

REPORT

The Davies-Bouldin Index (DB Index) is a metric for evaluating clustering algorithms, where a lower DB Index indicates that clusters are well-separated and compact. It essentially measures the average 'similarity' between clusters, where similarity is a ratio of within-cluster distances to between-cluster distances.

No of Clusters formed=5

Determination of the Number of Clusters

1. Using the Elbow Method:
 - The elbow method involves plotting the within-cluster sum of squares (WCSS) against the number of clusters. You look for a point where the decrease in WCSS becomes less rapid, which can be visually identified as an "elbow." This method was implemented in the `evaluate_clusters` function provided in your code, which plotted WCSS for a range from 2 to 10 clusters.
2. Silhouette Scores:
 - Alongside the WCSS, silhouette scores were also calculated for each potential number of clusters from 2 to 10. Silhouette scores measure how similar an object is to its own cluster compared to other clusters. A higher silhouette score indicates a model with better defined clusters.

These two methods were employed to help select an optimal number of clusters by providing a visual and quantitative way to assess the compactness and separation of clusters at various levels.

Comparative Analysis:

DB Index Scores:

- **K-Means Clustering:** 0.8525
- **Hierarchical Clustering:** 0.8851

A lower DB Index is preferable as it indicates a clustering configuration with better separation and more distinct clusters.

Reasons for Different Scores:

1. Algorithm Characteristics:

- **K-Means** aims to minimize the variance within each cluster, which often results in more compact clusters, particularly if the data naturally segregates into circular (spherical) groupings. K-Means performs well when clusters are distinct and well separated.
- **Hierarchical Clustering**, especially with Ward linkage, also aims to minimize variance within clusters. However, it builds clusters by hierarchically merging or splitting them based on distance metrics. This method can sometimes retain inherent data hierarchies better but might be less effective in minimizing variance at the global level compared to K-Means.

2. Cluster Shape and Size:

- **K-Means** often performs better when the clusters are roughly equal in terms of data points and are isotropic. In contrast, hierarchical clustering might not enforce such equality, leading to clusters of varying sizes and densities, which could explain the slightly higher DB Index score.

3. Sensitivity to Noise and Outliers:

- **Hierarchical Clustering** can sometimes be more sensitive to noise and outliers than K-Means. Since it builds clusters step by step, outliers can significantly influence the path of hierarchy formation, potentially leading to less optimal clustering (higher DB Index).

Which is Better?

Given that K-Means provided a **lower DB Index of 0.8525** compared to hierarchical clustering's 0.8851, K-Means would generally be considered better for this particular dataset based on this metric. The lower score suggests that K-Means managed to produce clusters that are more compact and better separated compared to those produced by hierarchical clustering.

Considerations:

- **Cluster Interpretation and Use Case:** Even though K-Means scored better in terms of the DB Index, hierarchical clustering offers insights into data structure that K-Means does not, such as the hierarchical relationships between clusters. Depending on the specific needs of your analysis or business case, this information might be valuable.
- **Validation:** It's beneficial to validate these findings with additional metrics and visualizations. Metrics like Silhouette Score or even a manual inspection of cluster

content can provide deeper insights into the quality of clustering beyond just the DB Index.

- **Experiment with Parameters:** Adjusting parameters like the number of clusters in K-Means or the linkage criterion in hierarchical clustering might yield different results, potentially improving the clustering performance of each method.