING

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load the dataset
data = pd.read_csv("Hackathon_Working_Data.csv")
# Create histograms for transaction amounts and quantities
plt.figure(figsize=(12, 6))
plt.subplot(1, 2, 1)
sns.histplot(data['BILL_AMT'], bins=30, kde=True)
plt.title('Distribution of Transaction Amounts')
```



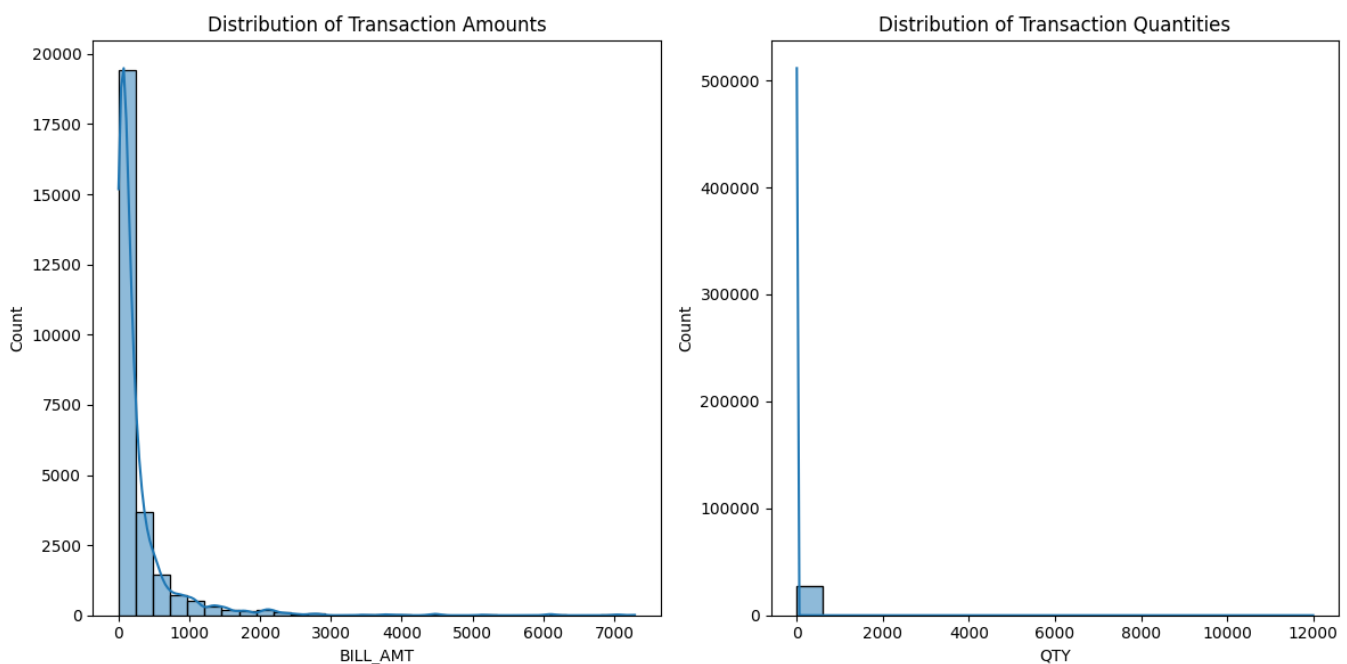
```
plt.subplot(1, 2, 2)
sns.histplot(data['QTY'], bins=20, kde=True)
plt.title('Distribution of Transaction Quantities')

plt.tight_layout()
plt.show()

# Box plot for transaction amounts
plt.figure(figsize=(8, 6))
sns.boxplot(y=data['BILL_AMT'])
plt.title('Boxplot of Transaction Amounts')
plt.show()

# Assuming 'MONTH' column represents the month
monthly_sales = data.groupby('MONTH').size()
```

OUTPUT:



✓
8s

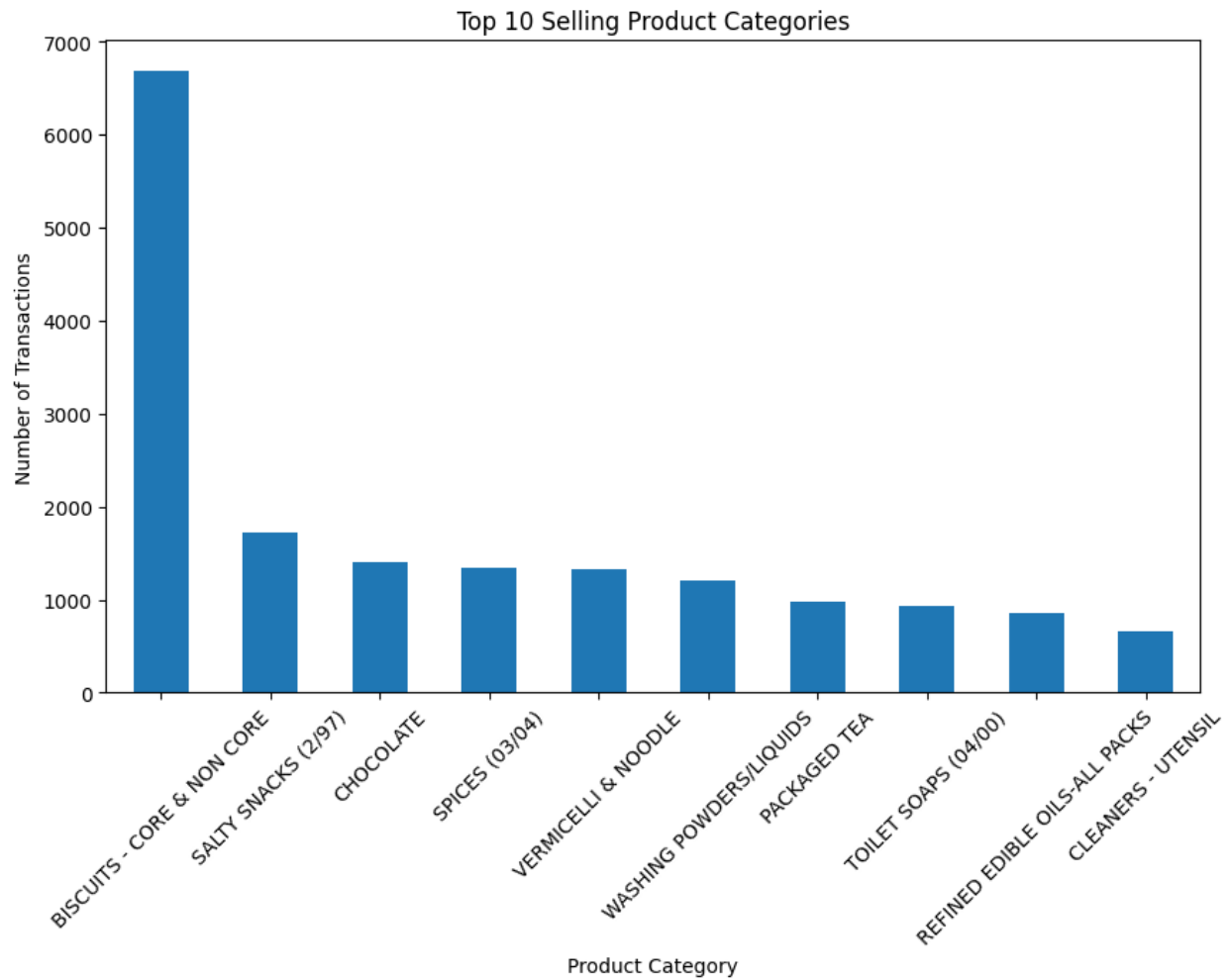


```
plt.figure(figsize=(10, 6))
monthly_sales.plot(kind='bar')
plt.title('Monthly Transaction Volumes')
plt.xlabel('Month')
plt.ylabel('Transaction Volume')
plt.xticks(rotation=45)
plt.show()

# Assuming 'GRP' column represents product categories
top_categories = data['GRP'].value_counts().nlargest(10)

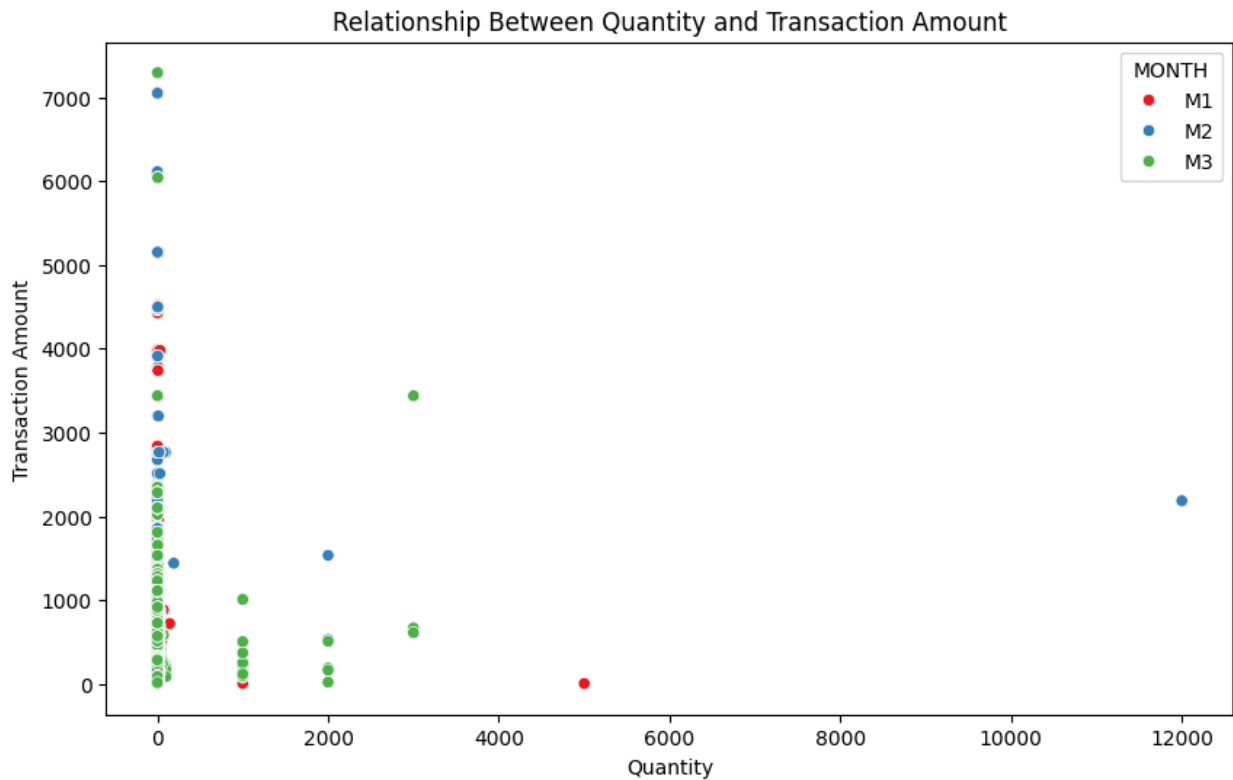
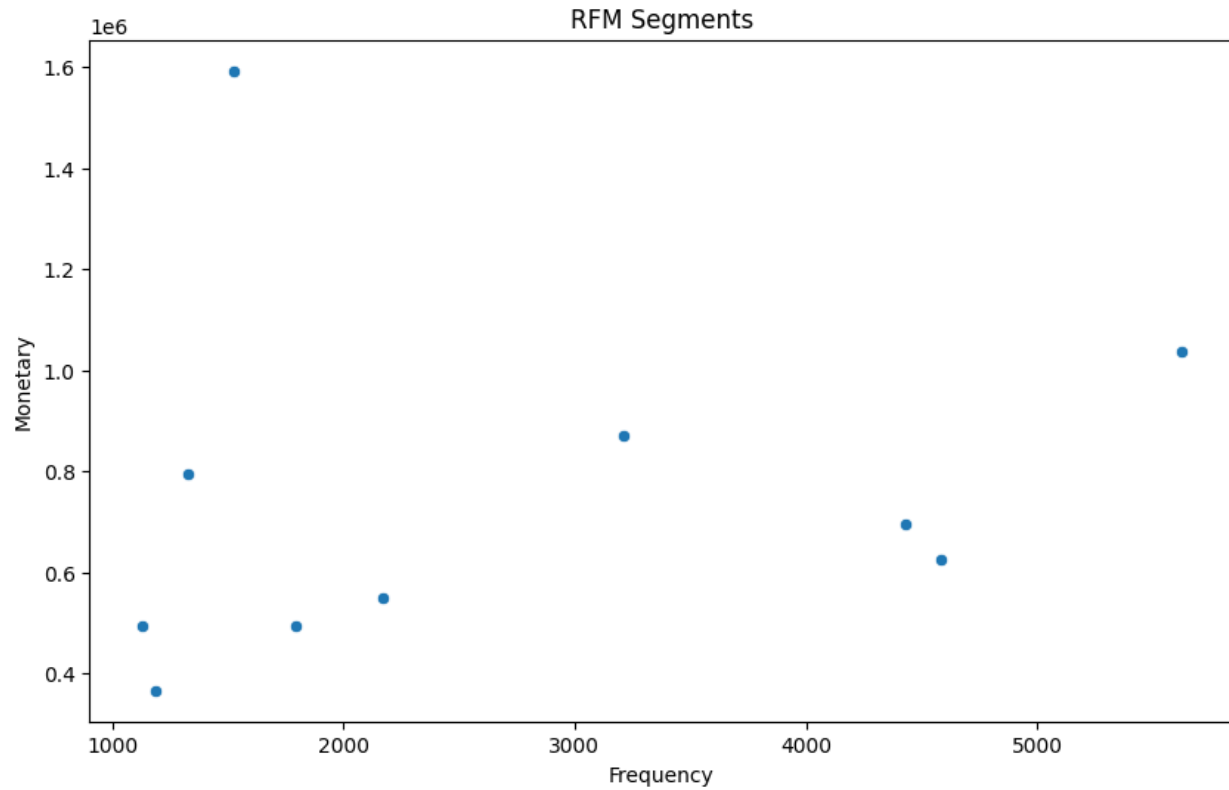
plt.figure(figsize=(10, 6))
top_categories.plot(kind='bar')
plt.title('Top 10 Selling Product Categories')
plt.xlabel('Product Category')
plt.ylabel('Number of Transactions')
plt.xticks(rotation=45)
plt.show()

# Conduct RFM analysis or clustering techniques
# Here, we'll use RFM analysis as an example
# Assuming 'STORECODE' represents unique customers
rfm_data = data.groupby('STORECODE').agg({
    'BILL_ID': 'count',          # Frequency
    'BILL_AMT': 'sum',          # Monetary
}).rename(columns={'BILL_ID': 'Frequency', 'BILL_AMT': 'Monetary'})
```



```
# Visualize RFM segments
plt.figure(figsize=(10, 6))
sns.scatterplot(data=rfm_data, x='Frequency', y='Monetary', palette='viridis')
plt.title('RFM Segments')
plt.xlabel('Frequency')
plt.ylabel('Monetary')
plt.show()

# Scatter plot to visualize relationship between quantity and transaction amount
plt.figure(figsize=(10, 6))
sns.scatterplot(data=data, x='QTY', y='BILL_AMT', hue='MONTH', palette='Set1')
plt.title('Relationship Between Quantity and Transaction Amount')
plt.xlabel('Quantity')
plt.ylabel('Transaction Amount')
plt.show()
```



CUSTOMER BEHAVIOR ANALYSIS:

For ideal_data:

+ Code + Text

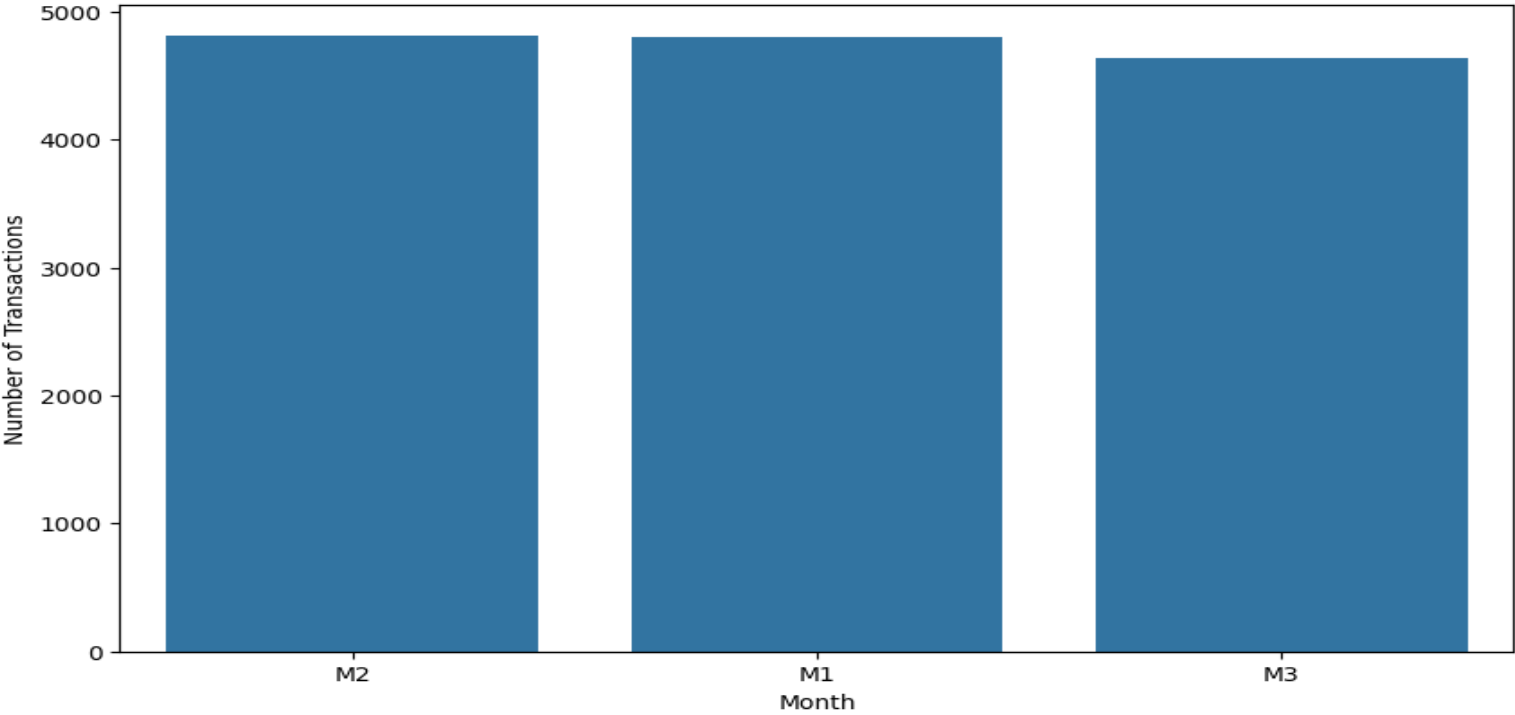
```
✓ 2s # For ideal_data
# Analyze purchasing patterns over time
plt.figure(figsize=(10, 6))
sns.countplot(data=ideal_data, x='MONTH', order=ideal_data['MONTH'].value_counts().index)
plt.title('Purchasing Patterns Over Time (Ideal Data)')
plt.xlabel('Month')
plt.ylabel('Number of Transactions')
plt.show()

# Identify popular product categories
plt.figure(figsize=(15, 10))
sns.countplot(data=ideal_data, y='GRP', order=ideal_data['GRP'].value_counts().index)
plt.title('Popular Product Categories (Ideal Data)')
plt.xlabel('Number of Transactions')
plt.ylabel('Product Category')
plt.xticks(fontsize=7)
plt.yticks(fontsize=7)
plt.show()

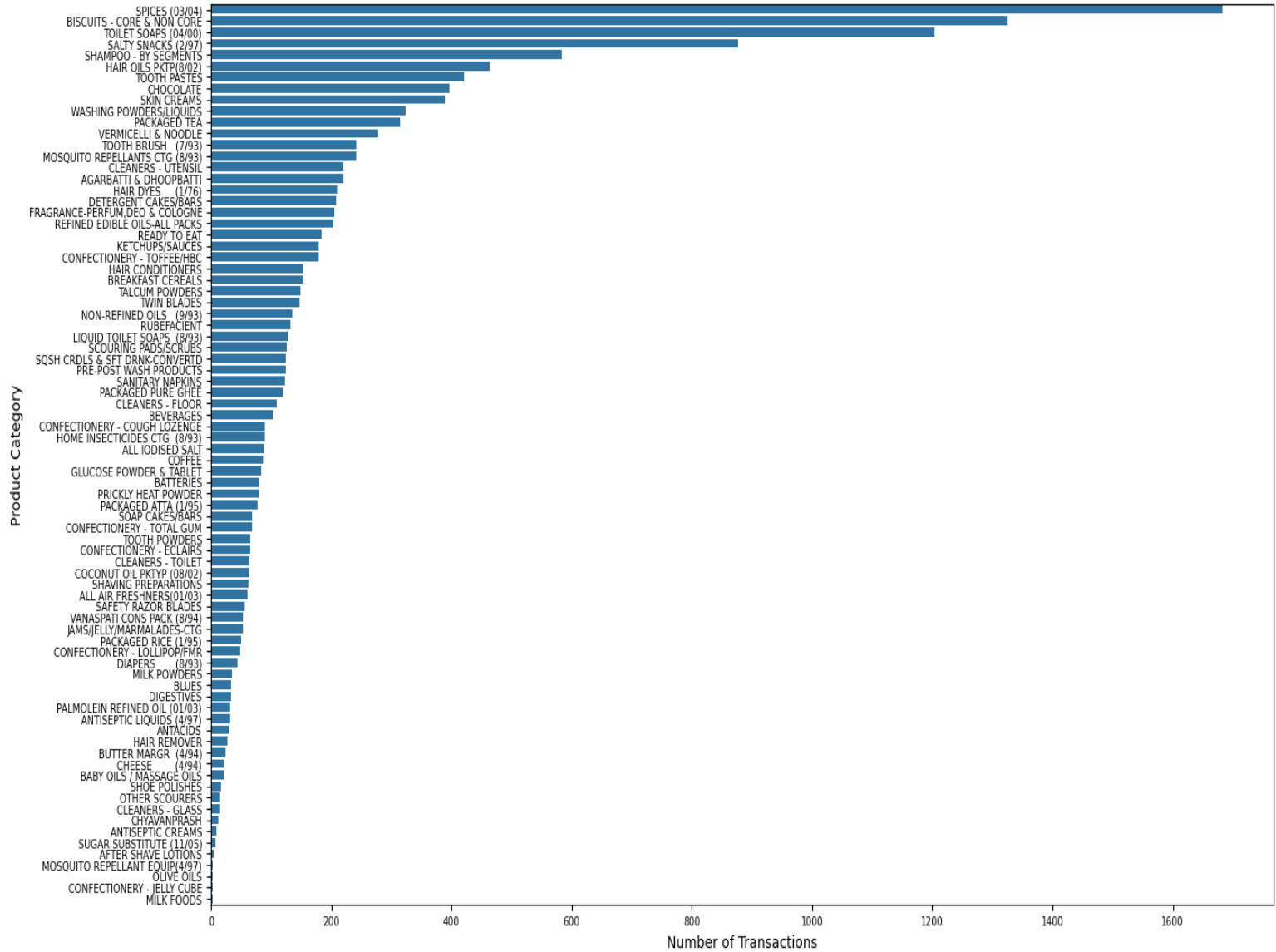
# Explore the relationship between quantity and value
plt.figure(figsize=(8, 6))
sns.scatterplot(data=ideal_data, x='QTY', y='VALUE')
plt.title('Relationship between Quantity and Value (Ideal Data)')
plt.xlabel('Quantity')
plt.ylabel('Value')
plt.show()

# Analyze the distribution of purchases across different stores
plt.figure(figsize=(10, 6))
sns.countplot(data=ideal_data, x='STORECODE', order=ideal_data['STORECODE'].value_counts().index)
plt.title('Distribution of Purchases Across Stores (Ideal Data)')
plt.xlabel('Store Code')
plt.ylabel('Number of Transactions')
plt.show()
```

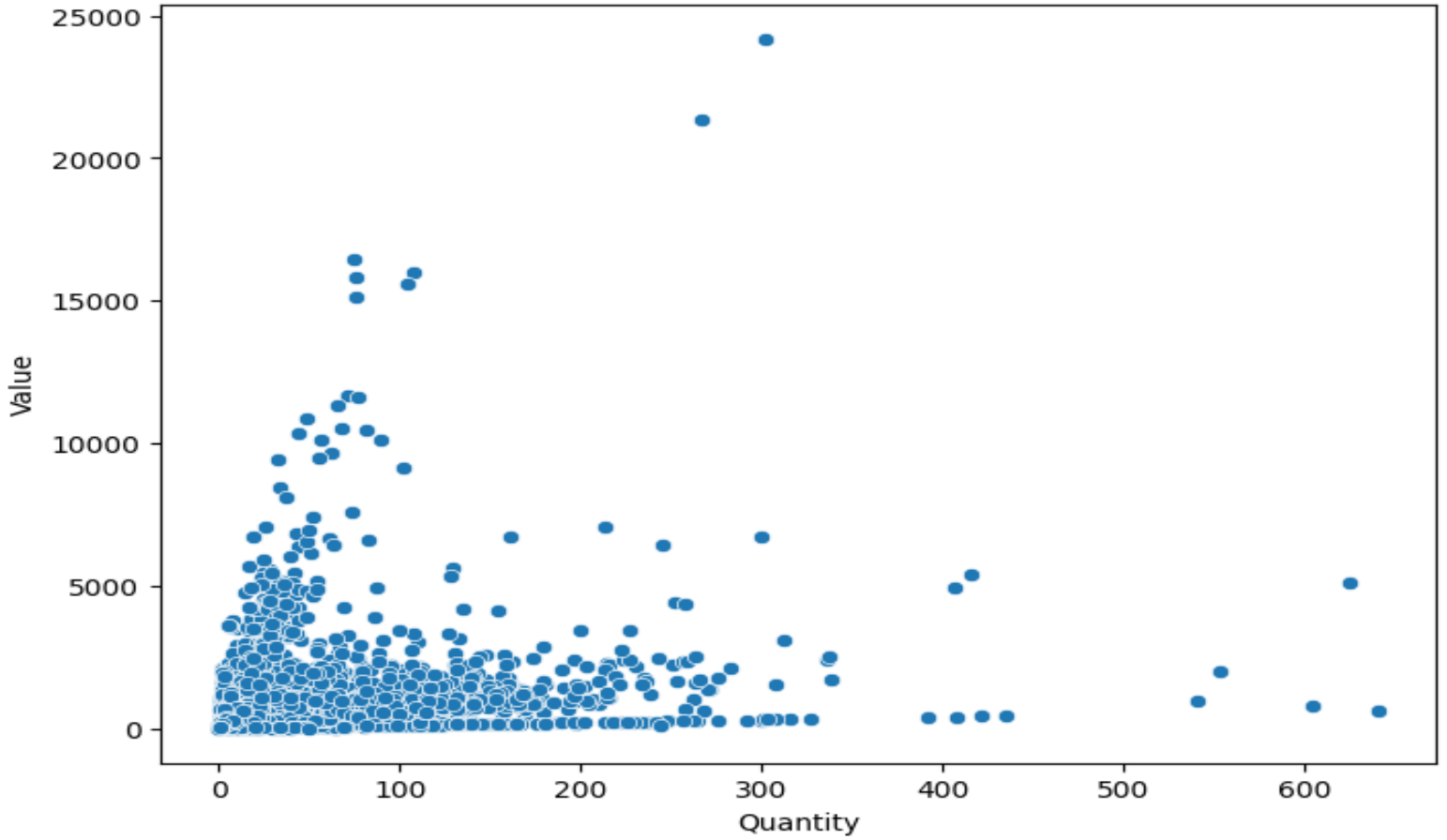
Purchasing Patterns Over Time (Ideal Data)



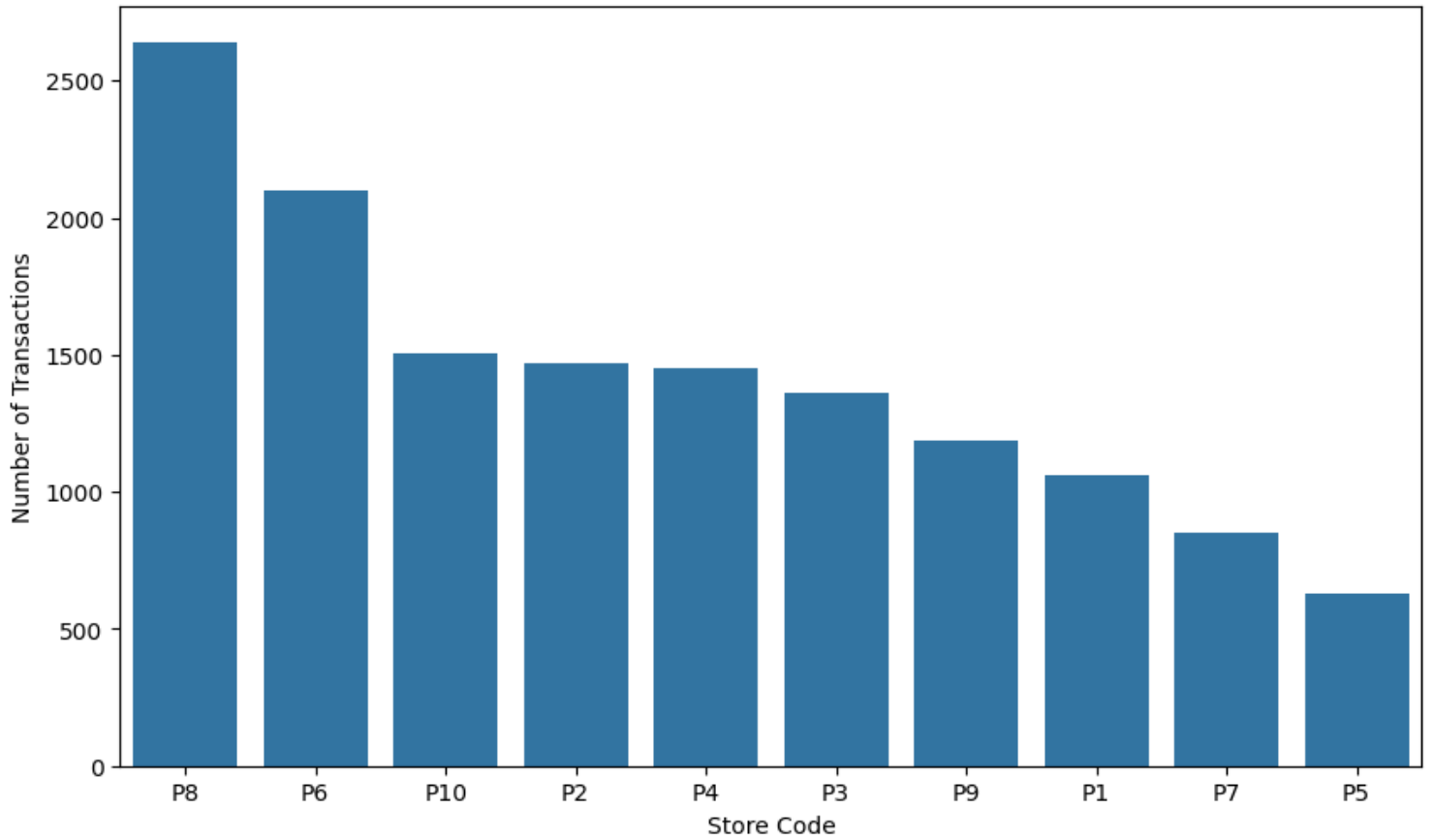
Popular Product Categories (Ideal Data)



Relationship between Quantity and Value (Ideal Data)



Distribution of Purchases Across Stores (Ideal Data)



For Working_data

✓
12s



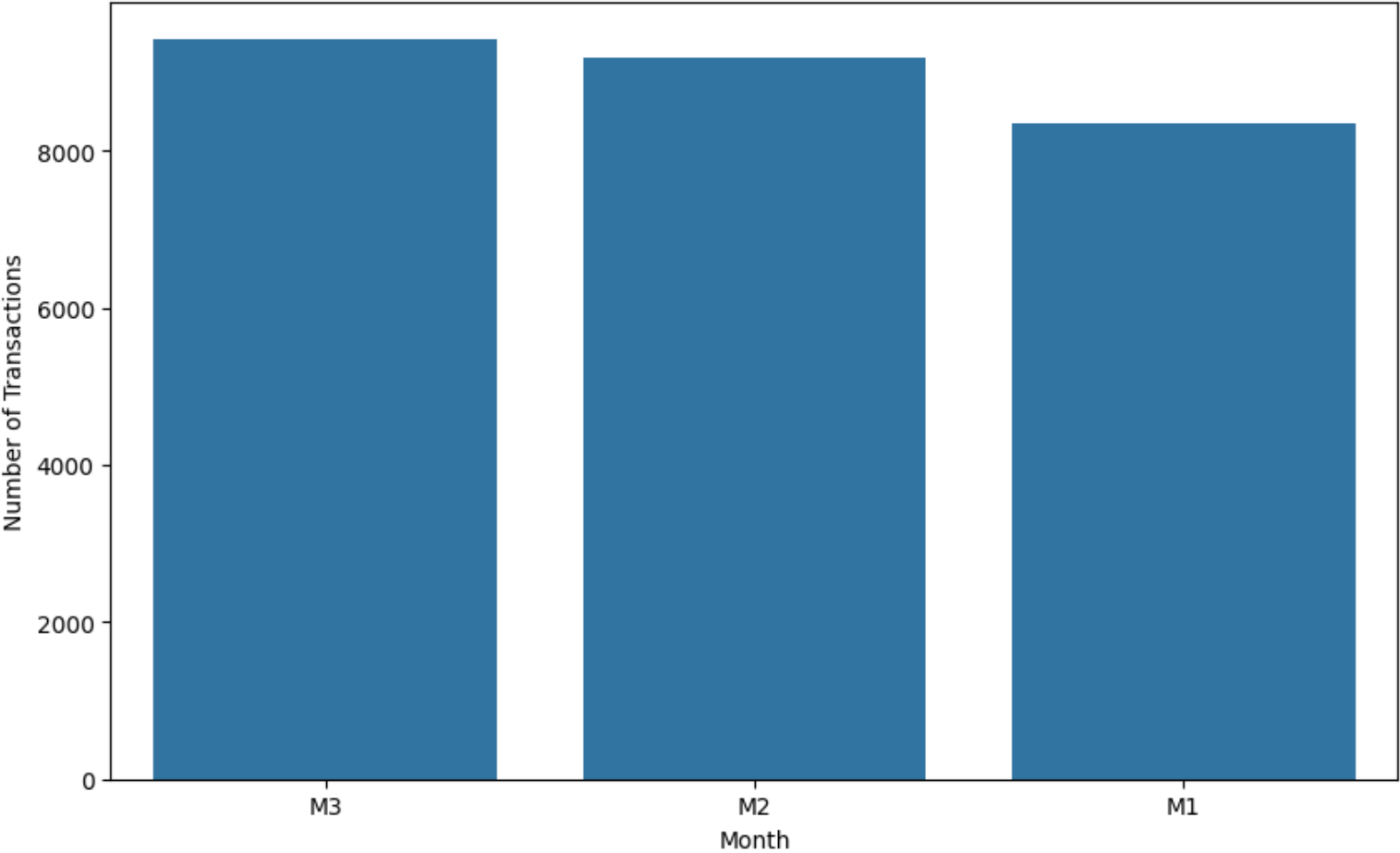
```
# For working_data
# Analyze purchasing patterns over time
plt.figure(figsize=(10, 6))
sns.countplot(data=working_data, x='MONTH', order=working_data['MONTH'].value_counts().index)
plt.title('Purchasing Patterns Over Time (Working Data)')
plt.xlabel('Month')
plt.ylabel('Number of Transactions')
plt.show()

# Identify popular product categories
plt.figure(figsize=(15, 10))
sns.countplot(data=working_data, y='GRP', order=working_data['GRP'].value_counts().index)
plt.title('Popular Product Categories (Working Data)')
plt.xlabel('Number of Transactions')
plt.ylabel('Product Category')
plt.xticks(fontsize=8)
plt.yticks(fontsize=7)
plt.show()

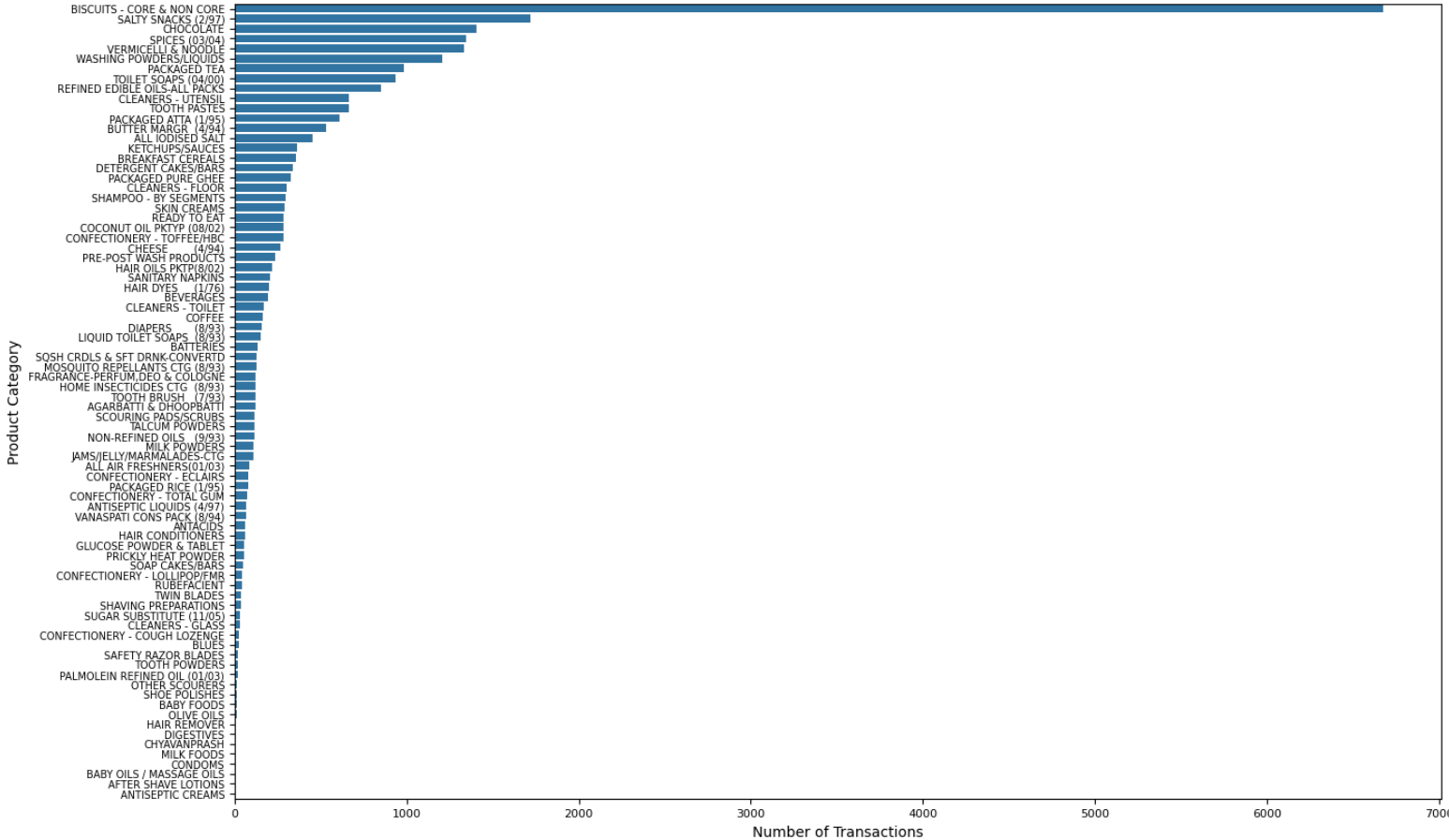
# Explore the relationship between quantity and value
plt.figure(figsize=(8, 6))
sns.scatterplot(data=working_data, x='QTY', y='VALUE')
plt.title('Relationship between Quantity and Value (Working Data)')
plt.xlabel('Quantity')
plt.ylabel('Value')
plt.show()

# Analyze the distribution of purchases across different stores
plt.figure(figsize=(10, 6))
sns.countplot(data=working_data, x='STORECODE', order=working_data['STORECODE'].value_counts().index)
plt.title('Distribution of Purchases Across Stores (Working Data)')
plt.xlabel('Store Code')
plt.ylabel('Number of Transactions')
plt.show()
```

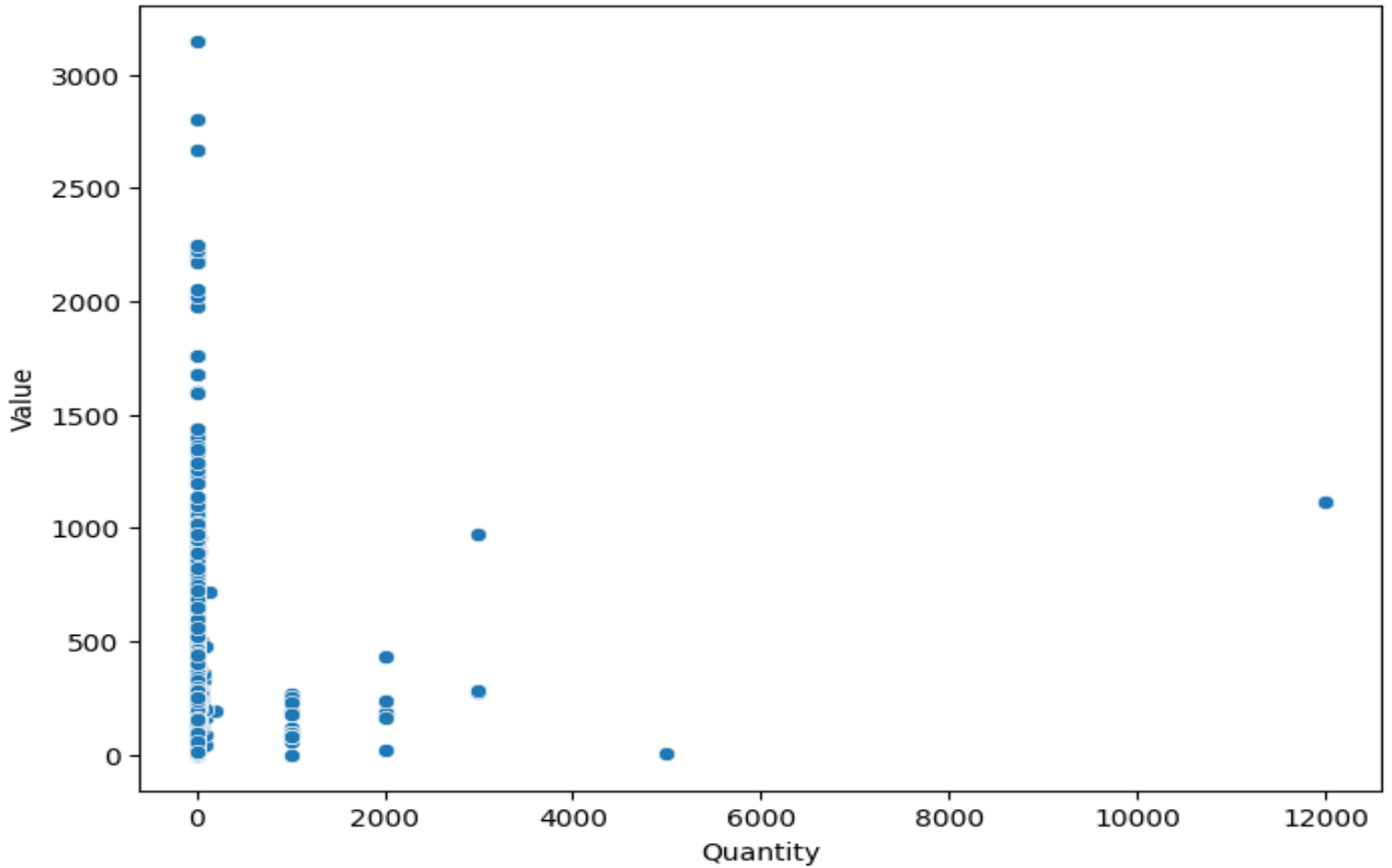
Purchasing Patterns Over Time (Working Data)



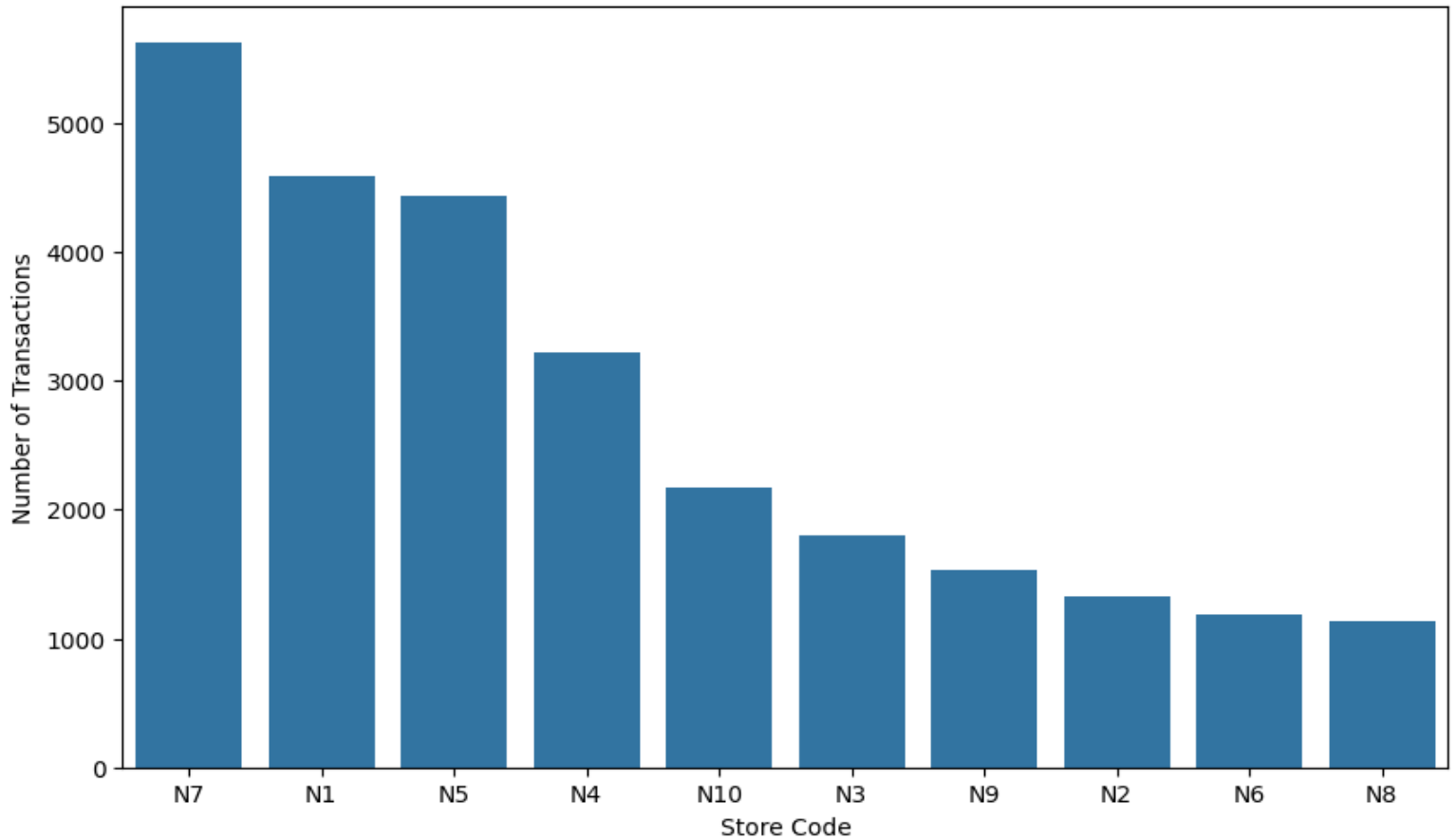
Popular Product Categories (Working Data)



Relationship between Quantity and Value (Working Data)



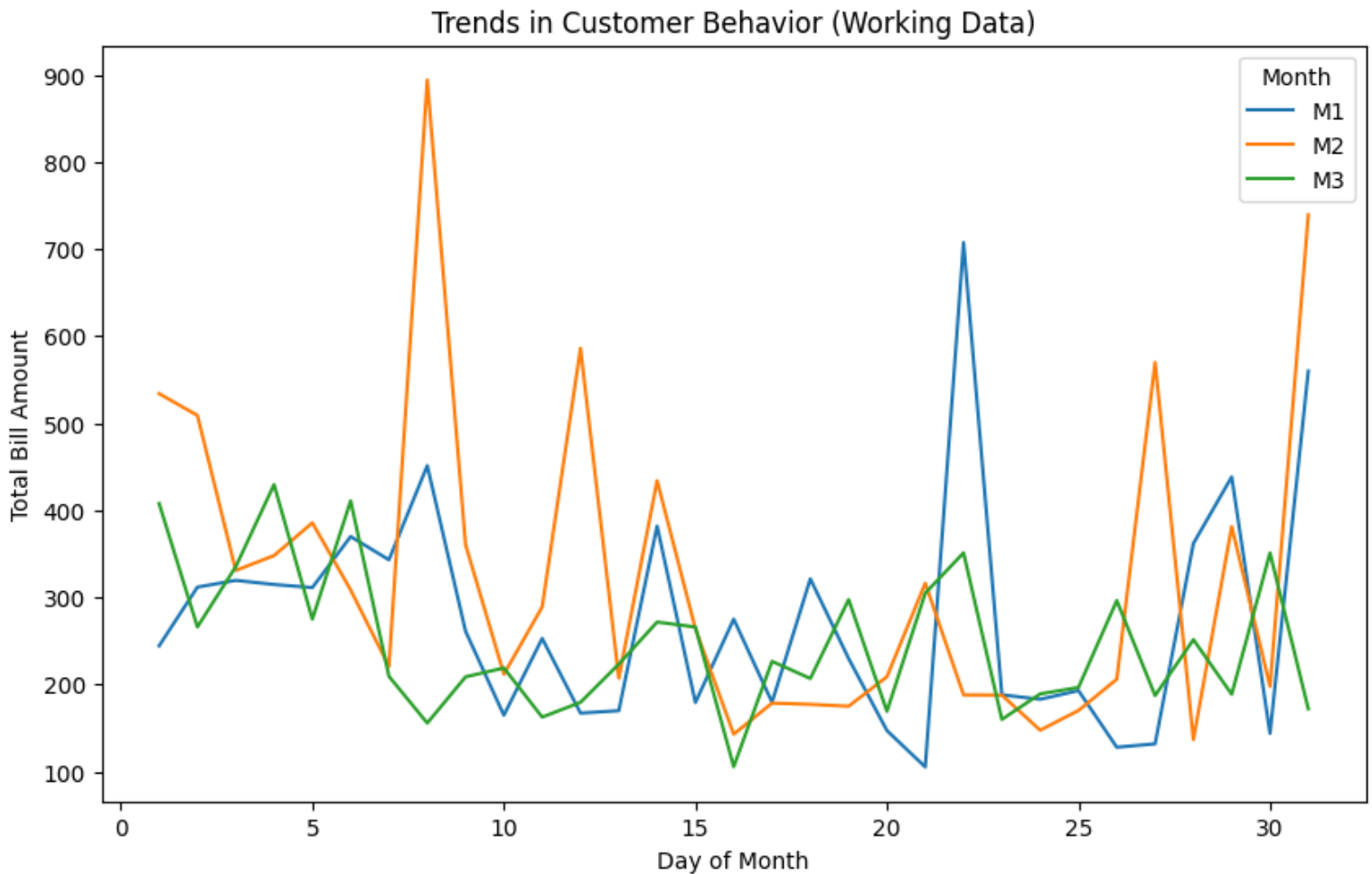
Distribution of Purchases Across Stores (Working Data)



VISUALIZATION AND REPORTING:

Customer Behavior:

```
✓ 1s ▶ plt.figure(figsize=(10, 6))
sns.lineplot(data=working_data, x='DAY', y='BILL_AMT', hue='MONTH', ci=None)
plt.title('Trends in Customer Behavior (Working Data)')
plt.xlabel('Day of Month')
plt.ylabel('Total Bill Amount')
plt.legend(title='Month')
plt.show()
```



Conclusion :

In conclusion, the analysis of customer behavior based on transactional data has provided valuable insights into various aspects of customer preferences, purchasing patterns, and overall behavior. Through exploratory data analysis (EDA), we have uncovered key trends, such as the distribution of sales across different stores, brands, and categories, as well as fluctuations in sales over time. By segmenting customers based on their purchasing behavior, we have identified distinct customer groups and their unique preferences.

Moreover, predictive modeling has allowed us to forecast total sales accurately, enabling better decision-making for inventory management and marketing strategies. These insights can be leveraged to tailor marketing campaigns, optimize product offerings, and enhance customer experience, ultimately driving business growth and profitability.

Moving forward, continuous monitoring of customer behavior and refinement of analytical techniques will be crucial for staying responsive to evolving market dynamics and customer needs. By adopting a data-driven approach, businesses can adapt proactively to changes in consumer behavior and maintain a competitive edge in the market.