

# Report on Clustering Results

Based on the results from different clustering algorithms, here's an overview of the clustering outcomes, including the number of clusters formed, DB Index value, and other relevant metrics.

## 1. KMeans Clustering

- **Number of Clusters:** 10
- **DB Index:** 0.6885
- **Additional Comments:**
  - KMeans is a centroid-based algorithm, and it works well when the clusters are well-separated and spherical in shape.
  - A DB Index value of 0.6885 is relatively moderate, indicating that the clusters formed are somewhat compact but there is still room for improvement.
  - KMeans tends to favor a fixed number of clusters, and the model found that 10 clusters provided the best separation in terms of DB Index.

## 2. Agglomerative Clustering

- **Number of Clusters:** 8
- **DB Index:** 0.6421
- **Additional Comments:**
  - Agglomerative Clustering is a hierarchical algorithm that forms a hierarchy of clusters by merging them based on their similarity.
  - The best DB Index value of 0.6421 was achieved with 8 clusters and the 'average' linkage method, which uses the average of the distances between all pairs of points in the two clusters being merged.
  - This indicates relatively compact clusters with moderate separation between them.

### **3. Gaussian Mixture Model (GMM)**

- **Number of Clusters:** 4
- **DB Index:** 0.7690
- **Additional Comments:**
  - GMM is a probabilistic model that assumes each cluster is generated from a Gaussian distribution, making it suitable for modeling clusters with varying shapes and densities.
  - The GMM method performed well, with a DB Index of 0.7690, indicating that the clustering formed by GMM is well-separated and compact.
  - The model chose 4 clusters as the best configuration, which might suggest that the data has fewer natural groupings compared to KMeans.

### **4. DBSCAN (Density-Based Spatial Clustering of Applications with Noise)**

- **Number of Clusters:** 3
- **DB Index:** 0.5154
- **Additional Comments:**
  - DBSCAN is a density-based clustering algorithm that groups together points that are closely packed and marks points that are in low-density regions as outliers (noise).
  - The algorithm formed 3 clusters, with a DB Index of 0.5154, the lowest among the models. This suggests that DBSCAN performed better in identifying distinct clusters with clear boundaries and low noise.
  - The DB Index value indicates good cluster separation, which is typical for density-based methods like DBSCAN.

### **5. MeanShift**

- **Number of Clusters:** 2 (Based on bandwidth of 0.1)
- **DB Index:** 0.3852
- **Additional Comments:**

- MeanShift is a non-parametric clustering technique that shifts the centroid of the data points iteratively towards regions of higher density.
- With a bandwidth of 0.1, the model identified 2 clusters, yielding the best DB Index of 0.3852.
- The relatively low DB Index value suggests that the clustering is well-separated, with compact and well-defined clusters.

## **Conclusion:**

- **Best DB Index:** The **Gaussian Mixture Model (GMM)** achieved the best DB Index (0.7690), indicating that its clustering is the most well-separated and compact among all methods. This suggests that the data has clusters that fit well with the Gaussian distribution assumption.
- **Density-Based Clustering: DBSCAN** (0.5154 DB Index) is effective for identifying clusters with clear density-based separation. It performed well with 3 clusters, indicating that the data likely has density-based structures rather than simple spherical shapes.
- **Parameter Sensitivity:** The **MeanShift** method (2 clusters) with a low bandwidth value yielded the best DB Index, indicating that clustering can be highly sensitive to the chosen bandwidth parameter.
- **Hierarchical Clustering: Agglomerative Clustering** (0.6421 DB Index) with 8 clusters and 'average' linkage provided a good separation but was not as effective as the other models in this particular case.
- **Cluster Count Differences:** KMeans preferred a larger number of clusters (10), while other algorithms, especially **Gaussian Mixture** and **MeanShift**, identified fewer clusters. This suggests that the optimal number of clusters depends on the underlying structure and distribution of the data.

In conclusion, **Gaussian Mixture Model (GMM)** and **DBSCAN** are the most effective clustering techniques for this dataset, with the lowest DB Index values, indicating better clustering performance.