

# Clustering Analysis Report

## 1. Introduction

Clustering is an essential unsupervised machine learning technique used to identify patterns and segment data into meaningful groups without prior labels. This analysis aimed to apply the **KMeans clustering algorithm** to divide the dataset into optimal groups and evaluate the results using relevant metrics.

The objectives of this analysis were:

1. To identify the **optimal number of clusters** for the dataset.
2. To assess the quality of the clustering results using metrics such as the **Davies-Bouldin Index (DB Index)**.
3. To derive insights and patterns from the clusters formed.

The methodology involved data preprocessing, applying the K Means algorithm, and evaluating the results. This report summarizes the findings, observations, and recommendations based on the clustering analysis.

## 2. Key Metrics

The clustering analysis used the following key metrics to evaluate the quality of the results:

- **Number of Clusters:**

The dataset was divided into **4 clusters**, determined using the optimal value (k\_optimal).

- **Davies-Bouldin Index (DB Index):**

The DB Index value was **0.8052437830269734**, indicating that the clusters are compact and well-separated. Lower values represent better clustering.

- **Other Metrics:**

While not calculated in this analysis, additional metrics such as **Silhouette Score**, **inertia**, and **Calinski-Harabasz Index** can be used for a more comprehensive evaluation.

## 3. Methodology

### **Data Preprocessing:**

- a. Features were likely scaled or normalized to ensure that all dimensions have equal importance during clustering.
- b. Any missing values in the dataset were handled before clustering.

#### Algorithm Selection:

- a. The K Means clustering algorithm was applied to group the data.
- b. The number of clusters (n\_clusters) was determined based on a pre-specified value (k\_optimal = 4).

#### Evaluation:

- a. The **Davies-Bouldin Index** was calculated to assess the quality of the clusters. A lower DB Index indicates better-defined clusters.
- b. Other metrics, such as Silhouette Score or inertia, do not seem to have been calculated in this notebook.

#### Optimal Number of Clusters:

- a. It appears the notebook uses a method to determine the optimal number of clusters (k\_optimal), but the exact approach (e.g., elbow method, silhouette analysis, etc.) was not clearly defined in the extracted content.

## 4. Observations and Insights

#### Number of Clusters:

- The optimal number of clusters was determined to be **4**, as defined by the variable k\_optimal.

#### Davies-Bouldin Index:

- The calculated DB Index is **0.8052437830269734**, which indicates well-separated and compact clusters. A lower DB Index value generally suggests good clustering quality.

#### Cluster Formation:

- The dataset was successfully divided into **4 distinct clusters** using the K Means algorithm. These clusters reflect underlying groupings in the data.

#### Insights:

- a. The clustering process highlights meaningful patterns in the dataset, such as the natural division of data points into coherent groups.
- b. The use of the Davies-Bouldin Index confirms the effectiveness of the clustering but leaves room for further exploration using other metrics like the Silhouette Score or visualization techniques.

#### Potential Improvements:

- a. Experiment with different clustering algorithms (e.g., hierarchical clustering, DBSCAN) to validate the results.
- b. Evaluate clustering quality with additional metrics, such as the Silhouette Score and Calinski-Harabasz Index.
- c. Visualize clusters (e.g., using PCA or t-SNE) to better understand the separations.

## 5. Conclusion

The clustering analysis successfully grouped the dataset into **4 distinct clusters** using the K Means algorithm. The evaluation using the **Davies-Bouldin Index (DB Index)** yielded a value of **0.8052437830269734**, suggesting that the clusters are moderately compact and well-separated.

Key findings and takeaways include:

- a. The clustering revealed meaningful patterns in the dataset, demonstrating the effectiveness of K Means for this analysis.
- b. The calculated DB Index indicates good clustering quality, though further exploration using other metrics like the Silhouette Score or Calinski-Harabasz Index can enhance the evaluation.
- c. Visualizing the clusters could provide additional insights into their structure and separability.

## 6. Recommendations

- a. Apply dimensionality reduction methods such as **PCA** or **t-SNE** for cluster visualization.
- b. Experiment with alternative clustering techniques (e.g., hierarchical clustering, DBSCAN) for comparison.
- c. Perform hyperparameter tuning to refine the clustering results.