

Customer Segmentation using KMeans Clustering

Objective

The main goal of the script is to segment customers into distinct clusters based on their purchasing behavior. This segmentation can be useful for various applications, such as targeted marketing, personalized recommendations, and customer retention strategies.

Workflow Overview

1. Data Loading and Preprocessing:

- The script loads three datasets:
 - **Customers.csv:** Contains information about the customers.
 - **Products.csv:** Contains product-related details (including price).
 - **Transactions.csv:** Contains transaction data such as the quantity of products purchased and total value.
- The CustomerID and ProductID columns are standardized across datasets to ensure consistent merging.

2. Data Merging:

- The datasets are merged into a single dataframe that combines transactional data, customer details, and product information.
- The merged data is then aggregated by CustomerID to create customer profiles with the following features:
 - **Quantity:** Total number of items purchased.
 - **TotalValue:** Total monetary value of the customer's purchases.
 - **Price (Average):** Average price of the purchased products.

3. Data Normalization:

- The customer profile features are normalized using StandardScaler to ensure they are on the same scale for clustering.

4. Optimal Clusters Determination (Elbow Method):

- The script uses the Elbow Method to find the optimal number of clusters by plotting the inertia (sum of squared distances from the cluster centroids).
- The elbow point (the point of diminishing returns in inertia reduction) helps identify the best number of clusters for the dataset.

5. KMeans Clustering:

- After determining the optimal number of clusters (set to 4 based on the Elbow Method), the script performs KMeans clustering.
 - Each customer is assigned to one of the clusters based on their purchasing behavior.
6. Evaluation (Davies-Bouldin Index):
- The Davies-Bouldin index is computed to evaluate the clustering performance. This index measures the average similarity ratio of each cluster with the one most similar to it; lower values indicate better clustering.
 - The Davies-Bouldin index is printed to assess the quality of the clusters.
7. Visualization:
- The customer clusters are visualized in a scatter plot using the first two principal components of the normalized data. Each point represents a customer, and the colors correspond to different clusters.
8. Output:
- The clustering results (CustomerID and their assigned cluster) are saved to a CSV file, Customer_Clusters.csv.
-

Insights and Applications

1. Segmentation Insights:
- By segmenting customers into distinct groups, businesses can tailor their strategies to each group. For example:
 - High-value customers may receive premium offers.
 - Low-value customers may be targeted with retention campaigns.
2. Marketing Strategy:
- Different clusters can be targeted with personalized campaigns. For example, a cluster of frequent buyers might be offered loyalty programs, while a cluster of low-spending customers might be incentivized with discounts.
3. Product Recommendations:
- Similar clusters can have similar preferences, and businesses can recommend products based on the purchasing behavior of other customers within the same cluster.
-

Evaluation Metrics

1. Davies-Bouldin Index:
- The Davies-Bouldin index provides a quantitative measure of cluster quality. A lower index value indicates that the clusters are more distinct and well-separated.

- The index is printed in the script, and a value closer to zero is ideal.
-

Output Example

The clustering results saved to Customer_Clusters.csv will have the following format:

CustomerID Cluster

C0001	0
C0002	1
C0003	2
C0004	0
C0005	3

This output helps identify which customers belong to which cluster and can be used for further analysis or marketing efforts.

Key Visualizations

1. Elbow Method for Optimal Clusters:

- A plot of inertia vs. the number of clusters helps visualize the elbow point, guiding the selection of the optimal number of clusters.

2. Customer Cluster Visualization:

- A scatter plot is generated using the first two features (principal components) of the normalized data, where customers are colored based on their cluster assignments. This allows us to visually assess how well the customers are grouped.
-

Strengths

- **Scalable Approach:** The script scales well to larger datasets, as KMeans is efficient and handles high-dimensional data well after normalization.
 - **Easy Interpretation:** The resulting clusters are easy to interpret and can be used directly for customer targeting and analysis.
 - **Visualization:** The use of visualizations helps in understanding the customer segmentation.
-

Suggestions for Improvement

1. Additional Feature Engineering:

- Incorporate more customer attributes (e.g., demographics, geographical location) into the clustering process to improve the segmentation quality.

- Consider using advanced techniques such as PCA (Principal Component Analysis) to reduce the dimensionality before clustering.

2. Cluster Profiling:

- After clustering, additional profiling of each cluster (e.g., average quantity, total value, product preferences) can provide deeper insights into the behavior of each segment.

3. Alternative Clustering Techniques:

- Other clustering methods like DBSCAN (Density-Based Spatial Clustering of Applications with Noise) or Hierarchical Clustering could be tested to compare results.

Conclusion

This script effectively segments customers based on their purchasing behavior using KMeans clustering. The approach allows businesses to gain valuable insights into customer preferences, which can be used for personalized marketing and product recommendations. By refining the clustering process and incorporating additional customer features, the segmentation could be made even more robust and insightful.