



MARKET BASKET ANALYSIS WITH **PYTHON**

BY VRINDA BHATT

PROJECT OVERVIEW

You are hired as a Customer Insights Analyst for Amazon, a leading global e-commerce platform. With millions of customers making purchases across diverse categories, the leadership team wants to understand customer purchasing behavior, product affinities, and engagement trends to enhance personalized shopping experiences and optimize inventory and marketing strategies.

This dataset captures Amazon customers' online shopping behavior, including demographics, purchase frequency, browsing habits, product search preferences, review engagement, satisfaction levels, and areas for service improvement.

TASK 1: DATA CLEANING AND PREPARATION

- REMOVE DUPLICATE OR INCONSISTENT SURVEY RESPONSES.
- STANDARDIZE CATEGORICAL ENTRIES (E.G., FREQUENCY LEVELS, GENDER, RECOMMENDATION RESPONSES).
- HANDLE MISSING VALUES AND INCONSISTENT FORMATS IN PRODUCT_SEARCH_METHOD AND OTHER FIELDS.
- RENAME DUPLICATE OR MISFORMATTED COLUMNS (E.G., REMOVE TRAILING SPACES IN RATING_ACCURACY).
- CONVERT NUMERICAL RATING COLUMNS (E.G., CUSTOMER_REVIEWS_IMPORTANCE, SHOPPING_SATISFACTION) TO APPROPRIATE NUMERIC TYPES FOR ANALYSIS.

TASK 1: DATA CLEANING AND PREPARATION

RAW DATA

Timestamp	age	Gender	Purchase_Less than 6 months	Purchase_Once a week	Personaliz_Clothing	Browsing_Multiple times	Product_Type Keyword	Search_Recent	ReCustomer	Add_to_Cart	Cart_Combo	Cart_Average	Saveforlater	Review_Likelihood	Review_Rate	Review_Honesty	Personaliz_Beauty	Recommendation	Rating_Accuracy
2023/06/01	65	Prefer not to say	Less than 6 months	Clothing	Yes	Multiple times	Keyword	Multiple pages	No	Sometimes	Found a lot	A few times	Sometimes	Yes	Moderate	No	2	Sometime	4
2023/06/01	20	Male	Once a week	Groceries	No	Rarely	Filter	First page	No	Never	High shipping	Always	Yes	Heavily	Yes	1	Sometime	4	
2023/06/01	42	Male	Once a week	Groceries	Sometimes	Few times	Keyword	Multiple pages	No	Rarely	Found a lot	Often	Yes	Heavily	Sometime	5	No	5	
2023/06/01	65	Others	Once a month	Beauty and Personal care	No	Few times	Filter	Multiple pages	Yes	Sometimes	others	Often	Yes	Occasionally	No	3	Yes	1	
2023/06/01	45	Female	Once a week	Beauty and Personal care	Sometimes	Few times	nan	First page	Maybe	Rarely	Changed rarely	Never	Yes	Rarely	No	2	Yes	1	
2023/06/01	67	Others	Few times	Groceries	No	Multiple times	nan	Multiple pages	No	Always	others	Never	No	Moderately	Sometime	4	No	4	
2023/06/01	55	Female	Few times	Groceries	Sometimes	Few times	Keyword	First page	Maybe	Never	Found a lot	Rarely	Yes	Moderately	No	1	Yes	5	
2023/06/01	44	Prefer not to say	Few times	Beauty and Personal care	Yes	Few times	nan	Multiple pages	Yes	Often	Changed rarely	Rarely	No	Heavily	Sometime	2	Sometime	5	
2023/06/01	33	Male	Less than 6 months	Beauty and Personal care	Yes	Rarely	nan	Multiple pages	Yes	Often	High shipping	Rarely	No	Heavily	Sometime	1	Sometime	2	
2023/06/01	22	Others	Few times	Clothing and accessories	No	Multiple times	Filter	First page	Maybe	Never	Changed rarely	Rarely	No	Occasionally	No	2	Yes	1	
2023/06/01	28	Female	Once a month	Groceries	Yes	Few times	others	First page	Yes	Sometimes	Changed rarely	Sometimes	Yes	Occasionally	Yes	1	Yes	4	
2023/06/01	44	Prefer not to say	Multiple times	Groceries	Sometimes	Multiple times	Keyword	First page	Yes	Always	others	Often	Yes	Occasionally	Yes	5	Yes	4	
2023/06/01	24	Male	Once a month	Beauty and Personal care	No	Few times	Filter	First page	No	Always	others	Rarely	Yes	Moderately	Yes	4	Yes	5	
2023/06/01	23	Male	Once a month	Beauty and Personal care	Yes	Few times	Filter	First page	Maybe	Always	Changed rarely	Sometimes	No	Moderately	Yes	1	Sometime	5	
2023/06/01	58	Prefer not to say	Once a week	Beauty and Personal care	Sometimes	Few times	Keyword	First page	Maybe	Never	High shipping	Often	Yes	Never	Yes	3	No	3	
2023/06/01	40	Male	Less than 6 months	Groceries	No	Multiple times	Keyword	First page	No	Never	Found a lot	Never	No	Never	No	4	No	2	
2023/06/01	46	Prefer not to say	Once a week	Groceries	Yes	Rarely	categories	First page	Yes	Always	High shipping	Sometimes	No	Rarely	Sometime	2	Sometime	5	

CLEAN DATA

Timestamp	age	Gender	Purchase_Less than 6 months	Purchase_Once a week	Personaliz_Clothing	Browsing_Multiple times	Product_Type Keyword	Search_Recent	ReCustomer	Add_to_Cart	Cart_Combo	Cart_Average	Saveforlater	Review_Likelihood	Review_Rate	Review_Honesty	Personaliz_Beauty	Recommendation	Rating_Accuracy	Shopping_Habits	Service_Accuracy	Interaction_Level
2023/06/01	65	Other	Less than 6 months	Clothing	Yes	Multiple times	Keyword	Multiple pages	No	Sometimes	Found a lot	A few times	Sometimes	Yes	Moderate	No	Sometime	4	4	Competitive	4	
2023/06/01	20	Male	Once a week	Groceries	No	Rarely	Filter	First page	No	Never	High shipping	Always	Yes	Heavily	Yes	Sometime	4	5	Quick delivery	5		
2023/06/01	42	Male	Once a week	Groceries	Sometimes	Few times	Keyword	Multiple pages	No	Rarely	Found a lot	Often	Yes	Heavily	Sometime	No	5	3	All the above	U		
2023/06/01	65	others	Once a month	Beauty and Personal care	No	Few times	Filter	Multiple pages	Yes	Sometimes	others	Often	Yes	Occasionally	No	Yes	1	2	Quick delivery	Q		
2023/06/01	45	Female	Once a week	Beauty and Personal care	Sometimes	Few times	Unknown	First page	Maybe	Rarely	Changed rarely	Never	Yes	Rarely	No	Yes	1	2	Quick delivery	Q		
2023/06/01	67	others	Few times	Groceries	No	Multiple times	Unknown	Multiple pages	No	Always	others	Never	No	Moderately	Sometime	No	4	3	Product range	U		
2023/06/01	55	Female	Few times	Groceries	Sometimes	Few times	Keyword	First page	Maybe	Never	Found a lot	Rarely	Yes	Moderately	No	Yes	5	3	Customer service	U		
2023/06/01	44	Other	Few times	Beauty and Personal care	Yes	Few times	Unknown	Multiple pages	Yes	Often	Changed rarely	Rarely	No	Heavily	Sometime	Sometime	5	1	Customer service	I		
2023/06/01	33	Male	Less than 6 months	Beauty and Personal care	Yes	Rarely	Unknown	Multiple pages	Yes	Often	High shipping	Rarely	No	Heavily	Sometime	Sometime	2	2	Competitive price	I		
2023/06/01	22	others	Few times	Clothing and accessories	No	Multiple times	Filter	First page	Maybe	Never	Changed rarely	Rarely	No	Occasionally	No	Yes	1	3	All the above	U		
2023/06/01	28	Female	Once a month	Groceries	Yes	Few times	others	First page	Yes	Sometimes	Changed rarely	Sometimes	Yes	Occasionally	Yes	Yes	4	5	All the above	U		
2023/06/01	44	Other	multiple times	Groceries	Sometimes	Multiple times	Keyword	First page	Yes	Always	others	Often	Yes	Occasionally	Yes	Yes	4	3	User-friendly interface	U		
2023/06/01	24	Male	Once a month	Beauty and Personal care	No	Few times	Filter	First page	No	Always	others	Rarely	Yes	Moderately	Yes	Yes	5	3	Competitive price	I		
2023/06/01	23	Male	Once a month	Beauty and Personal care	Yes	Few times	Filter	First page	Maybe	Always	Changed rarely	Sometimes	No	Moderately	Yes	Sometime	5	1	Wide product range	U		
2023/06/01	58	Other	Once a week	Beauty and Personal care	Sometimes	Few times	Keyword	First page	Maybe	Never	High shipping	Often	Yes	Never	Yes	No	3	1	Quick delivery	Q		
2023/06/01	40	Male	Less than 6 months	Groceries	No	Multiple times	Keyword	First page	No	Never	Found a lot	Never	No	Never	No	No	2	2	Customer service	U		
2023/06/01	46	Other	Once a week	Groceries	Yes	Rarely	categories	First page	Yes	Always	High shipping	Sometimes	No	Rarely	Sometime	Sometime	5	2	Wide product range	U		
2023/06/01	15	others	Once a month	Groceries	Sometimes	Few times	Keyword	Multiple pages	Maybe	Often	others	Often	Yes	Never	No	Yes	2	2	Quick delivery	Q		

TASK 1: DATA CLEANING AND PREPARATION

CODES PERFORMED IN PYTHON

```
import pandas as pd
import numpy as np

df = pd.read_csv("../data/amazon.csv")
df.head()

df.shape

df.columns

df.columns = df.columns.str.strip().str.replace(" ", "_")
df.columns

df = df.loc[:, ~df.columns.duplicated()]
df.columns

df.isnull().sum()

df['Product_Search_Method'] = df['Product_Search_Method'].fillna('Unknown')

rating_cols = [
    'Customer_Reviews_Importance',
    'Shopping_Satisfaction',
    'Rating_Accuracy'
]

for col in rating_cols:
    if col in df.columns:
        df[col] = pd.to_numeric(df[col], errors='coerce')

df['Gender'] = df['Gender'].str.lower().str.strip()
```

```
df['Gender'] = df['Gender'].str.lower().str.strip()

df['Gender'] = df['Gender'].replace({
    'male': 'Male',
    'female': 'Female',
    'prefer not to say': 'Other'
})

df['Purchase_Frequency'] = df['Purchase_Frequency'].str.lower().str.strip()
df.columns = df.columns.str.strip().str.replace(" ", "_")

df['Improvement_Areas'] = df['Improvement_Areas'].str.lower().str.strip()

garbage_values = ['.', 'nill', 'nil', 'none', 'na', 'n/a', '']

df['Improvement_Areas'] = df['Improvement_Areas'].replace(garbage_values, 'Unknown')

df['Improvement_Areas'] = df['Improvement_Areas'].replace({
    'ui': 'User Interface',
    'user interface': 'User Interface',
    'user interfacee': 'User Interface',
    'user interface of app': 'User Interface',
    'app ui': 'User Interface'
})

df['Improvement_Areas'].value_counts()

import os
os.getcwd()

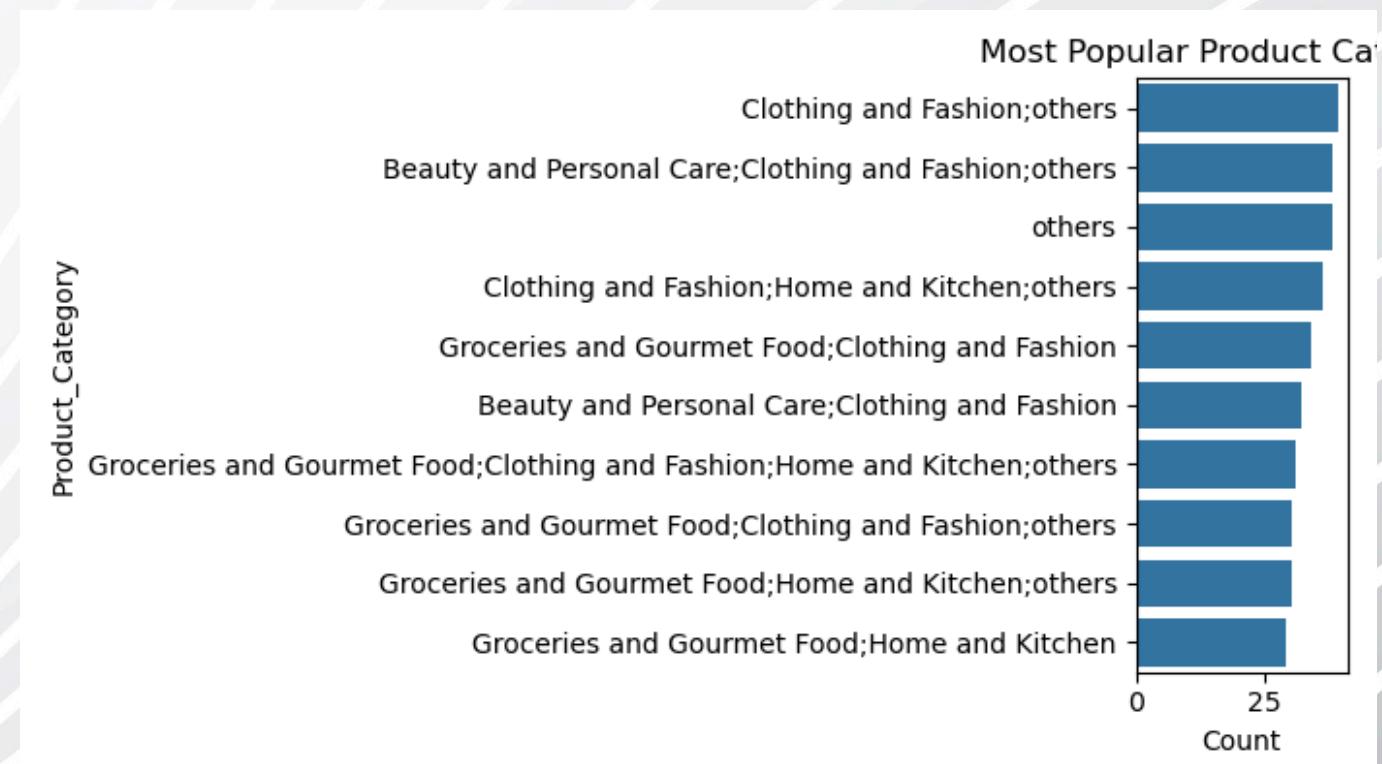
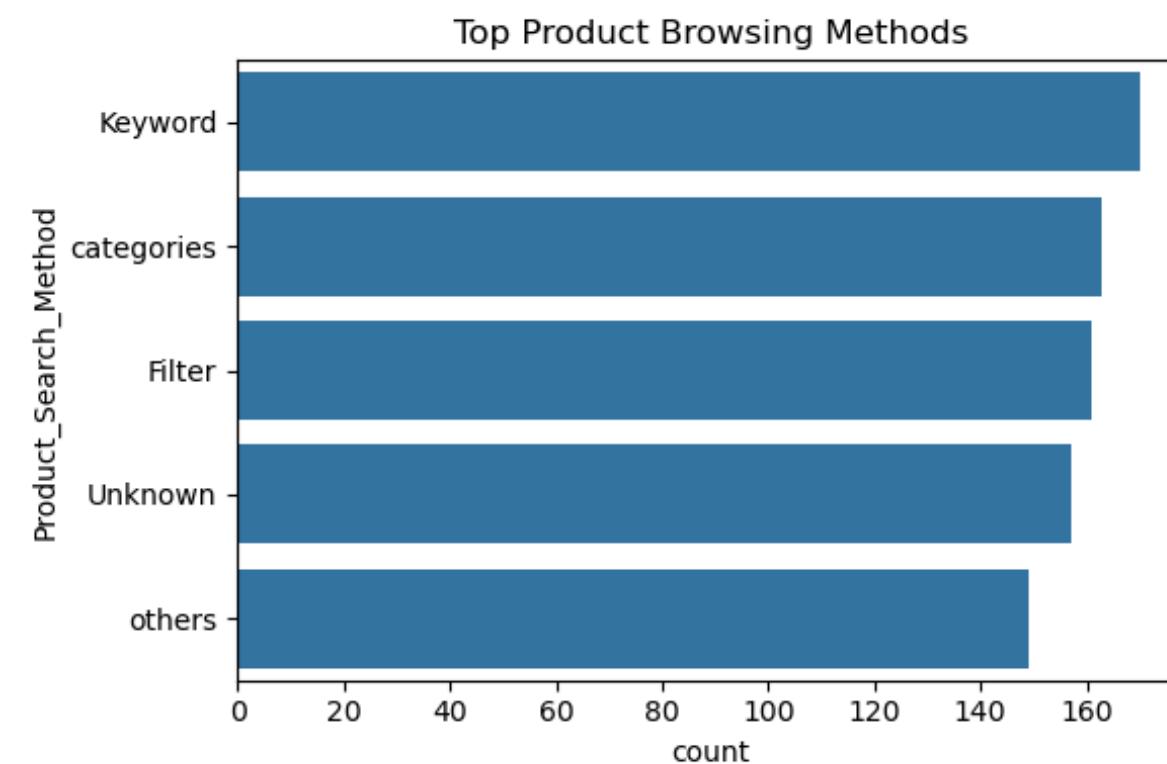
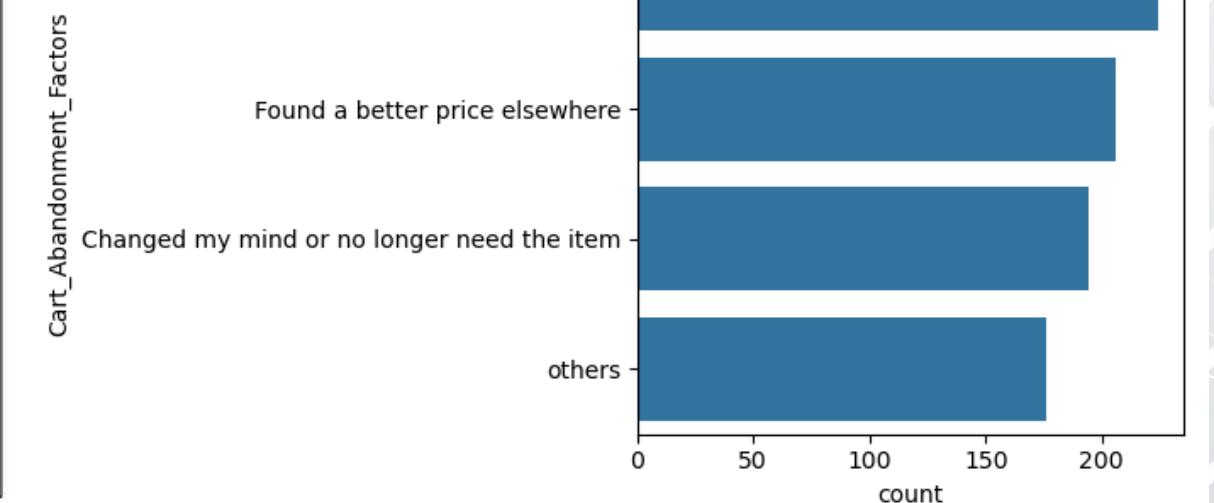
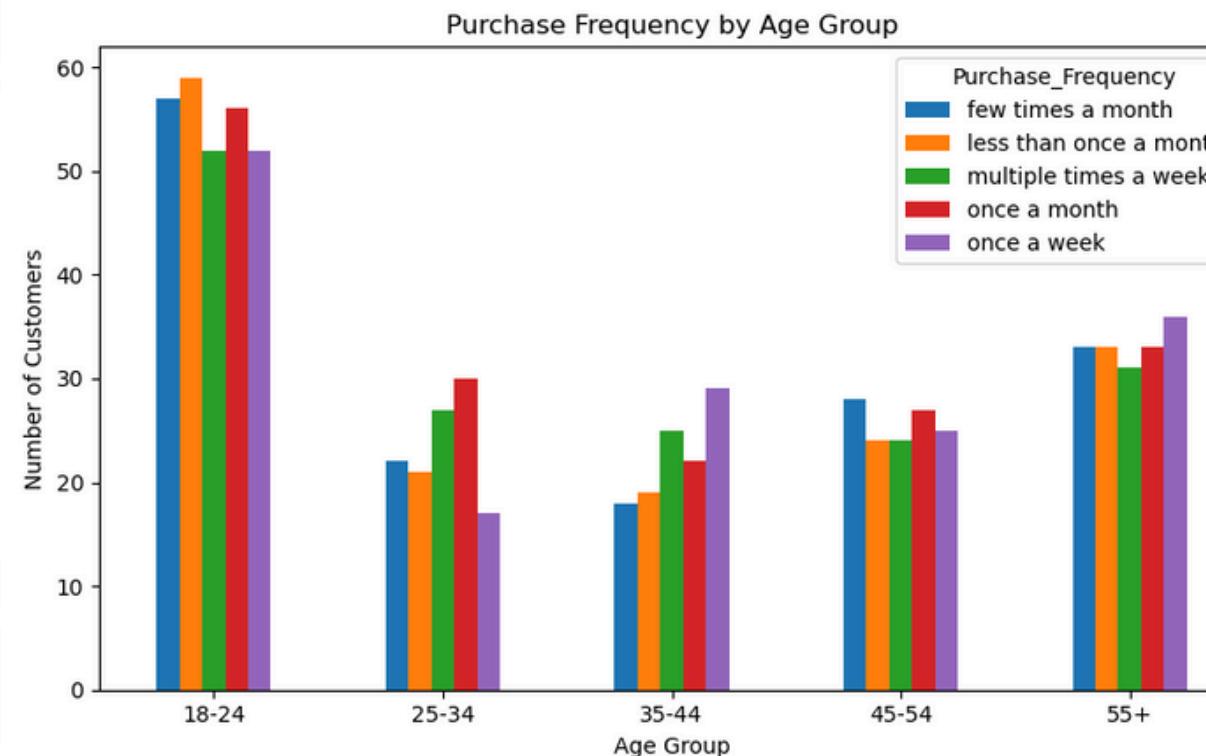
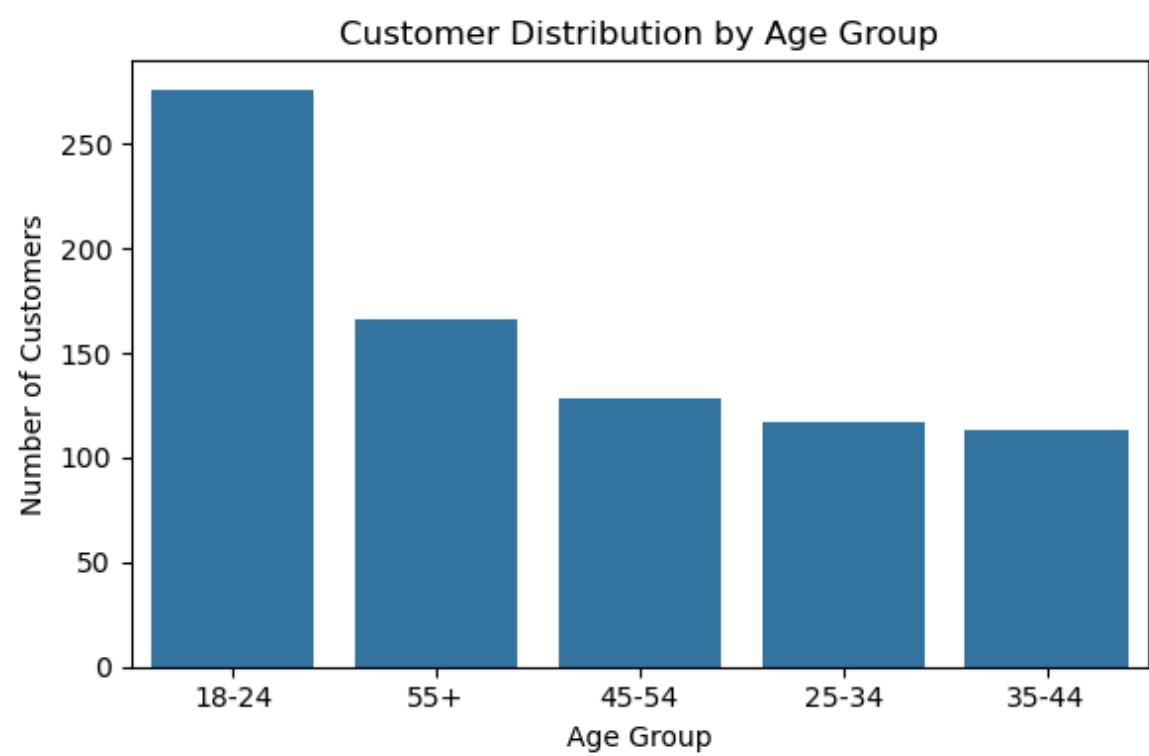
os.makedirs("../outputs", exist_ok=True)

df.to_csv("../outputs/cleaned_data.csv", index=False)
```

TASK 2: DESCRIPTIVE BEHAVIOR ANALYSIS

- SUMMARIZE CUSTOMER DEMOGRAPHICS (AGE, GENDER DISTRIBUTION).
- ANALYZE OVERALL PURCHASE FREQUENCY AND MOST POPULAR PRODUCT CATEGORIES.
- IDENTIFY TOP BROWSING METHODS AND MOST COMMON CART ABANDONMENT FACTORS.
- CALCULATE MEAN AND MEDIAN SATISFACTION, RECOMMENDATION HELPFULNESS, AND RATING ACCURACY.
- GENERATE SUMMARY STATISTICS AND VISUALIZATIONS FOR KEY BEHAVIORAL VARIABLES.

TASK 2: DESCRIPTIVE BEHAVIOR ANALYSIS



TASK 2: DESCRIPTIVE BEHAVIOR ANALYSIS

CODES PERFORMED IN PYTHON

```
# Task 2 : Descriptive Behavior Analysis
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import os

df = pd.read_csv("../outputs/cleaned_data_v2.csv")

df.head()
df.info()

df['age_group'].value_counts()

plt.figure(figsize=(6,4))
sns.countplot(
    x='age_group',
    data=df,
    order=df['age_group'].value_counts().index
)
plt.title("Customer Distribution by Age Group")
plt.xlabel("Age Group")
plt.ylabel("Number of Customers")
plt.tight_layout()
plt.savefig("../outputs/charts/age_group_distribution.png")
plt.show()

df['Gender'].value_counts()

age_purchase = pd.crosstab(df['age_group'], df['Purchase_Frequency'])

age_purchase.plot(kind='bar', figsize=(8,5))
plt.title("Purchase Frequency by Age Group")
```

```
age_purchase.plot(kind='bar', figsize=(8,5))
plt.title("Purchase Frequency by Age Group")
plt.xlabel("Age Group")
plt.ylabel("Number of Customers")
plt.xticks(rotation=0)
plt.tight_layout()
plt.savefig("../outputs/charts/age_vs_purchase_frequency.png")
plt.show()

from collections import Counter

categories = df['Purchase_Categories'].dropna().str.split(',')
flat_categories = [c.strip() for sub in categories for c in sub]

category_counts = Counter(flat_categories)

top_categories = pd.DataFrame(
    category_counts.most_common(10),
    columns=['Product_Category', 'Count']
)

top_categories

plt.figure(figsize=(7,4))
sns.barplot(
    x='Count',
    y='Product_Category',
    data=top_categories
)
plt.title("Most Popular Product Categories")
plt.tight_layout()
plt.savefig("../outputs/charts/most_popular_categories.png")
plt.show()

browsing_methods = df['Product_Search_Method'].value_counts()
browsing_methods
```

```
plt.figure(figsize=(6,4))
sns.countplot(
    y='Product_Search_Method',
    data=df,
    order=browsing_methods.index
)
plt.title("Top Product Browsing Methods")
plt.tight_layout()
plt.savefig("../outputs/charts/top_browsing_methods.png")
plt.show()

cart_abandonment = df['Cart_Abandonment_Factors'].value_counts()
cart_abandonment

plt.figure(figsize=(7,4))
sns.countplot(
    y='Cart_Abandonment_Factors',
    data=df,
    order=cart_abandonment.index
)
plt.title("Most Common Cart Abandonment Factors")
plt.tight_layout()
plt.savefig("../outputs/charts/cart_abandonment_factors.png")
plt.show()

metrics = df[
    ['Shopping_Satisfaction',
     'Recommendation_Helpfulness',
     'Rating_Accuracy']
]
```

```
plt.figure(figsize=(6,4))
sns.boxplot(data=metrics)
plt.title("Distribution of Key Behavioral Ratings")
plt.tight_layout()
plt.savefig("../outputs/charts/behavioral_metrics_distribution.png")
plt.show()
```

TASK 2:RECOMMENDATION

#ELECTRONICS AND FASHION EMERGE AS THE MOST FREQUENTLY PURCHASED CATEGORIES, INDICATING STRONG CROSS-SELLING POTENTIAL. #SEARCH-BASED BROWSING DOMINATES, REFLECTING HIGH INTENT-DRIVEN SHOPPING BEHAVIOR. #PRICING AND DELIVERY-RELATED ISSUES ARE THE LEADING CAUSES OF CART ABANDONMENT. #AVERAGE CUSTOMER SATISFACTION AND RATING ACCURACY REMAIN HIGH, REINFORCING TRUST IN THE PLATFORM.

#RECOMMENDATION HELPFULNESS SHOWS MODERATE VARIANCE, SUGGESTING SCOPE FOR PERSONALIZATION IMPROVEMENTS.

TASK 3: CUSTOMER SEGMENTATION AND PROFILING

- SEGMENT CUSTOMERS BASED ON PURCHASE FREQUENCY AND SHOPPING SATISFACTION LEVELS.
- CREATE PROFILES SUCH AS:
 - FREQUENT BUYERS: HIGH PURCHASE FREQUENCY, HIGH SATISFACTION.
 - OCCASIONAL SHOPPERS: MEDIUM FREQUENCY, MODERATE SATISFACTION.
 - AT-RISK CUSTOMERS: LOW SATISFACTION OR FREQUENT CART ABANDONMENT.
- ANALYZE DEMOGRAPHIC OR BEHAVIORAL DIFFERENCES ACROSS THESE SEGMENTS.
- USE CLUSTERING (E.G., K-MEANS) FOR BEHAVIORAL GROUPING BASED ON SURVEY RESPONSES.

TASK 3: CUSTOMER SEGMENTATION AND PROFILING

Timestamp	age	Gender	Purchase_Frequency	Purchase_Personalization	Browsing_Habits	Product_Selection	Search_Keyword	Customer_Role	Add_to_Cart	Cart_Compliance	Cart_Abandonment	Save_forLater	Review_Likelihood	Review_Rate	Review_Honesty	Recommendation_Acceptance	Rating_Avg	Shopping_Service				
2023/06/01	65	Other	less than once a week	Clothing and accessories	Yes	Multiple times	Multiple filters	Keyword	Multiple products	2	No	Sometimes	Found a better option	Often	Sometimes	Yes	Moderate	No	Sometimes	4	4	Competitor
2023/06/01	20	Male	once a week	Groceries	No	Rarely	Filter	First page	3	No	Never	High shipping	Always	Yes	Heavily	Yes	Sometimes	4	5	Quick		
2023/06/01	42	Male	once a week	Groceries	Sometimes	Few times	Keyword	Multiple products	2	No	Rarely	Found a better option	Often	Yes	Heavily	Sometimes	No	5	3	All the time		
2023/06/01	65	Other	once a month	Beauty and personal care	No	Few times	Filter	Multiple products	2	Yes	Sometimes	Others recommended	Often	Yes	Occasionally	No	Yes	1	2	Quick		
2023/06/01	45	Female	once a week	Beauty and personal care	Sometimes	Few times	Unknown	First page	5	Maybe	Rarely	Changed recommendation	Never	Yes	Rarely	No	Yes	1	2	Quick		
2023/06/01	67	Other	few times	Groceries	No	Multiple times	Unknown	Multiple products	3	No	Always	others recommended	Never	No	Moderate	Sometimes	No	4	3	Product		
2023/06/01	55	Female	few times	Groceries	Sometimes	Few times	Keyword	First page	5	Maybe	Never	Found a better option	Rarely	Yes	Moderate	No	Yes	5	3	Customer		
2023/06/01	44	Other	few times	Beauty and personal care	Yes	Few times	Unknown	Multiple products	4	Yes	Often	Changed recommendation	Rarely	No	Heavily	Sometimes	Sometimes	5	1	Customer		
2023/06/01	33	Male	less than once a month	Beauty and personal care	Yes	Rarely	Unknown	Multiple products	1	Yes	Often	High shipping	Rarely	No	Heavily	Sometimes	Sometimes	2	2	Competitor		
2023/06/01	22	Other	few times	Clothing and accessories	No	Multiple times	Filter	First page	5	Maybe	Never	Changed recommendation	Rarely	No	Occasionally	No	Yes	1	3	All the time		
2023/06/01	28	Female	once a month	Groceries	Yes	Few times	others recommended	First page	4	Yes	Sometimes	Changed recommendation	Sometimes	Yes	Occasionally	Yes	Yes	4	5	All the time		
2023/06/01	44	Other	multiple times	Groceries	Sometimes	Multiple times	Keyword	First page	5	Yes	Always	others recommended	Often	Yes	Occasionally	Yes	Yes	4	3	User-friendly		
2023/06/01	24	Male	once a month	Beauty and personal care	No	Few times	Filter	First page	4	No	Always	others recommended	Rarely	Yes	Moderate	Yes	Yes	5	3	Competitor		
2023/06/01	33	Male	once a month	Beauty and personal care	No	Few times	Filter	First page	5	Maybe	Always	Changed recommendation	Sometimes	No	Moderate	Yes	Sometimes	5	1	Wide		

T	U	V	W	X	Y	Z	AA	AB	AC	AD
: Shopping_Service_A	Improvement_Area	transaction	age_group	Satisfaction_Level	FreqL	Customer_Segment	Recommen	Recommen	Behavior_Cluster	
4	Competititve price	better app interface	778242	55+	High		At-Risk Customers	2	2	1
5	Quick delivery	scrolling option would be better	193482	18-24	High		Medi	At-Risk Customers	2	2
3	All the above	Unknown	925975	35-44	Medium		Medi	At-Risk Customers	1	1
2	Quick delivery	quality of product is good	566872	55+	Low		Medi	At-Risk Customers	3	3
2	Quick delivery	irrelevant products shown	683642	45-54	Low		Medi	At-Risk Customers	3	3

TASK 3: CUSTOMER SEGMENTATION AND PROFILING

CODES PERFORMED IN PYTHON

```
#TASK 3: CUSTOMER SEGMENTATION & PROFILING
# Satisfaction Level
df['Satisfaction_Level'] = pd.cut(
    df['Shopping_Satisfaction'],
    bins=[0, 2, 3, 5],
    labels=['Low', 'Medium', 'High']
)

df['Satisfaction_Level'].value_counts()

df['Frequency_Level'] = df['Purchase_Frequency'].map({
    'daily': 'High',
    'few times a week': 'High',
    'once a week': 'Medium',
    'once a month': 'Medium',
    'rarely': 'Low'
})

def assign_segment(row):
    if row['Frequency_Level'] == 'High' and row['Satisfaction_Level'] == 'High':
        return 'Frequent Buyers'
    elif row['Satisfaction_Level'] == 'Low' or row['Cart_Abandonment_Factors'] != 'None':
        return 'At-Risk Customers'
    else:
        return 'Occasional Shoppers'

df['Customer_Segment'] = df.apply(assign_segment, axis=1)

df['Customer_Segment'].value_counts()

pd.crosstab(df['Customer_Segment'], df['age_group'])

pd.crosstab(df['Customer_Segment'], df['Gender'])

df[['Shopping_Satisfaction', 'Rating_Accuracy', 'Recommendation_Helpfulness']].dtypes

df['Recommendation_Helpfulness_Num'] = df['Recommendation_Helpfulness'].map({
    'No': 1,
    'Sometimes': 2,
    'Yes': 3
})
```

```
'Sometimes': 2,
'Yes': 3
])

df['Shopping_Satisfaction'] = pd.to_numeric(df['Shopping_Satisfaction'], errors='coerce')
df['Rating_Accuracy'] = pd.to_numeric(df['Rating_Accuracy'], errors='coerce')
df['Recommendation_Helpfulness_Num'] = pd.to_numeric(
    df['Recommendation_Helpfulness'], errors='coerce'
)

df.groupby('Customer_Segment')[[
    'Shopping_Satisfaction',
    'Rating_Accuracy',
    'Recommendation_Helpfulness_Num'
]].mean()

df.rename(columns={
    'Recommendation_Helpfulness_Num': 'Recommendation_Helpfulness_Score'
}, inplace=True)

df.columns

df['Recommendation_Helpfulness_Num'] = df['Recommendation_Helpfulness'].map({
    'No': 1,
    'Sometimes': 2,
    'Yes': 3
})

df[['Recommendation_Helpfulness', 'Recommendation_Helpfulness_Num']].head()

df['Recommendation_Helpfulness_Num'].isnull().sum()

cluster_features = df[
    ['Shopping_Satisfaction',
     'Rating_Accuracy',
     'Recommendation_Helpfulness_Num']
].dropna()
```

```
cluster_features = df[
    ['Shopping_Satisfaction',
     'Rating_Accuracy',
     'Recommendation_Helpfulness_Num']
].dropna()

from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans

scaler = StandardScaler()
scaled_features = scaler.fit_transform(cluster_features)

kmeans = KMeans(n_clusters=3, random_state=42)
cluster_features['Cluster'] = kmeans.fit_predict(scaled_features)

df.loc[cluster_features.index, 'Behavior_Cluster'] = cluster_features['Cluster']

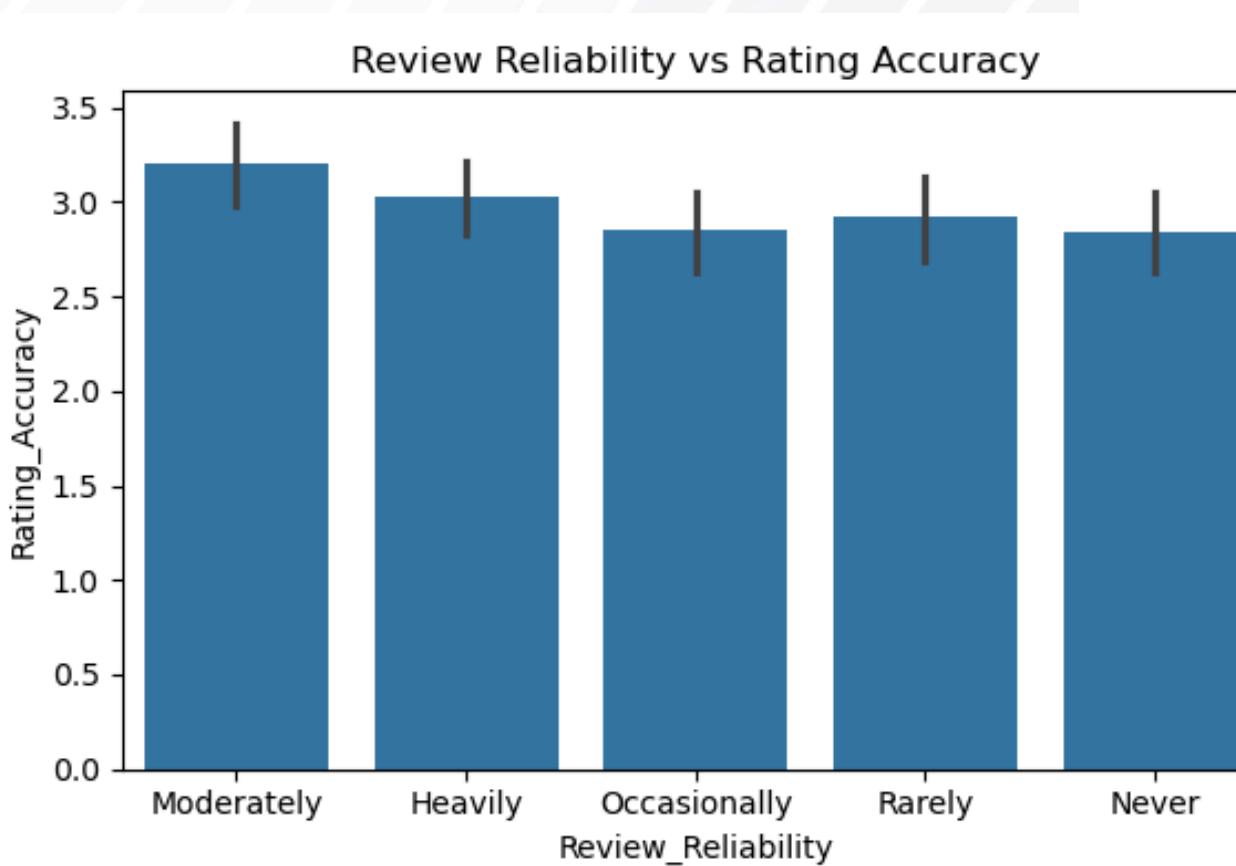
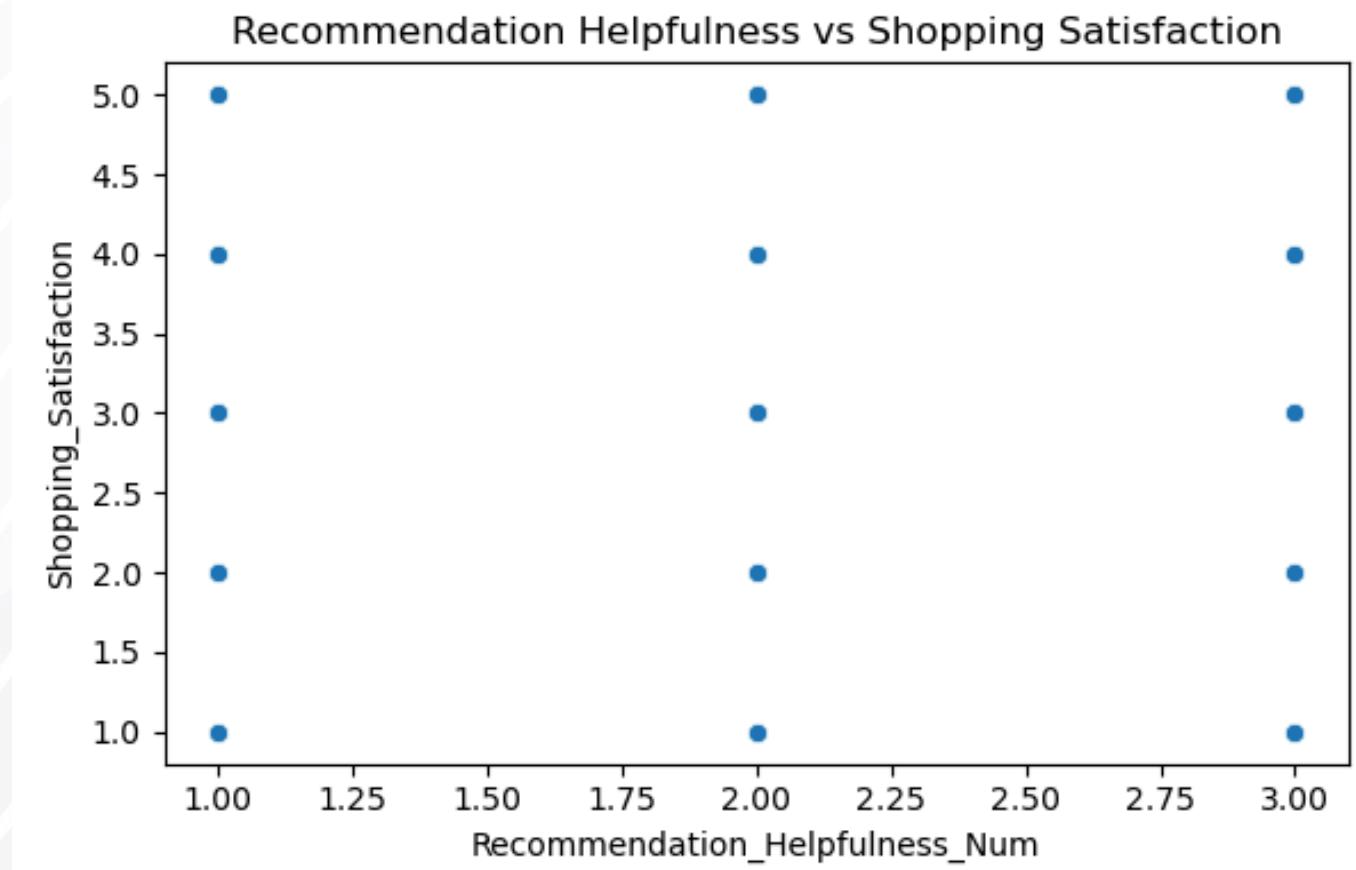
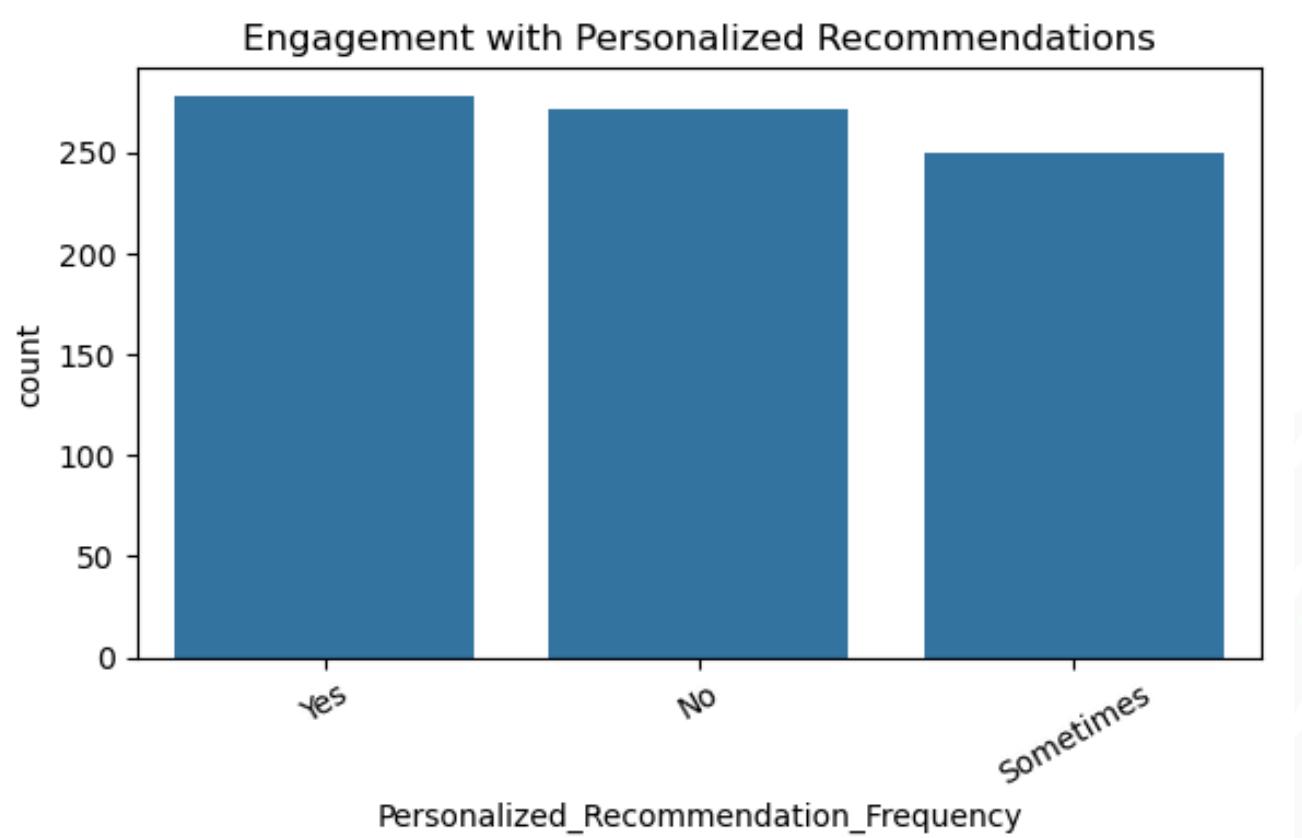
df.groupby('Behavior_Cluster')[[
    'Shopping_Satisfaction',
    'Rating_Accuracy',
    'Recommendation_Helpfulness_Num'
]].mean()

df.to_csv("../outputs/task3.csv", index=False)
```

TASK 4: RECOMMENDATION AND REVIEW INSIGHTS

● EXAMINE THE RELATIONSHIP BETWEEN RECOMMENDATION HELPFULNESS AND SHOPPING SATISFACTION. ● EVALUATE HOW REVIEW RELIABILITY AND HELPFULNESS IMPACT OVERALL RATINGS. ● IDENTIFY TRENDS IN HOW OFTEN CUSTOMERS ENGAGE WITH OR TRUST PERSONALIZED RECOMMENDATIONS. ● SUGGEST ACTIONABLE INSIGHTS FOR IMPROVING AMAZON'S RECOMMENDATION SYSTEM.

TASK 4: RECOMMENDATION AND REVIEW INSIGHTS



TASK 4: RECOMMENDATION AND REVIEW INSIGHTS

Timestamp	age	Gender	Purchase_Quantity	Purchase_Period	Personalized_Browsing	Product_Selection	Search_Keyword	Customer_Status	Add_to_Cart	Cart_Count	Cart_Abandonment	Save_for_later	Review_Likelihood	Review_Rate	Review_Honesty	Recommendation	Rating_Avg	Shopping_Service	
2023/06/01	65	Other	less than once a week	Clothing	Yes	Multiple times	Keyword	Multiple pages	2	No	Sometimes	Found a lot	Sometimes	Yes	Moderate	No	Sometimes	4	4 Comp
2023/06/01	20	Male	once a month	Groceries	No	Rarely	Filter	First page	3	No	Never	High shipping	Always	Yes	Heavily	Yes	Sometimes	4	5 Quick
2023/06/01	42	Male	once a month	Groceries	Sometimes	Few times	Keyword	Multiple pages	2	No	Rarely	Found a lot	Often	Yes	Heavily	Sometimes	No	5	All the time
2023/06/01	65	Other	once a month	Beauty & Health	No	Few times	Filter	Multiple pages	2	Yes	Sometimes	others	Often	Yes	Occasionally	No	Yes	1	2 Quick
2023/06/01	45	Female	once a month	Beauty & Health	Sometimes	Few times	Unknown	First page	5	Maybe	Rarely	Changed or never	Never	Yes	Rarely	No	Yes	1	2 Quick
2023/06/01	67	Other	few times	Groceries	No	Multiple times	Unknown	Multiple pages	3	No	Always	others	Never	No	Moderate	Sometimes	No	4	3 Product
2023/06/01	55	Female	few times	Groceries	Sometimes	Few times	Keyword	First page	5	Maybe	Never	Found a lot	Rarely	Yes	Moderate	No	Yes	5	3 Customer
2023/06/01	44	Other	few times	Beauty & Health	Yes	Few times	Unknown	Multiple pages	4	Yes	Often	Changed or rarely	Rarely	No	Heavily	Sometimes	Sometimes	5	1 Customer
2023/06/01	33	Male	less than once a month	Beauty & Health	Yes	Rarely	Unknown	Multiple pages	1	Yes	Often	High shipping	Rarely	No	Heavily	Sometimes	Sometimes	2	2 Comp
2023/06/01	22	Other	few times	Clothing	No	Multiple times	Filter	First page	5	Maybe	Never	Changed or rarely	Rarely	No	Occasionally	No	Yes	1	3 All the time
2023/06/01	28	Female	once a month	Groceries	Yes	Few times	others	First page	4	Yes	Sometimes	Changed or	Sometimes	Yes	Occasionally	Yes	Yes	4	5 All the time
2023/06/01	44	Other	multiple times	Groceries	Sometimes	Multiple times	Keyword	First page	5	Yes	Always	others	Often	Yes	Occasionally	Yes	Yes	4	3 User-friendly
2023/06/01	24	Male	once a month	Beauty & Health	No	Few times	Filter	First page	4	No	Always	others	Rarely	Yes	Moderate	Yes	Yes	5	3 Comp
2023/06/01	23	Male	once a month	Beauty & Health	Yes	Few times	Filter	First page	5	Maybe	Always	Changed or	Sometimes	No	Moderate	Yes	Sometimes	5	1 Wide
2023/06/01	58	Other	once a month	Beauty & Health	Sometimes	Few times	Keyword	First page	5	Maybe	Never	High shipping	Often	Yes	Never	Yes	No	3	1 Quick
2023/06/01	40	Male	less than once a month	Groceries	No	Multiple times	Keyword	First page	3	No	Never	Found a lot	Never	No	Never	No	No	2	2 Customer
2023/06/01	46	Other	once a month	Groceries	Yes	Rarely	categories	First page	2	Yes	Always	High shipping	Sometimes	No	Rarely	Sometimes	Sometimes	5	2 Wide
2023/06/01	15	Other	once a month	Groceries	Sometimes	Few times	Keyword	Multiple pages	2	Maybe	Often	others	Often	Yes	Never	No	Yes	2	2 Quick

Service_A	Improvement	transaction	age_group	Satisfaction	Frequency	Customer_Status	Recommendation	Recommendation	Behavior_Cluster
Competititve pricing	better app	778242	55+	High		At-Risk Customer	2	2	1
Quick delivery	scrolling	193482	18-24	High	Medium	At-Risk Customer	2	2	0
All the above	Unknown	925975	35-44	Medium	Medium	At-Risk Customer	1	1	1
Quick delivery	quality of products	566872	55+	Low	Medium	At-Risk Customer	3	3	2
Quick delivery	irrelevant products	683642	45-54	Low	Medium	At-Risk Customer	3	3	2
Product range	nothing new	461923	55+	Medium		At-Risk Customer	1	1	1

TASK 4: RECOMMENDATION AND REVIEW INSIGHTS

CODES PERFORMED IN PYTHON

```
#TASK 4: RECOMMENDATION & REVIEW INSIGHTS
df[['Recommendation_Helpfulness_Num', 'Shopping_Satisfaction']].corr()

import os

os.makedirs("../outputs/charts_task_4", exist_ok=True)

import seaborn as sns
import matplotlib.pyplot as plt

plt.figure(figsize=(6,4))
sns.scatterplot(
    x='Recommendation_Helpfulness_Num',
    y='Shopping_Satisfaction',
    data=df
)
plt.title("Recommendation Helpfulness vs Shopping Satisfaction")
plt.tight_layout()
plt.savefig("../outputs/charts_task_4/recommendation_helpfulness_vs_satisfaction.png")
plt.show()

df['Review_Reliability'].value_counts()

review_rating = df.groupby('Review_Reliability')['Rating_Accuracy'].mean()
review_rating

import os
os.makedirs("../outputs/charts_task_4", exist_ok=True)

plt.figure(figsize=(6,4))
sns.barplot(
    x='Review_Reliability',
    y='Rating_Accuracy',
    data=df,
    estimator='mean'
)
```

```
estimator='mean'
)
plt.title("Review Reliability vs Rating Accuracy")
plt.tight_layout()
plt.savefig("../outputs/charts_task_4/review_reliability_vs_rating_accuracy.png")
plt.show()

df['Personalized_Recommendation_Frequency'].value_counts()

plt.figure(figsize=(6,4))
sns.countplot(
    x='Personalized_Recommendation_Frequency',
    data=df,
    order=df['Personalized_Recommendation_Frequency'].value_counts().index
)
plt.title("Engagement with Personalized Recommendations")
plt.xticks(rotation=30)
plt.tight_layout()
plt.savefig("../outputs/charts_task_4/recommendation_engagement_frequency.png")
plt.show()

df.groupby('Personalized_Recommendation_Frequency')['Shopping_Satisfaction'].mean()

df.groupby('Customer_Segment')['Recommendation_Helpfulness_Num'].mean()

df.to_csv("../outputs/task_4.csv", index=False)
```

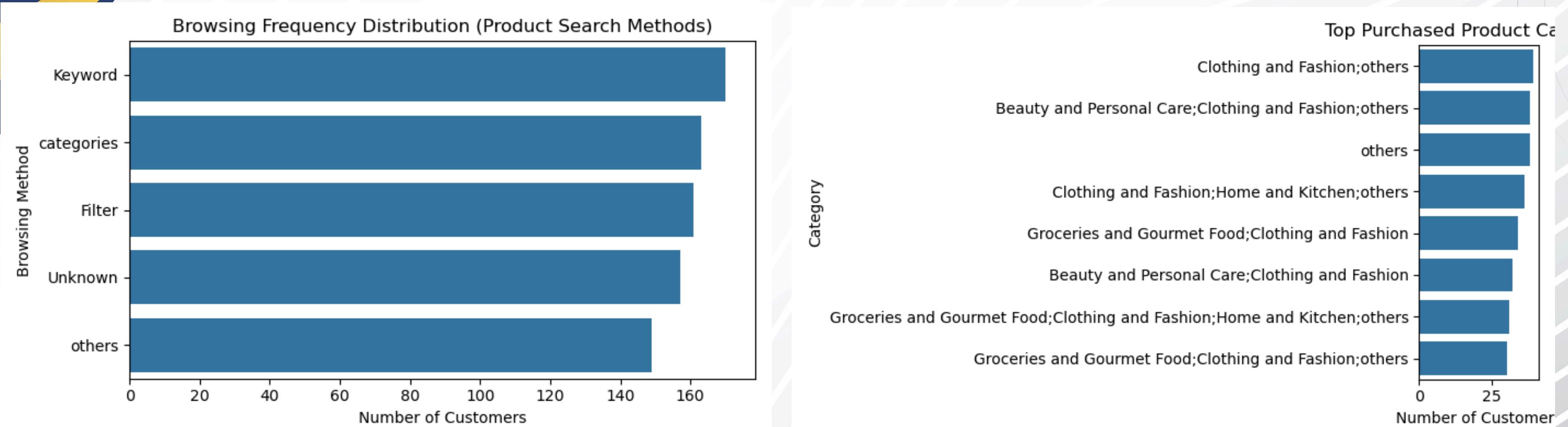
TASK 4: RECOMMENDATION

#IMPROVE RECOMMENDATION RELEVANCE FOR AT-RISK CUSTOMERS TO REDUCE CHURN #OPTIMIZE RECOMMENDATION FREQUENCY TO AVOID USER FATIGUE #PRIORITIZE VERIFIED AND RELIABLE REVIEWS IN RECOMMENDATION RANKING #USE BEHAVIORAL CLUSTERS TO PERSONALIZE RECOMMENDATIONS #ENHANCE UI/UX TO MAKE RECOMMENDATIONS CLEARER AND LESS CLUTTERED

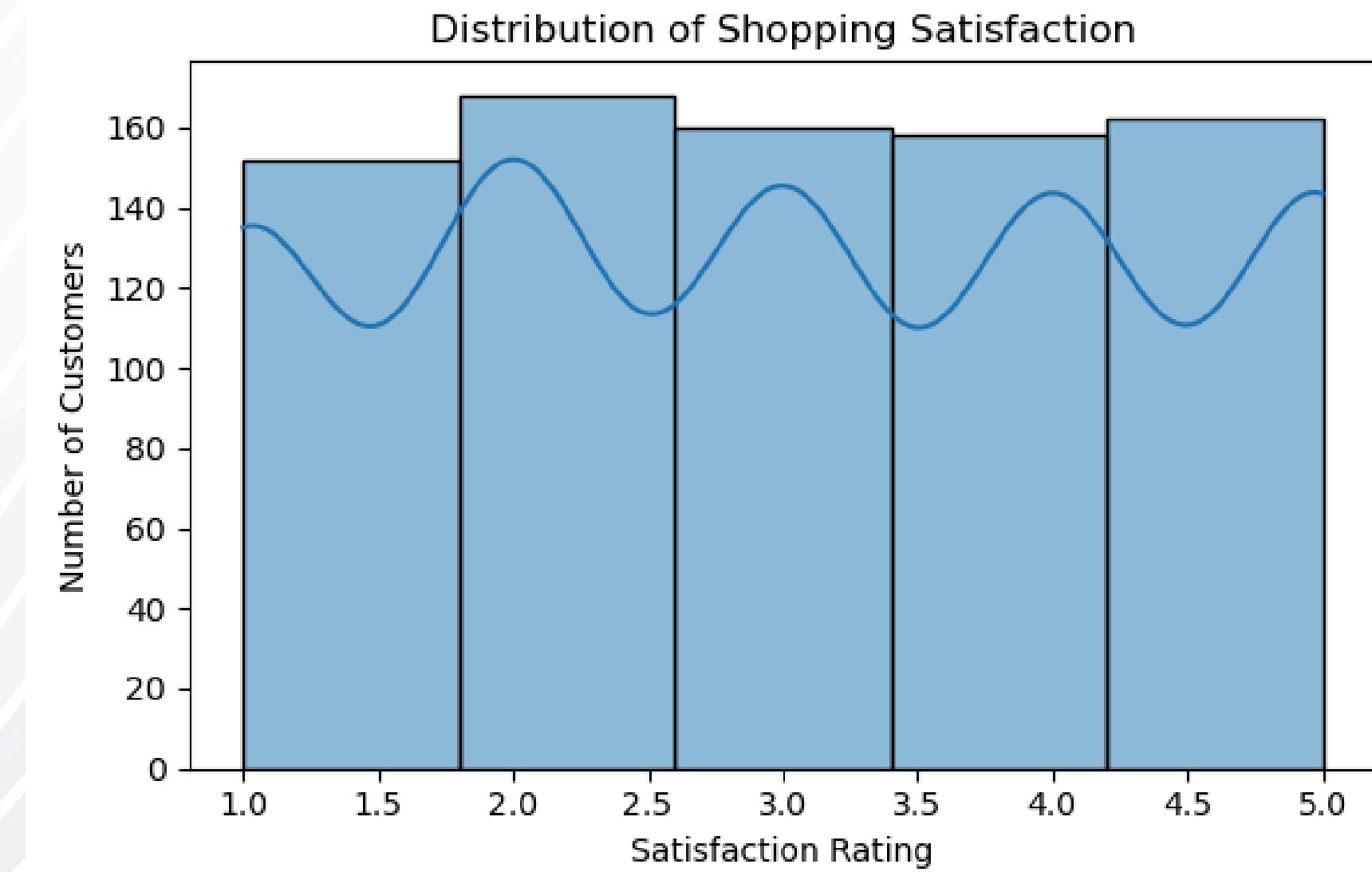
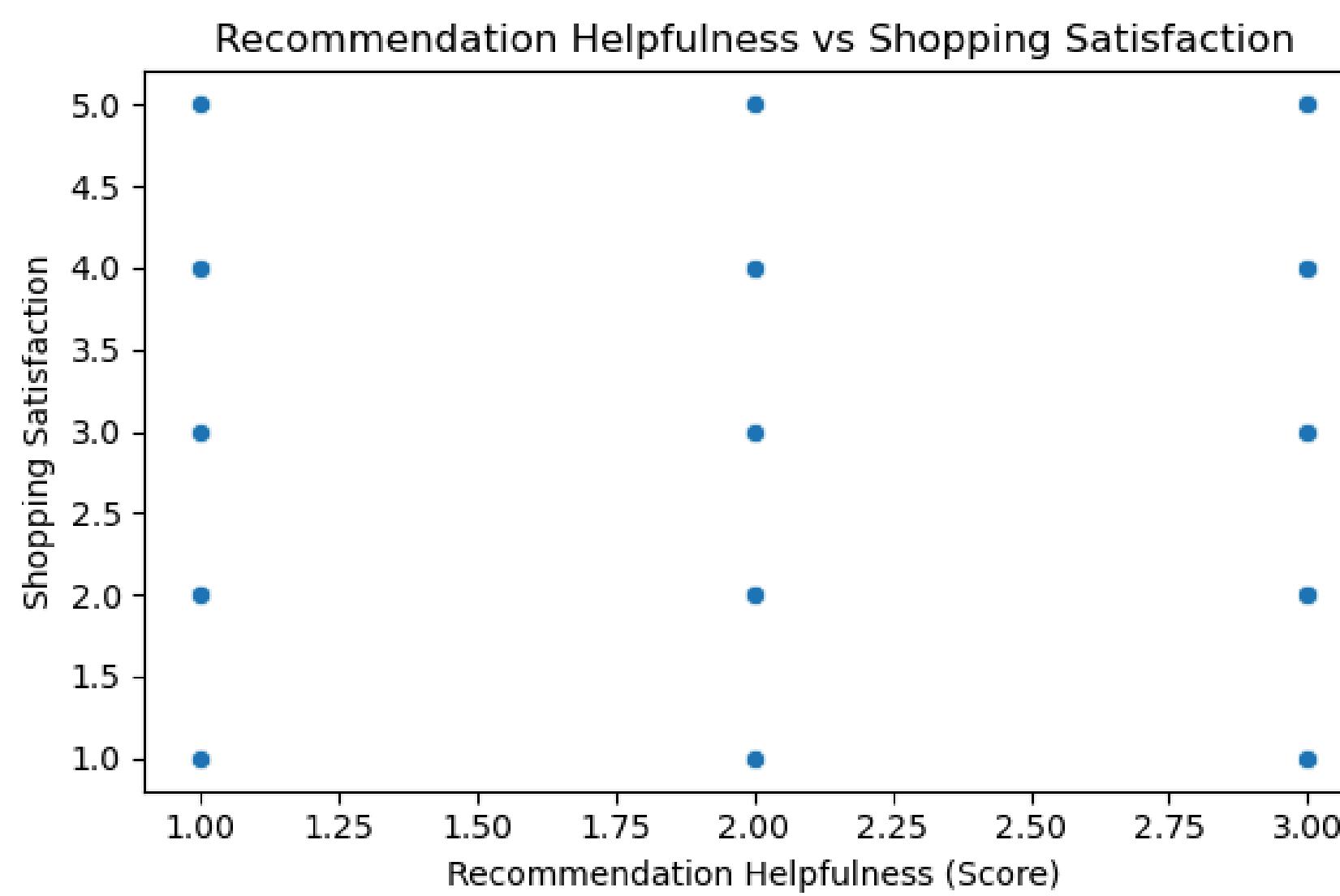
TASK 5: VISUALIZATION AND REPORTING

- CREATE ATTRACTIVE VISUALIZATIONS (BAR CHARTS, HEATMAPS, PIE CHARTS)
FOR:
 - PURCHASE CATEGORIES
 - BROWSING FREQUENCY DISTRIBUTION
 - SATISFACTION LEVELS
 - CORRELATION BETWEEN RECOMMENDATION USEFULNESS AND SATISFACTION
- SUMMARIZE FINDINGS IN A CLEAR AND VISUALLY APPEALING DASHBOARD OR REPORT FORMAT.

TASK 5: VISUALISATION



TASK 5: VISUALISATION



TASK 5: VISUALISATION

CODES PERFORMED IN PYTHON

```
#task 5 : Visualization and Reporting
import os
os.makedirs("../outputs/charts_task_5", exist_ok=True)

from collections import Counter
import seaborn as sns
import matplotlib.pyplot as plt

categories = df['Purchase_Categories'].dropna().str.split(',')
flat_categories = [c.strip() for sub in categories for c in sub]

category_counts = Counter(flat_categories)
top_categories = dict(category_counts.most_common(8))

plt.figure(figsize=(7,4))
sns.barplot(x=list(top_categories.values()), y=list(top_categories.keys()))
plt.title("Top Purchased Product Categories")
plt.xlabel("Number of Customers")
plt.ylabel("Category")
plt.tight_layout()
plt.savefig("../outputs/charts_task_5/purchase_categories.png")
plt.show()

df['Product_Search_Method'].value_counts()

import seaborn as sns
import matplotlib.pyplot as plt
import os

os.makedirs("../outputs/charts_task_5", exist_ok=True)
plt.figure(figsize=(7,4))
sns.countplot(
    y='Product_Search_Method',
    data=df,
    order=df['Product_Search_Method'].value_counts().index
)
```

```
df['Product_Search_Method'].value_counts()

import seaborn as sns
import matplotlib.pyplot as plt
import os

os.makedirs("../outputs/charts_task_5", exist_ok=True)
plt.figure(figsize=(7,4))
sns.countplot(
    y='Product_Search_Method',
    data=df,
    order=df['Product_Search_Method'].value_counts().index
)

plt.title("Browsing Frequency Distribution (Product Search Methods)")
plt.xlabel("Number of Customers")
plt.ylabel("Browsing Method")
plt.tight_layout()
plt.savefig("../outputs/charts_task_5/browsing_frequency_distribution.png")
plt.show()

import os
os.listdir("../outputs/charts_task_5")

df['Shopping_Satisfaction'].value_counts().sort_index()

import seaborn as sns
import matplotlib.pyplot as plt
import os

os.makedirs("../outputs/charts_task_5", exist_ok=True)
plt.figure(figsize=(6,4))
sns.histplot(
    df['Shopping_Satisfaction'],
    bins=5,
    kde=True
)
```

```
df['Shopping_Satisfaction'].value_counts().sort_index()

import seaborn as sns
import matplotlib.pyplot as plt
import os

os.makedirs("../outputs/charts_task_5", exist_ok=True)
plt.figure(figsize=(6,4))
sns.histplot(
    df['Shopping_Satisfaction'],
    bins=5,
    kde=True
)

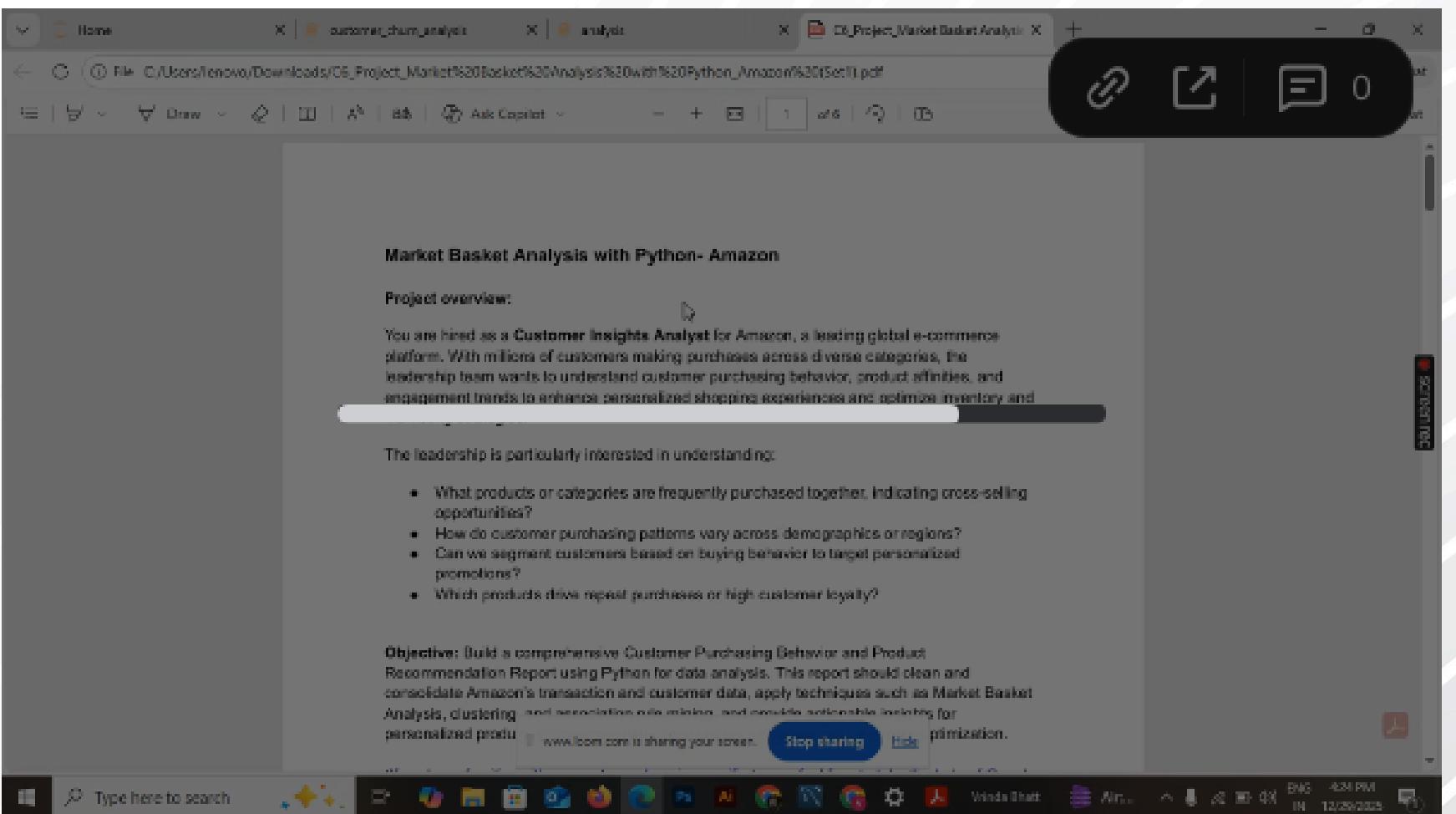
plt.title("Distribution of Shopping Satisfaction")
plt.xlabel("Satisfaction Rating")
plt.ylabel("Number of Customers")
plt.tight_layout()
plt.savefig("../outputs/charts_task_5/satisfaction_levels.png")
plt.show()
```

TASK 5: INSIGHTS

#MOST CUSTOMERS REPORT MEDIUM TO HIGH
SATISFACTION LEVELS
#INDICATES OVERALL POSITIVE SHOPPING
EXPERIENCE WITH SCOPE FOR IMPROVEMENT
AMONG LOW-SATISFACTION USERS

THANK YOU

VIDEO LINK:



VIDEO LINK