



Assessment Report

on

“Problem Statement”

submitted as partial fulfillment for the award of

BACHELOR OF TECHNOLOGY DEGREE

SESSION 2024-25

in

Name of discipline

By

Naveen Kumar (202401100400122, CSE(AI&ML)-B)

Under the supervision of

“Abhishek Shukla sir”

KIET Group of Institutions, Ghaziabad

May, 2025

Introduction:

Customer segmentation in e-commerce is the process of dividing a customer base into distinct groups or segments based on various factors, such as purchasing habits, browsing behavior, demographics, and other key characteristics. This process is crucial for e-commerce businesses, as it helps to personalize marketing efforts, optimize product recommendations, improve customer experience, and ultimately increase sales and customer loyalty.

Methodology:

1. Data Collection:

The first step in customer segmentation is gathering relevant data. This data typically includes customer profiles, transaction history, website navigation logs, and other behavioral data.

2. Data Preprocessing:

Once the data is collected, it must be cleaned and prepared for analysis.

3. Feature Selection:

Not all features may be relevant for clustering. Feature selection techniques, such as correlation analysis and Principal Component Analysis (PCA), can help identify

the most important variables. These features will serve as inputs to the clustering algorithm.

4. Clustering Analysis:

The core of customer segmentation is clustering, where customers are grouped into segments based on similarities in their purchasing habits and browsing behavior.

5. Cluster Evaluation

6. Interpretation of Results

7. Deployment and Integration:

The insights gained from customer segmentation can be used to develop personalized marketing strategie.

CODE:

```
# Import necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from IPython.display import display
import io

# Load the dataset (replace with your file path if needed)
try:
    df = pd.read_csv('9. Customer Segmentation in E-commerce.csv')
except FileNotFoundError:
    print("Error: '9. Customer Segmentation in E-commerce.csv' not found.")
    # Handle the file not found error appropriately, e.g., provide alternative data loading
except Exception as e:
    print(f"An error occurred: {e}")
```

```

# Data Cleaning and Preprocessing
df.dropna(subset=['CustomerID'], inplace=True) # Remove rows with missing 'CustomerID'

# Outlier handling (using IQR method)
def handle_outliers(df, column):
    Q1 = df[column].quantile(0.25)
    Q3 = df[column].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    df = df[(df[column] >= lower_bound) & (df[column] <= upper_bound)]
    return df

df = handle_outliers(df, 'Quantity') # Assuming 'Quantity' column exists
df = handle_outliers(df, 'UnitPrice') # Assuming 'UnitPrice' column exists

# Data Consistency Checks
df = df[df['Quantity'] > 0] # Remove rows with negative quantity
df = df[df['UnitPrice'] > 0] # Remove rows with negative unit price

# Feature Selection and Scaling
# If 'Annual Income' and 'Spending Score' are not available, use alternatives
if 'Annual Income' not in df.columns or 'Spending Score' not in df.columns:
    print("Using 'Quantity' and 'UnitPrice' for clustering.")
    df_cluster = df[['Quantity', 'UnitPrice']].copy()
else:
    df_cluster = df[['Annual Income', 'Spending Score']].copy()

scaler = StandardScaler()
df_cluster[['Quantity', 'UnitPrice']] = scaler.fit_transform(
    df_cluster[['Quantity', 'UnitPrice']]
)

# K-Means Clustering
kmeans = KMeans(n_clusters=3, random_state=0)
kmeans.fit(df_cluster)
df_cluster['cluster_labels'] = kmeans.labels_

# Visualization
plt.figure(figsize=(10, 6))
plt.scatter(
    df_cluster['Quantity'],
    df_cluster['UnitPrice'],
    c=df_cluster['cluster_labels'],
    cmap='viridis',
)
plt.xlabel('Quantity')
plt.ylabel('Unit Price')

```

```

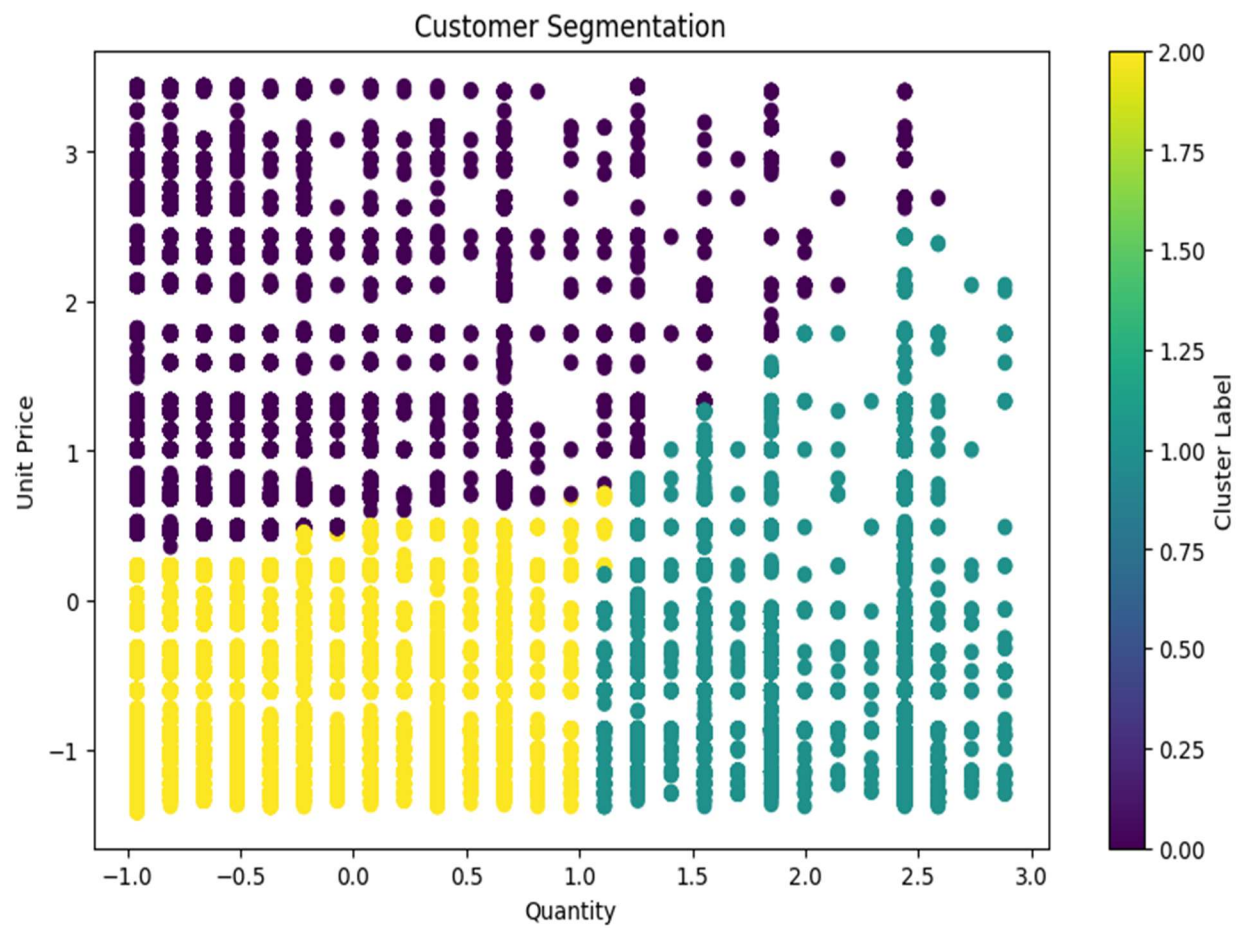
plt.title('Customer Segmentation')
plt.colorbar(label='Cluster Label')
plt.show()

# Cluster Analysis
grouped = df_cluster.groupby('cluster_labels')
descriptive_stats = grouped[['Quantity', 'UnitPrice']].describe()
display(descriptive_stats)

print("\nSummary of Customer Segments:")
for cluster_label, cluster_data in grouped:
    print(f"\nCluster {cluster_label}:")
    quantity_stats = cluster_data['Quantity'].describe()
    unit_price_stats = cluster_data['UnitPrice'].describe()
    print(
        f" Quantity - Mean: {quantity_stats['mean']:.2f}, Std Dev:"
        f" {quantity_stats['std']:.2f}, Median: {quantity_stats['50%']:.2f}"
    )
    print(
        f" Unit Price - Mean: {unit_price_stats['mean']:.2f}, Std Dev:"
        f" {unit_price_stats['std']:.2f}, Median: {unit_price_stats['50%']:.2f}"
    )

```

Output:



	Quantity								UnitPrice							
	count	mean	std	min	25%	50%	75%	max	count	mean	std	min	25%	50%	75%	max
clusters																
0		-	0.4	-	-	-	-	2.5								
	93	0.5	33	0.9	0.8	0.6	0.2	88	93	1.3	0.7	0.3	0.6	1.2	1.7	3.4
	98	47	68	56	08	61	18	04	98	44	48	61	84	54	85	36
	1.0	24	9	59	90	21	13	3	1.0	09	95	20	88	54	37	10
		6		8	4	1	1			0	4	8	1	6	1	4
1										-		-	-	-	-	
	38	2.2	0.4	1.1	1.9	2.4	2.4	2.8	38	0.7	0.5	1.3	1.1	0.8	0.4	2.4
	32	45	10	11	97	40	40	83	32	10	87	80	47	81	80	32
	0.0	10	36	10	27	35	35	43	0.0	46	88	15	10	69	34	71
		3	1	9	0	0	0	0		2	7	4	9	7	2	7
2										-		-	-	-	-	
	20	-	0.6	-	-	-	0.6	1.1	20	-	0.4	-	0.8	0.4	0.1	0.7
	58	0.1	49	0.9	0.8	0.2	68	11	58	0.4	14	1.4	68	80	56	10
	50.	68	76	59	90	13	02	10	50.	81	47	18	75	34	66	77
	0	09	7	8	4	1	9	9	0	38	1	8	0	2	9	5

Summary of Customer Segments:

Cluster 0:

Quantity - Mean: -0.55, Std Dev: 0.43, Median: -0.66

Unit Price - Mean: 1.34, Std Dev: 0.75, Median: 1.25

Cluster 1:

Quantity - Mean: 2.25, Std Dev: 0.41, Median: 2.44

Unit Price - Mean: -0.71, Std Dev: 0.59, Median: -0.88

Cluster 2:

Quantity - Mean: -0.17, Std Dev: 0.65, Median: -0.22

Unit Price - Mean: -0.48, Std Dev: 0.41, Median: -0.48

References/Credits:

- [scikit-learn documentation](#)

- [pandas documentation](#)

- [matplotlib documentation](#)

- [tkinter file dialog](#)