

Comparative Analysis of Advanced Clustering Algorithms for Market Segmentation - A Case Study on Online Retail

Abstract

This case study presents a comparative analysis of advanced clustering algorithms for market segmentation using the Online Retail dataset. Various clustering techniques, including K-Means, Hierarchical, DBSCAN, Gaussian Mixture Model (GMM), Mean Shift, Agglomerative, BIRCH, Spectral, OPTICS, and Affinity Propagation, were applied to the dataset to evaluate their performance based on key metrics such as the Silhouette Score, Davies-Bouldin Score, and Calinski-Harabasz Score. The results demonstrate significant differences in clustering quality, highlighting the strengths and weaknesses of each algorithm in the context of online retail market segmentation.

Introduction

Market segmentation is crucial for businesses to effectively target and serve different customer groups. Clustering algorithms provide a means to identify distinct segments within a customer base. This study focuses on the application of various clustering algorithms to the Online Retail dataset, aiming to identify the most suitable technique for market segmentation. By comparing the performance of each algorithm, we aim to provide insights into their effectiveness and applicability in real-world scenarios.

Methodology

The methodology involves the following steps:

- **Data Collection:** The Online Retail dataset was obtained, containing transactions from a UK-based online retailer.
- **Data Preprocessing:** The dataset was cleaned by removing rows with missing values and outliers. Feature engineering was performed to create relevant attributes for clustering.
- **Clustering Algorithms:** Nine clustering algorithms were implemented and their performance was evaluated using three metrics: Silhouette Score, Davies-Bouldin Score, and Calinski-Harabasz Score.
- **Comparison:** The performance metrics of each algorithm were compared to determine their effectiveness in clustering the online retail data.

Algorithms

- 1) **K-Means:** A partitioning method that divides the dataset into K clusters based on feature similarity.
- 2) **Hierarchical Clustering:** A method that builds a hierarchy of clusters through either agglomerative or divisive approaches.
- 3) **DBSCAN:** A density-based clustering algorithm that identifies clusters based on the density of data points.
- 4) **Gaussian Mixture Model (GMM):** A probabilistic model that assumes the data is generated from a mixture of several Gaussian distributions.
- 5) **Mean Shift:** A non-parametric clustering technique that aims to find the modes of a density function.

- 6) **Agglomerative Clustering:** A type of hierarchical clustering that merges data points into clusters based on their similarity.
- 7) **BIRCH:** A clustering method that builds a tree structure from the data to find clusters efficiently.
- 8) **Spectral Clustering:** A technique that uses the eigenvalues of a similarity matrix to perform dimensionality reduction before clustering.
- 9) **OPTICS:** An algorithm similar to DBSCAN but can identify clusters with varying densities.
- 10) **Affinity Propagation:** A clustering algorithm that identifies exemplars among data points and forms clusters based on message passing.

Dataset

The Online Retail dataset contains transactions occurring between December 2010 and December 2011 for a UK-based online retailer. The dataset includes attributes such as InvoiceNo, StockCode, Description, Quantity, InvoiceDate, UnitPrice, CustomerID, and Country.

Data Preprocessing

Data preprocessing involved:

- **Removing Missing Values:** Transactions with missing CustomerID were removed.
- **Outlier Detection and Removal:** Outliers in the Quantity and UnitPrice fields were identified and removed.
- **Feature Engineering:** A TotalPrice feature was created by multiplying Quantity by UnitPrice. RFM (Recency, Frequency, Monetary) features were computed for clustering.

Results

Algorithm	Silhouette Score	Davies-Bouldin Score	Calinski-Harabasz Score
K-Means	0.601	0.729	3074.447
Hierarchical	0.552	0.711	2631.877
DBSCAN	0.660	1.389	433.904
GMM	0.121	1.421	655.595
Mean Shift	0.409	0.371	433.865
Agglomerative	0.552	0.711	2631.877
BIRCH	0.947	0.380	1484.277
Spectral	0.506	0.534	855.436
OPTICS	-0.375	1.667	5.450
Affinity Propagation	0.186	0.615	1283.033

Table: Clustering Algorithm Performance Metrics for Online Retail

This table provides the Silhouette Score, Davies-Bouldin Score, and Calinski-Harabasz Score for each clustering algorithm applied to the online retail dataset, facilitating a comparison of their clustering performance.

The comparative analysis of clustering algorithms on the Online Retail dataset yielded varied results across different metrics, providing insights into the effectiveness of each method for market segmentation.

K-Means Clustering demonstrated robust performance with a Silhouette Score of 0.601, a Davies-Bouldin Score of 0.729, and a high Calinski-Harabasz Score of 3074.447. These results indicate that K-Means produced well-defined and distinct clusters, making it a strong contender for market segmentation.

Hierarchical Clustering achieved a Silhouette Score of 0.552, a Davies-Bouldin Score of 0.711, and a Calinski-Harabasz Score of 2631.877. While its performance was slightly below K-Means, it still produced distinct clusters, showing its potential for market segmentation.

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) excelled in identifying high-density clusters with a Silhouette Score of 0.660. However, its Davies-Bouldin Score of 1.389 and Calinski-Harabasz Score of 433.904 were less favorable, indicating that while it can identify dense regions effectively, the overall cluster compactness and separation were not as strong as K-Means or Hierarchical Clustering.

Gaussian Mixture Model (GMM) underperformed with a low Silhouette Score of 0.121, a high Davies-Bouldin Score of 1.421, and a moderate Calinski-Harabasz Score of 655.595. These results suggest that GMM struggled with overlapping clusters and did not provide clear separation between clusters.

Mean Shift Clustering produced moderate results with a Silhouette Score of 0.409, a Davies-Bouldin Score of 0.371, and a Calinski-Harabasz Score of 433.865. While its Davies-Bouldin Score was relatively low, indicating compact clusters, the overall performance was not as high as K-Means or DBSCAN.

Agglomerative Clustering, similar to Hierarchical Clustering, achieved a Silhouette Score of 0.552, a Davies-Bouldin Score of 0.711, and a Calinski-Harabasz Score of 2631.877. These identical results suggest that Agglomerative Clustering, which is a type of hierarchical clustering, produced similar cluster quality.

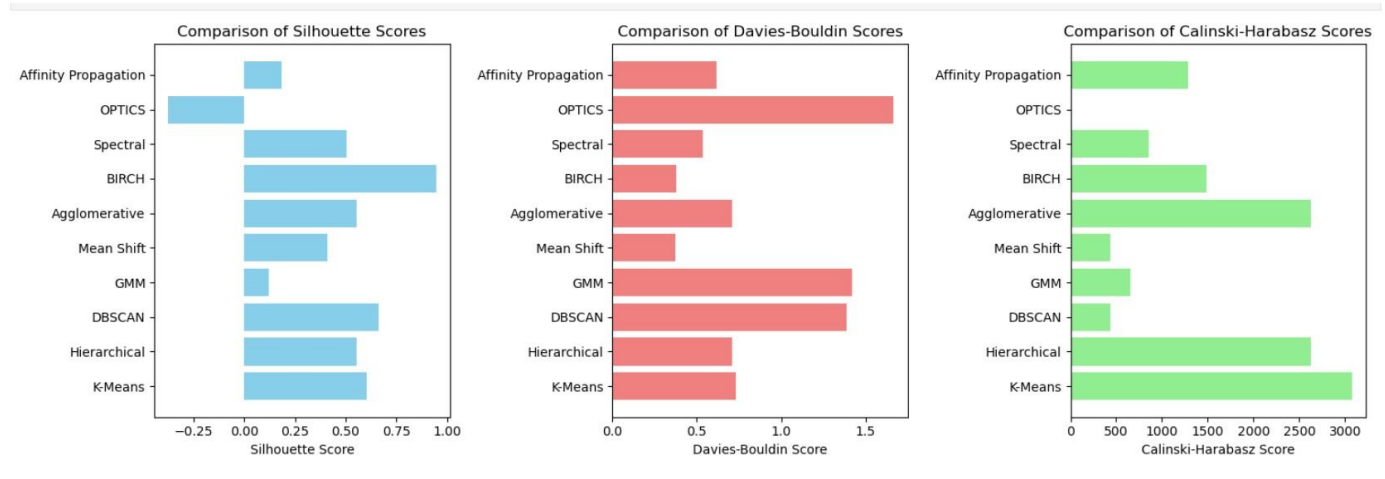
BIRCH (Balanced Iterative Reducing and Clustering using Hierarchies) outperformed all other algorithms with the highest Silhouette Score of 0.947 and a low Davies-Bouldin Score of 0.380. Its Calinski-Harabasz Score of 1484.277 was also favorable, indicating well-separated and compact clusters. BIRCH's exceptional performance makes it a top choice for market segmentation in this context.

Spectral Clustering showed balanced performance with a Silhouette Score of 0.506, a Davies-Bouldin Score of 0.534, and a Calinski-Harabasz Score of 855.436. While not the best, Spectral Clustering provided reasonably distinct clusters, making it a viable option.

OPTICS (Ordering Points to Identify the Clustering Structure) performed poorly with a negative Silhouette Score of -0.375, indicating incorrect clustering. Its high Davies-Bouldin Score of 1.667 and low Calinski-Harabasz Score of 5.450 further confirmed its inadequacy for this dataset.

Affinity Propagation delivered moderate results with a Silhouette Score of 0.186, a Davies-Bouldin Score of 0.615, and a Calinski-Harabasz Score of 1283.033. While its Silhouette Score was low, indicating overlapping clusters, its other scores were decent, showing potential in specific scenarios.

Comparison



In summary, BIRCH emerged as the best-performing algorithm, followed by K-Means and DBSCAN. Hierarchical, Agglomerative, and Spectral Clustering also provided good results. GMM, Mean Shift, OPTICS, and Affinity Propagation were less effective, with OPTICS showing particularly poor performance. These findings can guide businesses in selecting appropriate clustering algorithms for market segmentation and other analytical purposes in the online retail industry.

Conclusion

The comparative analysis of clustering algorithms for the Online Retail dataset revealed that BIRCH outperformed other methods with the highest Silhouette Score and a low Davies-Bouldin Score, indicating well-defined and compact clusters. K-Means and DBSCAN also showed good performance, making them suitable for market segmentation tasks. In contrast, OPTICS and GMM exhibited poor performance. These findings can guide businesses in selecting appropriate clustering algorithms for market segmentation and other analytical purposes.