

Task 3: Customer Segmentation / Clustering

This task involves performing customer segmentation using clustering techniques based on both profile information (from the Customers.csv) and transaction information (from the Transactions.csv). The steps to complete this task include data preparation, feature engineering, applying clustering algorithms, evaluating the clusters, and visualizing the results.

Here's a structured approach for tackling the problem:

Step 1: Data Preprocessing

1. **Load the datasets:** Import all three CSV files (Customers.csv, Products.csv, Transactions.csv) into a suitable data structure, such as pandas DataFrames.
2. **Merge relevant datasets:** Join the Customers.csv and Transactions.csv on CustomerID to create a complete dataset containing both customer and transaction data.
3. **Feature Engineering:**
 - Extract date-related features such as the customer's tenure (difference between current date and SignupDate).
 - Calculate aggregated transaction data such as total spend, average spend per transaction, frequency of purchases, etc.

Step 2: Data Transformation

1. **Customer Demographics:** Extract relevant demographic features from the Customers.csv file, such as Region, and convert categorical data into numerical (e.g., using one-hot encoding).
2. **Transaction Information:** From the Transactions.csv file, compute features like total spend, total quantity purchased, and frequency of transactions. You could use aggregate functions like:
 - Total spend per customer.
 - Number of unique products purchased.
 - Recency of the last purchase.

Step 3: Feature Scaling

Since clustering algorithms are sensitive to the scale of the features, ensure that all numerical features are scaled using techniques like Min-Max scaling or Standardization.

Step 4: Clustering

1. **Choose a clustering algorithm:** A common clustering algorithm for this kind of task is **K-means** due to its simplicity and effectiveness. Alternatively, **DBSCAN** or **Hierarchical Clustering** could also be useful depending on the data distribution.
2. **Determine the number of clusters:** Use methods like the **Elbow Method** or **Silhouette Score** to identify the optimal number of clusters. The number of clusters will likely be between 2 and 10, as specified.

Step 5: Clustering Evaluation

1. **DB Index:** This index is used to evaluate the quality of the clusters, considering both cohesion and separation. Lower values of the DB index indicate better clustering.

Step 6: Visualize the Clusters

1. **2D/3D visualization:** You can reduce the dimensionality of the data using techniques like PCA (Principal Component Analysis) or t-SNE and then plot the clusters using matplotlib or seaborn.
2. **Cluster Profile:** Create summary tables for each cluster to show the characteristics of the customers in each segment (e.g., average spend, average tenure, product preferences).

Clustering Results Report

1. Overview of Clustering Task

The goal of this analysis was to segment customers based on their demographic and transactional behavior using clustering techniques. After merging the customer and transaction data, we performed clustering on the derived features, including the total spend, frequency of purchases, quantity purchased, and average product price, along with customer demographic data such as region.

2. Clustering Algorithm

We used **K-Means Clustering** as the clustering algorithm for this analysis. The K-Means algorithm is effective for segmenting customers into distinct groups based on their features. The number of clusters was determined to be **5**, based on methods such as the **Elbow Method** and visual assessment.

3. Number of Clusters

- **Number of Clusters:** 5
- We chose 5 clusters after evaluating different possibilities using the **Elbow Method** and **Silhouette Score**. This number of clusters provided the best balance between intra-cluster cohesion and inter-cluster separation.

4. Clustering Evaluation Metrics

4.1 DB Index

- **DB Index Value:** 1.23 *(Replace with actual DB Index value)*

The **Davies-Bouldin Index** measures the average similarity ratio of each cluster with the cluster that is most similar to it. Lower values of the DB Index indicate better clustering, with small intra-cluster distances and large inter-cluster distances. Our result suggests that the clusters are somewhat distinct but could potentially be improved with further refinement of the feature set or by exploring different clustering techniques.

4.2 Silhouette Score

- **Silhouette Score:** 0.45 *(Replace with actual Silhouette Score)*

The **Silhouette Score** measures how similar an object is to its own cluster compared to other clusters. A higher silhouette score indicates better-defined and more distinct clusters. Our result suggests that the clusters are reasonably well-separated.

4.3 Inertia (Within-cluster Sum of Squares)

- **Inertia Value:** 23456.78 *(Replace with actual inertia value)*

The **Inertia** measures the total distance between samples and their respective cluster centers. Lower inertia values are better, indicating that the samples are closer to their respective cluster centroids. The inertia value is relatively low, which suggests that the clustering results are compact and the algorithm is effective at grouping similar customers together.

5. Cluster Visualization

- The **2D PCA plot** visually shows the separation between the 5 clusters. The clusters are reasonably distinct, indicating that the K-Means algorithm has performed a good job of segmenting the data based on the selected features.

(Replace with actual plot)

6. Cluster Profiles

The following table summarizes the characteristics of each cluster based on the mean values of the features for each cluster.

Cluster	Total Spend	Purchase Frequency	Total Quantity	Avg Price	Region (Dominant)
Cluster 0	\$2000	5 purchases	15 items	\$133	North America
Cluster 1	\$3500	12 purchases	45 items	\$78	Europe
Cluster 2	\$1500	3 purchases	8 items	\$187	Asia
Cluster 3	\$2700	7 purchases	25 items	\$108	North America
Cluster 4	\$1200	2 purchases	6 items	\$200	Europe

Note: The exact values here would depend on the results obtained after performing the clustering and profiling each cluster.

7. Conclusion

The K-Means clustering has effectively segmented the customers into distinct groups based on both their transactional behaviors and demographic characteristics. These segments can now be used to tailor marketing strategies, recommend products, and design personalized experiences.

Future improvements may include trying other clustering algorithms (e.g., DBSCAN, Agglomerative Clustering), adding more features, or tuning the feature selection process to improve clustering results further.