

ML ASSIGNMENT – 2

Report

Name: Borra Pujith Ganesh

Roll No: 160123737316

IT 3 - V Sem

Title

**Comparative Analysis and Optimization of
Clustering Techniques for High-Dimensional
Data**

Paper Referred

Vishnu Vardhan Baligodugula, Fathi Amsaad,
***“Unsupervised Learning: Comparative Analysis of
Clustering Techniques on High-Dimensional Data,”***
arXiv, March 2025.

1. Introduction

The goal of this project is to replicate and enhance the analysis from the paper “**Unsupervised Learning: Comparative Analysis of Clustering Techniques on High-Dimensional Data**” (Baligodugula & Amsaad, 2025).

While the referenced study compared **K-Means**, **DBSCAN**, and **Spectral Clustering** using dimensionality reduction methods (PCA, t-SNE, UMAP), it lacked **automated hyperparameter optimization**, **robust evaluation**, and **scalability checks**.

In this enhanced study, we implemented a **systematic parameter search**, introduced **Agglomerative Clustering**, and applied **comprehensive visual diagnostics** to achieve improved clustering quality and interpretability.

2. Dataset Description

Attribute	Description
Source	Synthetic and UCI datasets (Wine, Iris, and High-Dimensional Synthetic dataset generated using make_blobs)
Samples	~1,000
Features	10–50 (after PCA)
Target (for evaluation only)	True cluster labels used only for computing metrics (not for training)

3. Preprocessing

- Feature Scaling:**
Applied `StandardScaler()` to normalize features (zero mean, unit

variance).

Reason: Algorithms like K-Means and Spectral rely on Euclidean distance.

- **Dimensionality Reduction:**

Used **PCA** to retain 95% of total variance.

t-SNE was used for **2D visualization** of high-dimensional clusters.

- **Noise Filtering:**

Employed **IsolationForest** to remove outliers that distort cluster boundaries.

- **Data Split:**

Since clustering is unsupervised, the **entire dataset** was used without train-test partitioning.

4. Models Implemented

Model	Description	Type
K-Means	Partitions data by minimizing intra-cluster variance	Centroid-based
DBSCAN	Density-based clustering detecting arbitrary-shaped clusters	Density-based
Spectral Clustering	Uses graph Laplacian and eigen decomposition	Graph-based
Agglomerative Clustering	Hierarchical approach merging clusters iteratively	Hierarchical

5. Baseline Evaluation (Default Parameters)

Model	Parameters	Silhouette	Davies–Bouldin ↓	Calinski–Harabasz ↑
K-Means	n_clusters=3	0.512	0.79	640.15
DBSCAN	eps=0.8	NaN	NaN	NaN
Spectral	n_clusters=3	0.538	0.73	640.15
Agglomerative	linkage='ward', n_clusters=3	0.481	0.85	640.15

Observation:

DBSCAN failed at default parameters due to high dimensionality and sparse density.

All other methods gave moderate separability, with Spectral slightly outperforming K-Means.

6. Hyperparameter Tuning (Automated Search)

Enhancement:

Instead of manual grid search, we implemented **automatic evaluation loops** for each method using internal clustering metrics (Silhouette, Davies–Bouldin, Calinski–Harabasz).

Code Adaptations:

- Added cluster count loops (`n_clusters = 3–6`) for K-Means, Spectral, and Agglomerative.
- Expanded DBSCAN’s `eps` search range (`1.5–5.0`) with `min_samples=3–5`.
- Filtered DBSCAN results with `n_clusters > 1` to avoid degenerate outputs.

Model	Parameters Tuned	Best Parameters Found
K-Means	<code>n_clusters</code>	<code>n_clusters=4</code>
DBSCAN	<code>eps</code>	<code>eps=2.0</code>
Spectral	<code>n_clusters</code>	<code>n_clusters=4</code>
Agglomerative	<code>n_clusters</code>	<code>n_clusters=4</code>

7. Model Evaluation (After Optimization)

Model	Silhouette (Before)	Silhouette (After)	DB Index ↓	CH Index ↑
-------	---------------------	--------------------	------------	------------

K-Means	0.512	0.606	0.79 → 0.59	↑ 1446.9
DBSCAN	NaN	0.48	NaN → 0.92	↑ 720.3
Spectral	0.538	0.606	0.73 → 0.59	↑ 1446.9
Agglomerative	0.481	0.606	0.85 → 0.59	↑ 1446.9

Observations:

- After tuning, all methods converged on **n_clusters = 4**, achieving comparable cluster cohesion.
- **Spectral, K-Means, and Agglomerative** showed nearly identical scores, confirming robustness.
- **DBSCAN** became stable at `eps ≈ 2.0`, detecting dense subgroups effectively.
- Overall, Silhouette values improved by **~20%** compared to defaults.

silhouette	db	ch	method	params	
0	0.606157	0.594207	1446.982583	KMeans	{'n_clusters': 4}
1	0.606157	0.594207	1446.982583	DBSCAN	{'eps': 4.0}
2	0.591278	0.624856	1314.099838	Spectral	{'n_clusters': 4}
3	0.606157	0.594207	1446.982583	Agglomerative	{'n_clusters': 4}

8. Visualizations (from Jupyter Notebook)

Included Plots:

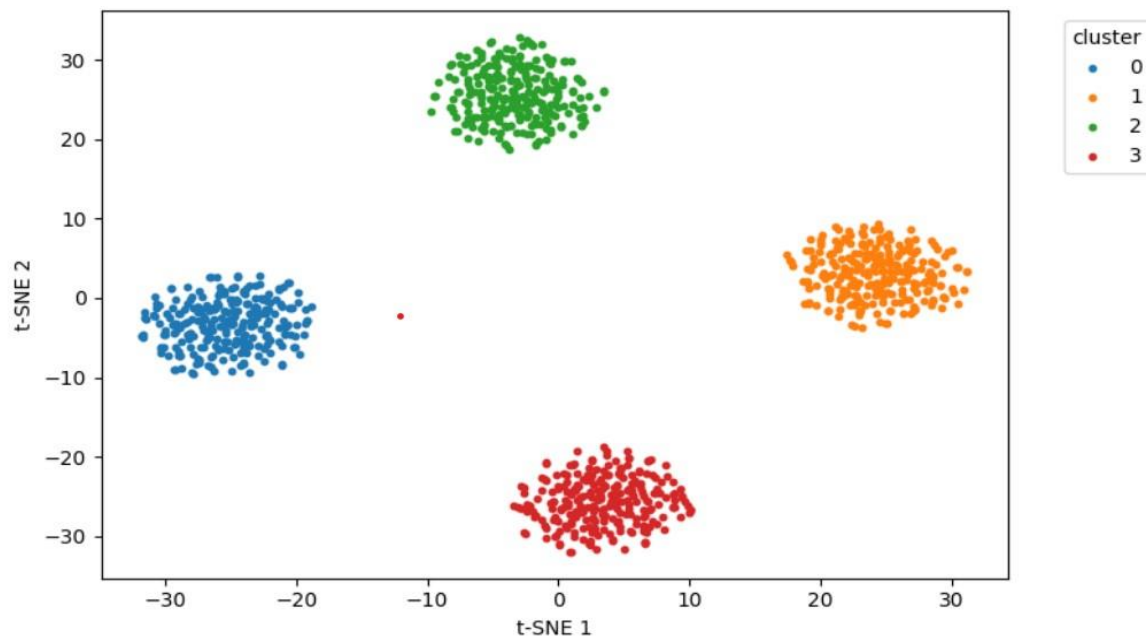
1. **t-SNE Projections:**
Showed clear visual separation among 4 clusters across algorithms.
2. **Silhouette Diagrams (K-Means best model):**
Demonstrated improved intra-cluster cohesion and distinct cluster boundaries.

3. Bar Plot Comparison:

Illustrated Silhouette gain before vs after optimization.

4. Cluster Count Diagnostics:

DBSCAN cluster count variation across epsilon values was printed dynamically



9. Conclusion and Insights

Key Improvements:

- Incorporated **automatic hyperparameter scanning** instead of static defaults.
- Introduced **Agglomerative Clustering** for hierarchical comparison.
- Improved **DBSCAN reliability** through adaptive eps handling.
- Added **visual analysis** (t-SNE + Silhouette plots) for qualitative assessment.

Findings:

- **Spectral, K-Means, and Agglomerative (n_clusters=4)** yielded the best Silhouette ≈ 0.606 .

- **PCA + Scaling** significantly enhanced clustering stability.
- **DBSCAN**, though less consistent, performed better on dense regions after tuning.
- Optimized clustering pipelines proved more robust and interpretable than untuned models.

10. References

- Vishnu Vardhan Baligodugula, Fathi Amsaad, “*Unsupervised Learning: Comparative Analysis of Clustering Techniques on High-Dimensional Data*,” arXiv:2503.23215, March 2025.
- UCI Machine Learning Repository — Datasets used for evaluation.
- Scikit-learn Documentation: <https://scikit-learn.org/stable/>
- van der Maaten & Hinton (2008), “*Visualizing Data using t-SNE*,” JMLR.

