

```
!pip install pandas matplotlib seaborn scikit-learn
```

```
Requirement already satisfied: pandas in /usr/local/lib/python3.11/dist-packages (2.2.2)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.11/dist-packages (3.10.0)
Requirement already satisfied: seaborn in /usr/local/lib/python3.11/dist-packages (0.13.2)
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.11/dist-packages (1.6.1)
Requirement already satisfied: numpy>=1.23.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (1.26.4)
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.11/dist-packages (from pandas) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.11/dist-packages (from pandas) (2024.2)
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.11/dist-packages (from pandas) (2025.1)
Requirement already satisfied: contourpy>=1.0.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.3.1)
Requirement already satisfied: cycler>=0.10 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (4.55.5)
Requirement already satisfied: kiwisolver>=1.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (1.4.8)
Requirement already satisfied: packaging>=20.0 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=8 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (11.1.0)
Requirement already satisfied: pyparsing>=2.3.1 in /usr/local/lib/python3.11/dist-packages (from matplotlib) (3.2.1)
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.13.1)
Requirement already satisfied: joblib>=1.2.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in /usr/local/lib/python3.11/dist-packages (from scikit-learn) (3.5.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.11/dist-packages (from python-dateutil>=2.8.2->pandas) (1.17.0)
```

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
customers = pd.read_csv('Customers.csv')
products = pd.read_csv('Products.csv')
transactions = pd.read_csv('/content/Transaction.csv')

print(customers.head(), customers.info())
print(products.head(), products.info())
print(transactions.head(), transactions.info())

customers['SignupDate'] = pd.to_datetime(customers['SignupDate'])
transactions['TransactionDate'] = pd.to_datetime(transactions['TransactionDate'])

merged = transactions.merge(customers, on='CustomerID').merge(products, on='ProductID')

top_products = merged.groupby('ProductName')['Quantity'].sum().sort_values(ascending=False).head(5)
print("Top 5 Selling Products:\n", top_products)

sales_by_date = merged.groupby('TransactionDate')['TotalValue'].sum()
plt.figure(figsize=(10, 6))
plt.plot(sales_by_date.index, sales_by_date.values)
plt.title('Total Sales Over Time')
plt.xlabel('Date')
plt.ylabel('Total Sales')
plt.show()

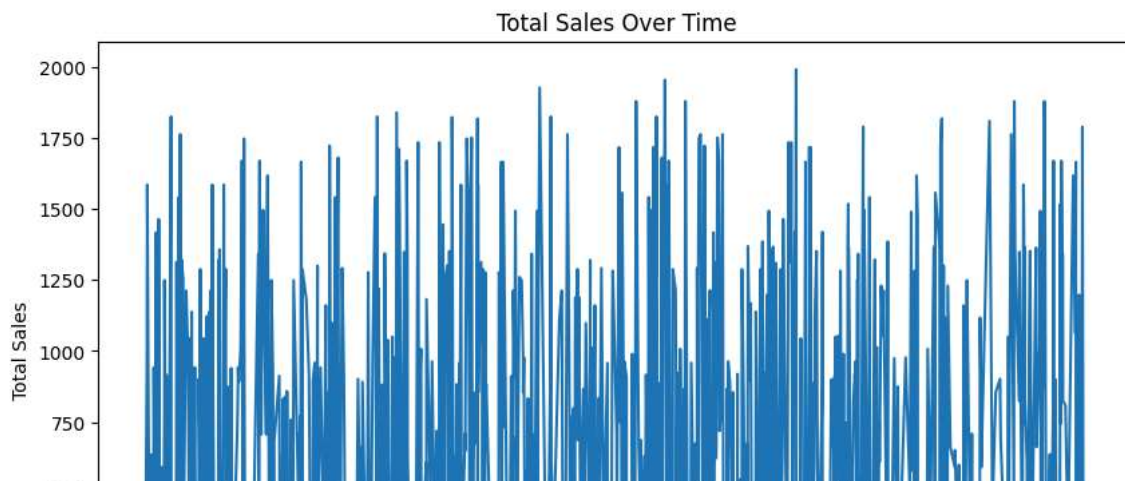
merged.to_csv('merged_dataset.csv', index=False)
```

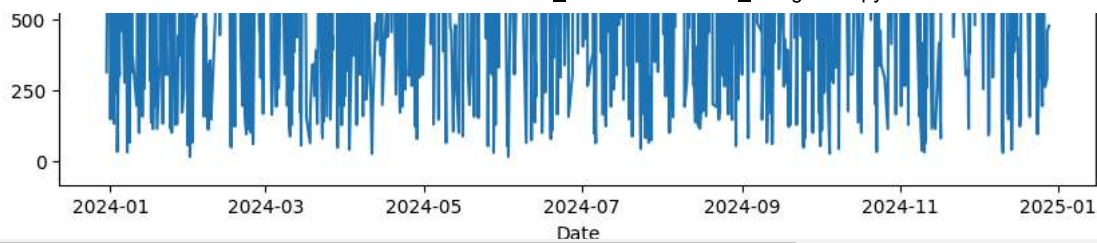
```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   CustomerID  200 non-null   object
 1   CustomerName 200 non-null   object
 2   Region      200 non-null   object
 3   SignupDate  200 non-null   object
dtypes: object(4)
memory usage: 6.4+ KB
   CustomerID  CustomerName      Region  SignupDate
0      C0001  Lawrence Carroll  South America  2022-07-10
1      C0002  Elizabeth Lutz      Asia  2022-02-13
2      C0003  Michael Rivera  South America  2024-03-07
3      C0004  Kathleen Rodriguez  South America  2022-10-09
4      C0005  Laura Weber      Asia  2022-08-15 None
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 100 entries, 0 to 99
Data columns (total 4 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   ProductID  100 non-null   object
 1   ProductName 100 non-null   object
 2   Category    100 non-null   object
 3   Price       100 non-null   float64
dtypes: float64(1), object(3)
memory usage: 3.3+ KB
   ProductID  ProductName      Category  Price
0      P001  ActiveWear Biography      Books  169.30
1      P002  ActiveWear Smartwatch  Electronics  346.30
2      P003  ComfortLiving Biography      Books  44.12
3      P004  BookWorld Rug      Home Decor  95.69
4      P005  TechPro T-Shirt      Clothing  429.31 None
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1000 entries, 0 to 999
Data columns (total 7 columns):
 #   Column      Non-Null Count  Dtype
---  ---
 0   TransactionID  1000 non-null  object
 1   CustomerID     1000 non-null  object
 2   ProductID      1000 non-null  object
 3   TransactionDate 1000 non-null  object
 4   Quantity       1000 non-null  int64
 5   TotalValue     1000 non-null  float64
 6   Price          1000 non-null  float64
dtypes: float64(2), int64(1), object(4)
memory usage: 54.8+ KB
   TransactionID  CustomerID  ProductID      TransactionDate  Quantity \
0      T00001      C0199      P067  2024-08-25 12:38:23      1
1      T00112      C0146      P067  2024-05-27 22:23:54      1
2      T00166      C0127      P067  2024-04-25 07:38:55      1
3      T00272      C0087      P067  2024-03-26 22:55:37      2
4      T00363      C0070      P067  2024-03-21 15:10:10      3

   TotalValue  Price
0      300.68  300.68
1      300.68  300.68
2      300.68  300.68
3      601.36  300.68
4      902.04  300.68 None
Top 5 Selling Products:
   ProductName
ActiveWear Smartwatch      100
SoundWave Headphones       97
HomeSense Desk Lamp        81
ActiveWear Rug              79
SoundWave Cookbook         78
Name: Quantity, dtype: int64

```





```

from sklearn.metrics.pairwise import cosine_similarity
import numpy as np

customer_features = merged.groupby('CustomerID').agg({
    'Quantity': 'sum',
    'TotalValue': 'sum',
    'Category': lambda x: ' '.join(x)
}).reset_index()

category_encoded = pd.get_dummies(customer_features['Category'])
customer_features = pd.concat([customer_features, category_encoded], axis=1).drop(columns=['Category'])

features = customer_features.drop(columns=['CustomerID'])
similarity_matrix = cosine_similarity(features)

lookalike = {}
for i, customer_id in enumerate(customer_features['CustomerID']):
    similar_indices = np.argsort(similarity_matrix[i])[::-1][1:4]
    similar_customers = [(customer_features.iloc[j]['CustomerID'], similarity_matrix[i][j]) for j in similar_indices]
    lookalike[customer_id] = similar_customers

lookalike_data = []
for cust_id, similarities in lookalike.items():
    row = [cust_id] + [item for sublist in similarities for item in sublist]
    lookalike_data.append(row)

columns = ['CustomerID', 'Lookalike1', 'Score1', 'Lookalike2', 'Score2', 'Lookalike3', 'Score3']
lookalike_df = pd.DataFrame(lookalike_data, columns=columns)
lookalike_df.to_csv('Vaishnavi_GaneshChaudhari_Lookalike.csv', index=False)

from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import davies_bouldin_score
import matplotlib.pyplot as plt

merged = pd.read_csv('merged_dataset.csv')
print(merged.columns)
merged = transactions.merge(customers, on='CustomerID').merge(products, on='ProductID')

features = merged.groupby('CustomerID').agg({
    'Quantity': 'sum',
    'TotalValue': 'sum'
}).reset_index()

scaler = StandardScaler()
scaled_features = scaler.fit_transform(features.drop(columns=['CustomerID']))

kmeans = KMeans(n_clusters=4, random_state=42)
clusters = kmeans.fit_predict(scaled_features)

features['Cluster'] = clusters

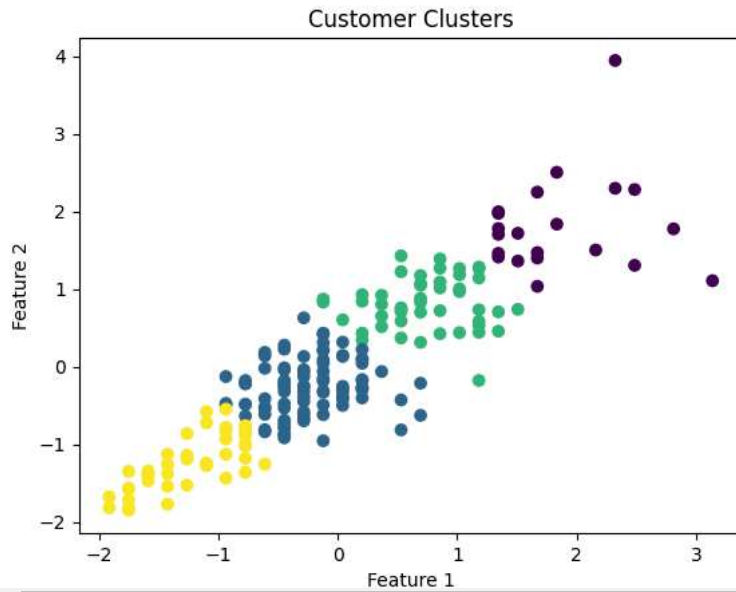
db_index = davies_bouldin_score(scaled_features, clusters)
print(f'Davies-Bouldin Index: {db_index}')

plt.scatter(scaled_features[:, 0], scaled_features[:, 1], c=clusters, cmap='viridis')
plt.title('Customer Clusters')
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.show()

features.to_csv('Vaishnavi_GaneshChaudhari_Clustering.csv', index=False)

```

```
Index(['TransactionID', 'CustomerID', 'ProductID', 'TransactionDate',
      'Quantity', 'TotalValue', 'Price_X', 'CustomerName', 'Region',
      'SignupDate', 'ProductName', 'Category', 'Price_Y'],
      dtype='object')
Davies-Bouldin Index: 0.7212797181816303
```



```
from sklearn.cluster import KMeans
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import davies_bouldin_score
import matplotlib.pyplot as plt

print(merged.columns)

features = merged.groupby('CustomerID').agg({
    'Quantity': 'sum',
    'TotalValue': 'sum'
}).reset_index()

scaler = StandardScaler()
scaled_features = scaler.fit_transform(features.drop(columns=['CustomerID']))

kmeans = KMeans(n_clusters=4, random_state=42)
clusters = kmeans.fit_predict(scaled_features)

features['Cluster'] = clusters

db_index = davies_bouldin_score(scaled_features, clusters)
print(f'Davies-Bouldin Index: {db_index}')

plt.scatter(scaled_features[:, 0], scaled_features[:, 1], c=clusters, cmap='viridis')
plt.title('Customer Clusters')
plt.xlabel('Feature 1')
plt.ylabel('Feature 2')
plt.show()

features.to_csv('Vaishnavi_GaneshChaudhari_Clustering.csv', index=False)
```