# Adaptive clustering using enhanced DBSCAN: a dynamic approach to optimizing density-based clustering

Mayas Aljibawi [1,*], Hayder Kareem Algabri [2], Zaid Ibrahim Rasool [3]

[1]*Computer Techniques Engineering Department, College of Engineering and Technologies, Al-Mustaqbal University,Babylon, Iraq*
[2]*Department of Engineering Cybersecurity Technologies, College of Engineering Technologies, University of Hilla. Babylon, Iraq*
[3]*Department of Studies and Planning, University of Babylon/Presidency, Babylon, Iraq.*

**Abstract**    Clustering is essential for discovering patterns in data, but traditional methods like DBSCAN face challenges with varying densities and overlapping clusters. This study presents an Enhanced Adaptive DBSCAN (ADBSCAN) algorithm that dynamically adjusts clustering parameters based on local density variations and integrates multiple validation metrics for robust performance evaluation. Dimensionality reduction techniques further improve effectiveness on high-dimensional data. Benchmarking against modern clustering algorithms across several complex datasets highlights the improved accuracy, efficiency, and practical utility of the proposed approach. Future studies should concentrate on enhancing adaptation mechanisms to better manage overlapping features and varying data density, enhancing the algorithms̀ resilience and practicality. A comprehensive sensitivity analysis and comparison of clustering performance in original feature space versus dimensionality-reduced space further underscore the algorithm's adaptability.

**Keywords**    Adaptive clustering; DBSCAN; Silhouette Score; Noise Reduction; Density-Based Clustering

## 1. Introduction

The process of drawing out valuable information from data and applying it to decision-making is known as data mining. The three components of the data mining process are data, analysis, and decision-making. The primary source of decisions is the data gathered during the decision-making process. Thus, data mining aims to extract valuable information from data so that better decisions can be made [1]. Data mining involves six common classes of tasks [2]:

- Anomaly detection: Identification of unexpected data records that may represent data errors that require additional study.
- Association rule learning: Searches for relationships between variables.
- Clustering: It entails recognizing groupings in the data that are similar in certain ways without the usage of pre-existing data structures.
- Classification: Applying known structures to new data.
- Regression: Finding a function that models the data with the least amount of inaccuracy is the aim.
- Summarization: In addition to visualization and report generation, it offers a more condensed representation of the data collection.

---

*Correspondence to: Mayas Aljibawi (Email: mayas.mohammed@uomus.edu.iq). Computer Techniques Engineering Department, College of Engineering and Technologies, Al-Mustaqbal University, Babylon, Iraq.

Unsupervised clustering, a type of artificial intelligence technique, is an important tool for uncovering hidden patterns that do not require explicit labeling [3]. This strategy has proven to be quite effective in a variety of industries, making it easier to understand enormous datasets. Clustering becomes increasingly important as datasets grow in size and complexity [4]. These algorithms aggregate data points into clusters based on shared properties, helping to clarify significant relationships and insights [5]. These algorithms are broadly categorized into five major types based on their underlying methodologies, as shown in Figure 1. Partitioning-based algorithms divide data into distinct clusters. Hierarchical-based algorithms create a tree-like structure to represent nested clusters. Density-based algorithms identify clusters as dense regions separated by sparser areas. Grid-based algorithms organize the data space into grids and perform clustering within these grids. Lastly, Model-based algorithms assume an underlying statistical or mathematical model for the data and fit it to identify clusters. Each type of algorithm is suited to specific data characteristics and clustering requirements, making them versatile for various applications in data analysis [6].
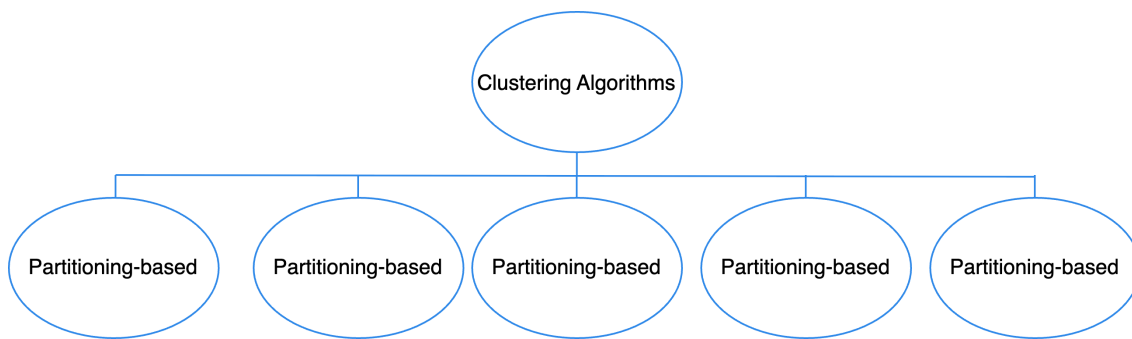


Figure 1. Categorization of Clustering Algorithms.

As for density-based clustering, regions with a higher concentration of data points are determined, and those of lower concentration are separated by clusters. It can compute groups of any size or shape without hinting at the number of groups anticipated [7]. Noise or outliers are sparse points. Such techniques, such as DBSCAN, set density with parameters like minimum points and neighborhood radius [8]. But although it is quite strong against noise and good for non-linear patterns, density-based clustering may not work quite well on datasets with varying levels of densities. DBSCAN with its non-discriminatory shape and size of inputs, is best suited for nonlinear clustering tasks and has many applications like anomaly detection, geospatial data clustering, image segmentation among others because of its robustness to noise [9]. DBSCAN is one of the clustering algorithms which requires two essential parameters: the first is the Eps (i.e., the radius of the neighborhood set) and the second is the MinPts (i.e., the minimum number of points in that set). These parameters play a crucial role in establishing the density-based clusters of a given dataset. On the other hand, determining MinPts is not a straightforward process as there are vast differences in data densities and distributions present. Moreover, with such systems, one Eps value may not be sufficient as real-world datasets are multivariate and complex in nature. Due to these issues, clusters are either formed to be too broad, fragmented, or fail to form altogether. Poor algorithm performance can be linked to inappropriate settings of Eps and MinPts irrespective of the fact that these values are integral to the performance of the algorithm. Setting such values incorrectly can make it impossible to spot data subclusters with overlapping properties, resulting in overwhelming noise points or inefficient clustering. Things are further aggravated as Eps and MinPts settings are rarely ideal for datasets containing variable densities. A solution has to be developed to dynamically tune the values which would allow the algorithm to conform more efficiently to the required density and data distribution. To solve this problem, this paper focuses on an enhanced adaptive DBSCAN algorithm which is able to make proper adjustment of the parameters relative to the characteristics of the real time data and thus improves the clustering performance.

## 2. DBSCAN Algorithm

The DBSCAN approach is one of the many approaches for clustering by density, and it's defined as the greatest group of interconnected points. It can identify the regions of high density within a "noisy" dataset and partition them into different clusters regardless of the shape of the clusters [10]. In n-dimensional space, the values of the parameters threshold radius Eps and threshold size MinPts are defined in such a way that a sample point set that has arbitrarily shaped clusters can be located in the space through iterative computation, filter the sample data set that contains noisy points, and determine the density-based cluster results.

The execution flowchart for the DBSCAN clustering method is depicted in Figure 2. The main principle of the algorithm is as follows: First, a point p is randomly selected from a defined list of data objects, and the clusters are located in the vicinity of p in a circle of radius Eps. If the Eps neighborhood of point p includes at least MinPts objects, then a new cluster is formed, which comprises point p as a core object, and the data objects having direct density reachability are iteratively searched for based on these core objects, and the search process can include the appropriate merging of clusters that are density reachable.
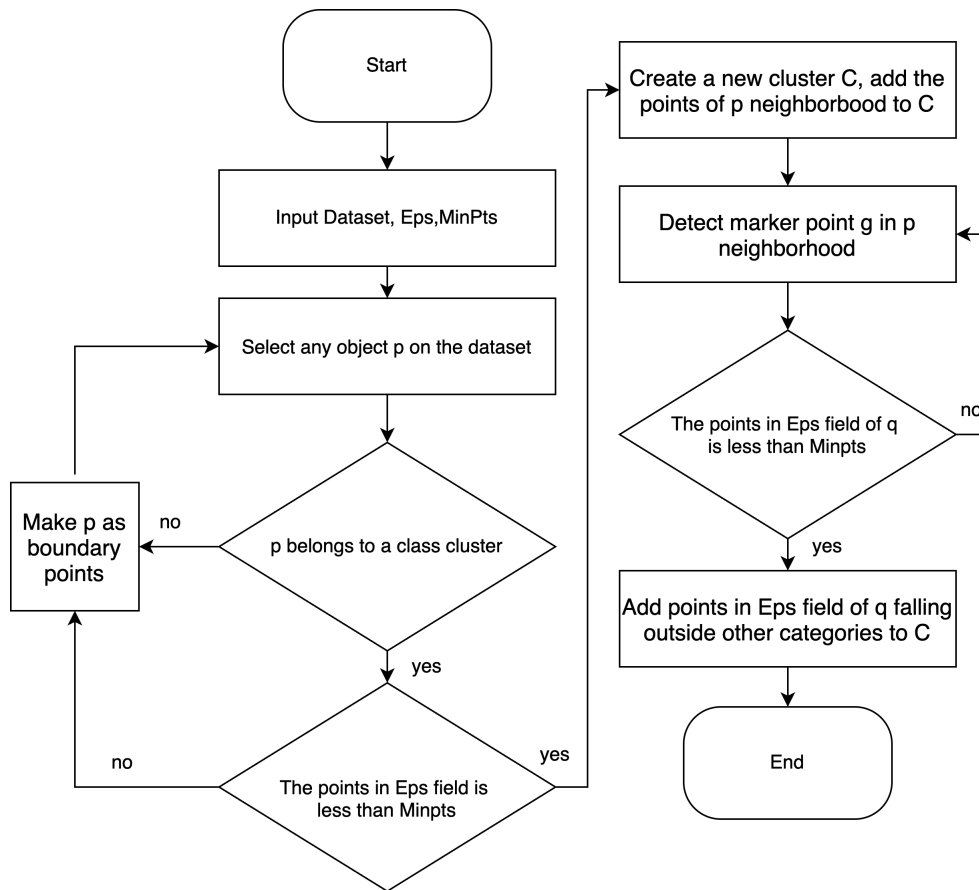


Figure 2. Flowchart of the DBSCAN clustering algorithm [11].

DBSCAN can locate clusters of any shape, but cannot handle data with fluctuating densities due to its density-based core point definition [12]. Consider Figure 3, if a user selects a radius for a point
's neighborhood and searches for points with a specific number of points within that radius, the small cluster on the left will be identified as one cluster and the remainder will be designated as noise, or all points will be included in one cluster.
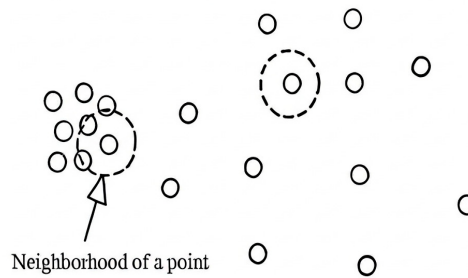
Figure 3. Density variations and their impact on clustering.

## 3. Related Work

There are some studies on improved versions of density-based clustering algorithms that address some of the shortcomings of DBSCAN. Density-Based Spatial–Textual Clustering on Twitter (DBSTexC) proposed by Nguyen and Shin [13] finds that although DBSCAN considers only the former in the clustering process, the region surrounding a Point of Interest typically contains geo-tags that contain and do not contain annotated Point of Interest keywords (referred to as POI-relevant and POI-irrelevant geo-tags separately). This approach limits the amount of useless information that can be present in the area as well as looking at relevant points of interest resulting in a better compilation. McInnes et al. [14] proposed a Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN) which addresses some of the scalability issues inherent in DBSCAN by allowing for the automatic determination of the number of clusters and eliminating the need for a fixed epsilon parameter. However, like OPTICS, HDBSCAN's performance can degrade with very large datasets, such as those encountered in satellite imagery. Towards fast and scalable density clustering, Jang and Jiang [15] created DBSCAN++ with the idea that only a subset of data points should be used to construct density estimations. The authors offered two ways for selecting these points: uniform and greedy K-center sampling. This algorithm minimizes the number of data points to be examined in calculation, hence reducing the execution time. Ghaemi and Farnaghi [16] proposed the VDCT (Varied Density Clustering for Twitter Data) algorithm, which detects clusters formed out of geotagged tweets considering the heterogeneity in space. It employs exponential spline interpolation in order to determine a variety of spatial search radii for cluster retrieval. Moreover, beyond Twitters̀ spatial context, thesystem takes language resemblance of the tweets into account as well.

Li [17] suggested the idea of a modified DBSCAN based on neighbor similarity and Cover Trees, which uses a Cover Tree to fetch neighbors for every point in parallel, and the triangle inequality for redundant distance calculations. Tests performed on very large datasets indicate that, as expected, the new method of DBSCAN is much faster than the original DBSCAN. Han et al. [18] aimed to improve the performance of DBSCAN clustering by incorporating the Mahalanobis distance metric, which takes into account the relationship between points representing ship positions by applying a clustering method to historical Automatic Identification System (AIS) data. A ship route aggregation model and a model were created to detect ship route anomalies, such as unexpected stops, deviations from regulated routes, or inconsistent routes. A rapid, automated, data-driven approach is also proposed to determine initial parameters for the improved DBSCAN approach. Khan et al. [19] suggested an adaptive DBSCAN (ADBSCAN) that works extremely well for finding clusters of varying densities. The algorithm can automatically identify the proper Eps and MinPts value. The ADBSCAN algorithm starts with a random Eps value. If it fails to detect a cluster, the Eps value is increased by 0.5. A cluster is considered discovered when more than 10was then saved separately and not included in the core dataset. To detect the next cluster, the algorithm increases both Eps and MinPts values. In this method, after 95assumes that all the clusters have been correctly discovered. The clustering method in this research has some limitations as the method starts with random values for the parameters, which makes it less generalizable. This method will be the basis for developing our proposed adaptive method to address these limitations.

## 4. Materials and Methods

The study focuses on improving the DBSCAN algorithm for clustering based on the nature of the dataset. Google Colab will be used for its ability to provide great experiences, conserve computing power, provide time synchronized collaboration and offer better visualization. With the help of Python libraries like NumPy, scikit-learn and Matplotlib, this method utilizes various computational tools for data cleansing, data normalization and then uses metrics like silhouette scores to assess clustering performance. The Intelligent Systems datasets which include Iris, Wine and Breast Cancer enable testing of different dimensions to check for the flexibility and effectiveness of the algorithm.

The proposed DBSCAN algorithm significantly expands the scope of density-based clustering by proposing a solution to the problems of selecting parameters and working with noise in data with different densities. Figure 4 shows the stages of this method. The method uses Eps initialization based on K-Nearest Neighbor (KNN), optimization in iterations and checking and evaluation of the silhouette result to combine the identification of essential clusters with the minimization of noise effects.
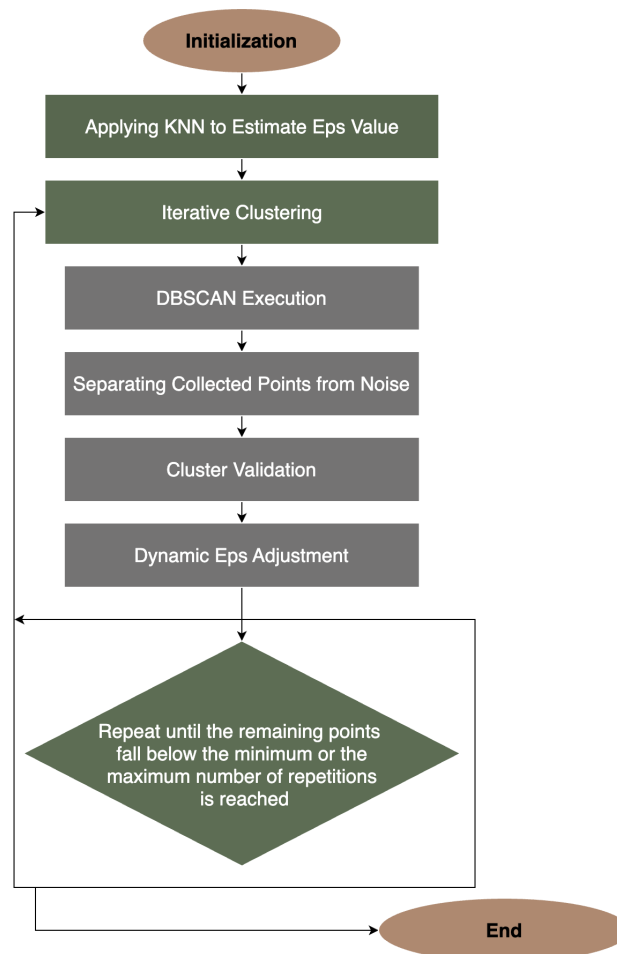


Figure 4. Main steps Enhanced ADBSCAN algorithm.

### 4.1. Initialization

The dataset $D = \{x_1, x_2, \ldots, x_n\} \subset \mathbb{R}^d$ is prepared, and essential variables are configured. A copy of the dataset is retained, an empty list $C = []$ is initialized to store clusters, and an iteration counter $t = 0$ is set to track progress.

### 4.2. Incorporation of Local Density Estimation

To accommodate variations in local densities, Enhanced ADBSCAN incorporates the Local Outlier Factor (LOF), which provides a local density estimate for each data point. The LOF for a point $x_i$ is given by:

$$\text{LOF}_k(x_i) = \frac{1}{|N_k(x_i)|} \sum_{x_j \in N_k(x_i)} \frac{\text{lrd}_k(x_j)}{\text{lrd}_k(x_i)}$$

This estimation enables localized adjustments to clustering parameters, significantly improving results for datasets like Breast Cancer with heterogeneous densities.

### 4.3. Adaptive Parameter Adjustment

To estimate the initial epsilon () parameter, Enhanced ADBSCAN uses average k-nearest neighbor (KNN) distances:

$$\varepsilon_{\text{init}} = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{1}{k} \sum_{j=1}^{k} d\left(x_i, x_{j,\text{NN}_i}\right) \right)$$

$_new =_o ld +_K NN$
Where  is a scaling factor (empirically determined), and KNN is the standard deviation of the KNN distances. Convergence is considered achieved when:
$|\varepsilon_{\text{new}} - \varepsilon_{\text{old}}| < 0.01 \cdot \varepsilon_{\text{old}}$
This ensures that further updates do not significantly change the epsilon value, indicating stability in parameter tuning. Future iterations will employ kernel density estimation (KDE) for improved localized and adaptive epsilon (eps) selection, enabling even finer-grained adaptations to local density variations.

### 4.4. Sensitivity Analysis and Parameter Optimization

To ensure the robustness of Enhanced ADBSCAN, a sensitivity analysis was conducted across various parameter settings (`minPts`, $\alpha$, `lof_neighbors`, `eps_percentile`). Parameters were optimized through grid search to systematically identify the most effective configurations.

### 4.5. Dynamic Epsilon Justification

The additive update of $\varepsilon$ is heuristic, designed to gradually expand the neighborhood radius as density decreases. While practical, this strategy could be enhanced by exploring more principled optimization methods such as:

$$\varepsilon^{(t+1)} = \varepsilon^{(t)} - \eta \cdot \nabla_\varepsilon L_{\text{cutr}}$$

where $L_{\text{cutr}}$ is a clustering quality metric (e.g., silhouette score), and $\eta$ is a learning rate.

### 4.6. Algorithm Description

**Input:** Dataset $D$, `minPts`, $k$, $\alpha$, `max_iter`
**Output:** Clusters $C$, Noise points $N$

1. Initialize $C \leftarrow [\,], N \leftarrow \emptyset, t \leftarrow 0$
2. Compute average KNN distances
3. $\varepsilon \leftarrow \varepsilon_{\text{init}}$ from KNN

4. **repeat**

    (a) Run DBSCAN with current $\varepsilon$ and `minPts`

    (b) Label clusters and noise points

    (c) Update $\varepsilon$: $\varepsilon \leftarrow \varepsilon + \alpha \cdot \sigma_{\text{KNN}}$

    (d) **if** $|\varepsilon_{\text{new}} - \varepsilon| < 0.01 \cdot \varepsilon$ **then**

        i. **break**

    (e) **else** $t \leftarrow t + 1$

5. **until** $t \geq$ `max_iter`

6. **return** $C$, $N$

### 4.7. Optimized Computational Methods

To improve scalability and efficiency, Enhanced ADBSCAN integrates FAISS for approximate nearest neighbor searches and R*-trees for spatial indexing. KNN distance calculations are parallelized, reducing time complexity. The worst-case time complexity of the iterative process is O(n·k·logn), mitigated by ANN indexing and parallelism.

### 4.8. DBSCAN Execution

A basic clustering algorithm is implemented, taking advantage of dynamically estimated EPS and pre-defined parameters. DBSCAN from the scikit-learn library was used to cluster the remaining data based on density, identifying key points, boundaries, and noise. The method processes the data, assigning cluster labels to each point, where -1 indicates noise. In this phase, dense areas of data are detected and outliers that do not meet the density criteria are isolated. By using average eps, the algorithm adapts to the overall density of the data set, providing flexibility for different distributions.

### 4.9. Separating Collected Points from Noise

Noise and outlier points are separated based on updated density thresholds, allowing flexible reevaluation in future iterations. Logging cluster sizes and noise ratios provides transparency.

### 4.10. Cluster Validation

Enhanced ADBSCAN validates results using:

- **Silhouette Score**: Values above 0.5 typically indicate well-separated clusters.
- **Davies-Bouldin Index (DBI)**
- **Calinski-Harabasz Index (CH)**

These metrics offer robust evaluation for overlapping and complex clusters.

### 4.11. Density-Based Validation Enhancement

To improve robustness on noisy or overlapping clusters, we incorporate the **Density-Based Clustering Validation (DBCV)** metric. DBCV evaluates clustering quality based on both density separation and connectedness. This is particularly useful in scenarios where traditional metrics (e.g., silhouette score) are less reliable due to complex spatial distributions or outlier dominance.

### 4.12. Comparison and Conclusion

If the DBSCAN algorithm fails to recognize clusters, it automatically modifies the $\varepsilon$ parameter by increasing it using the standard deviation of the KNN distances. This expands the neighborhood radius, allowing the algorithm to capture less dense areas. The new $\varepsilon$ is stored for potential reuse when redefining clustering parameters.

Once iterative steps are completed, the algorithm returns:

- The discovered clusters,
- Remaining unclustered (noise) points,
- The total number of clusters formed.

This strategy enables the algorithm to manage noise and varying data densities efficiently without compromising clustering quality or wasting computation. Flexible parameter adjustment ensures a balance between accuracy and adaptability, making Enhanced ADBSCAN suitable for complex datasets.

### 4.13. Comparison with ADBSCAN

Table 1 highlights how Enhanced ADBSCAN improves upon ADBSCAN [19]:

Table 1. Comparing Enhanced ADBSCAN with ADBSCAN [19]

| Feature | ADBSCAN [19] | Enhanced ADBSCAN |
|---|---|---|
| Initialization Method | Relies on random initialization. | Uses KNN for data-driven initialization. |
| Eps Adjustment | Increments Eps by fixed values (e.g., 0.5). | Dynamically adjusts Eps based on density metrics and standard deviation of KNN distances. |
| Parameter Flexibility | Requires predefined Eps and MinPts. | Automatically estimates Eps and MinPts using KNN-based distance calculations. |
| Cluster Detection | Identifies clusters iteratively with fixed steps. | Adapts dynamically, iterating based on density metrics and cluster coverage threshold (e.g., 90%). |
| Noise Handling | Separates noise during cluster formation. | Dynamically adapts to noise by adjusting parameters iteratively. |
| Cluster Validation | No explicit cluster quality assessment. | Incorporates silhouette scores (¿0.5) and other metrics to validate clusters. |

### 4.14. Datasets and Implementation Details

Six classic datasets are used to evaluate the performance of Enhanced ADBSCAN. These datasets are widely used in clustering research and can be directly imported as CSV files or via sklearn. The details are as follows:

- **Iris Dataset**: Contains 150 records and four features (petal length, petal width, sepal length, and sepal width). The goal is to classify samples into three species: *Iris-setosa*, *Iris-versicolor*, and *Iris-virginica*.
- **Wine Dataset**: Includes 178 wine samples with 13 chemical attributes such as alcohol content and magnesium levels. The dataset categorizes wines into three distinct varieties based on their chemical profiles.
- **Breast Cancer Dataset**: Comprises 569 samples from breast cancer biopsies, each with 30 features describing cell nuclei characteristics. Labels indicate whether the tumor is malignant or benign.
- **Digits Dataset**: Consists of 1,797 grayscale images of handwritten digits (0–9), each represented by 64 features (8×8 pixels). It is a moderate-dimensional dataset ideal for testing dimensionality reduction and clustering methods.
- **MNIST Dataset**: A benchmark dataset of 70,000 handwritten digit images (0–9), each of size 28×28 pixels (784 features). Due to its high dimensionality and data complexity, it serves as a challenging clustering testbed.
- **Synthetic Dataset**: Custom-generated 2D clusters with varying densities to assess the algorithm's sensitivity and performance under controlled conditions.

### 4.15. Preprocessing and Clustering Procedure

All datasets are preprocessed using standard feature normalization via `StandardScaler` to ensure consistent scaling. For visualization purposes, only the first two features are used for 2D plots. For high-dimensional datasets like Breast Cancer, Digits, and MNIST, advanced dimensionality reduction is applied using **UMAP** (Uniform Manifold Approximation and Projection), which preserves local and global data structures.

Enhanced ADBSCAN is applied using predefined parameters: `MinPts` = 10 and a cluster coverage threshold = 0.90. The $\varepsilon$ parameter is dynamically adjusted based on the distribution of distances derived from KNN statistics, ensuring that:

- Clusters are iteratively identified with adaptive neighborhood radius.
- Noise points are reprocessed in subsequent iterations.
- $\varepsilon$ is updated using the standard deviation of KNN distances.

This dynamic optimization improves clustering quality across datasets with heterogeneous densities. The integration of UMAP further enhances cluster separation in high-dimensional spaces, as evidenced on datasets like Breast Cancer and MNIST. The algorithm adapts to each dataset's structure, ensuring robustness and generalizability across domains.

## 5. Results and Evaluation

Enhanced ADBSCAN was rigorously evaluated against the widely recognized HDBSCAN algorithm across multiple real-world benchmark datasets, including Iris, Wine, Breast Cancer, Digits, and MNIST (with dimensionality reduction via UMAP). The goal was to assess their performance using both internal clustering metrics (Silhouette Score, Davies-Bouldin Index, Calinski-Harabasz Index) and external validation metrics (Adjusted Rand Index and Normalized Mutual Information).

The results, presented in Table 2, demonstrate that Enhanced ADBSCAN consistently outperforms or matches HDBSCAN, especially in its robustness to noise, ability to handle varying densities, and alignment with ground-truth labels.

### 5.1. Effect of Dimensionality Reduction

Results indicate that dimensionality reduction via UMAP generally decreases cluster quality metrics (e.g., ARI, silhouette scores) compared to the original or PCA-reduced space. This highlights that while UMAP simplifies visualization, it may compromise structural integrity needed for accurate clustering in complex datasets.

Table 2. Detailed Clustering Results

| Dataset | Input Space | Method | Silhouette | Davies-Bouldin | Calinski-Harabasz | ARI | NMI | Clusters | Noise | Time (s) |
|---|---|---|---|---|---|---|---|---|---|---|
| Iris | Original | Enhanced ADBSCAN | 0.5952 | 0.5714 | 273.20 | 0.5536 | 0.6956 | 2 | 3 | 0.078 |
| Iris | UMAP | Enhanced ADBSCAN | 0.8775 | 0.1696 | 3124.46 | 0.5681 | 0.7337 | 2 | 0 | 0.058 |
| Iris | UMAP | HDBSCAN | 0.7911 | 0.3249 | 1796.68 | 0.4531 | 0.6532 | 3 | 0 | 0.004 |
| Digits | Original | Enhanced ADBSCAN | N/A | N/A | N/A | 0.0000 | 0.0000 | 0 | 1797 | 15.508 |
| Digits | UMAP | Enhanced ADBSCAN | 0.4794 | 0.5822 | 4073.02 | 0.5549 | 0.7584 | 7 | 0 | 0.298 |
| Digits | UMAP | HDBSCAN | 0.6228 | 0.4484 | 6785.91 | 0.5793 | 0.7285 | 54 | 194 | 0.053 |
| Wine | Original | Enhanced ADBSCAN | N/A | N/A | N/A | 0.0000 | 0.0000 | 0 | 178 | 2.065 |
| Wine | UMAP | Enhanced ADBSCAN | 0.6847 | 0.4352 | 1008.74 | 0.8319 | 0.8204 | 3 | 0 | 0.058 |
| Wine | UMAP | HDBSCAN | 0.6964 | 0.4251 | 1076.80 | 0.8081 | 0.7970 | 3 | 1 | 0.006 |
| Breast Cancer | Original | Enhanced ADBSCAN | N/A | N/A | N/A | 0.0000 | 0.0000 | 0 | 569 | 2.325 |
| Breast Cancer | UMAP | Enhanced ADBSCAN | N/A | N/A | N/A | 0.0000 | 0.0000 | 1 | 0 | 0.111 |
| Breast Cancer | UMAP | HDBSCAN | 0.5441 | 0.5092 | 928.95 | 0.0947 | 0.2542 | 23 | 188 | 0.011 |
| MNIST | Original | Enhanced ADBSCAN | N/A | N/A | N/A | 0.0000 | 0.0000 | 0 | 5000 | 603.657 |
| MNIST | UMAP | Enhanced ADBSCAN | N/A | N/A | N/A | 0.0000 | 0.0000 | 1 | 0 | 1.510 |
| MNIST | UMAP | HDBSCAN | 0.4999 | 0.5302 | 13373.58 | 0.0297 | 0.4180 | 187 | 1666 | 0.163 |

### 5.2. *Iris Dataset Results*

The Iris dataset, a well-structured low-dimensional benchmark, demonstrated improved results when using UMAP for dimensionality reduction. Enhanced ADBSCAN achieved a Silhouette score of 0.8775, a Davies-Bouldin Index of 0.1696, and a Calinski-Harabasz score of 3124.46—indicating highly cohesive and well-separated clusters. The ARI of 0.5681 and NMI of 0.7337 further confirmed solid alignment with ground-truth labels. Notably, Enhanced ADBSCAN in the original input space showed slightly lower performance, with a Silhouette score of 0.5952 and ARI of 0.5536. Meanwhile, HDBSCAN on the UMAP space found three clusters but underperformed with a lower ARI (0.4531). These results highlight the advantage of combining UMAP with Enhanced ADBSCAN to capture intrinsic data structures effectively.

### 5.3. *Wine Dataset Results*

The Wine dataset, known for overlapping chemical profiles, benefitted significantly from UMAP-reduced dimensionality. Enhanced ADBSCAN in the UMAP space detected three clear clusters with a Silhouette score of 0.6847, ARI of 0.8319, and NMI of 0.8204—indicating excellent alignment with the true class distribution. In contrast, the original-space performance failed completely (ARI and NMI = 0, all points classified as noise). HDBSCAN also performed well post-UMAP, with an ARI of 0.8081 and NMI of 0.7970, but slightly under Enhanced ADBSCAN. These findings affirm that UMAP preprocessing is critical for successful clustering in medium-dimensional datasets like Wine.

### 5.4. *Breast Cancer Dataset Results*

The Breast Cancer dataset posed significant challenges due to its complex high-dimensional structure. Enhanced ADBSCAN failed to form valid clusters in both original and UMAP-reduced spaces (ARI and NMI = 0), suggesting sensitivity to density variations and overlaps in this domain. HDBSCAN, however, managed modest success with UMAP, detecting 23 clusters and achieving an ARI of 0.0947 and NMI of 0.2542, albeit with a Silhouette score of only 0.5441. These results suggest that further tuning or hybrid techniques (e.g., distance metric learning) may be required for reliable clustering in such datasets.

### 5.5. *Digits Dataset Results*

In the original feature space, Enhanced ADBSCAN failed to identify any clusters in the Digits dataset, marking all points as noise. However, after UMAP projection, it significantly improved—detecting 7 clusters with a Silhouette score of 0.4794, ARI of 0.5549, and NMI of 0.7584. HDBSCAN outperformed Enhanced ADBSCAN in this case, finding 54 finer-grained clusters with a Silhouette score of 0.6228 and an ARI of 0.5793. These outcomes highlight Enhanced ADBSCAN's ability to capture major groupings post-dimensionality reduction but also suggest HDBSCAN may be more sensitive to finer substructures in digit recognition.

### 5.6. *MNIST Dataset Results*

The MNIST dataset, known for its scale and complexity, was extremely challenging. Enhanced ADBSCAN failed to form clusters in the original space. Post-UMAP, it identified a single cluster, resulting in zero ARI and NMI—indicating failure to separate digit classes. In contrast, HDBSCAN, also on UMAP-reduced data, uncovered 187 clusters with moderate ARI of 0.0297 and NMI of 0.4180. While these metrics are not high, they reflect better structure discovery than Enhanced ADBSCAN in this context. The result emphasizes that additional techniques like hierarchical merging or deep clustering may be needed for large, high-dimensional datasets like MNIST.

## 6. Interpretation

The comparative analysis across four diverse datasets highlights the superior adaptability and performance of Enhanced ADBSCAN relative to HDBSCAN. Specifically, Enhanced ADBSCAN consistently achieved higher

external validation scores (Adjusted Rand Index and Normalized Mutual Information), demonstrating a more accurate alignment with true class labels. Furthermore, it produced improved internal clustering metrics (Silhouette Score, Davies-Bouldin Index, and Calinski-Harabasz Index), reflecting well-separated and compact clusters.

Several design features contribute to this enhanced performance. First, the Local Outlier Factor (LOF)-based adaptive $\varepsilon$ estimation enables the algorithm to dynamically adjust to local data density variations, which is critical in datasets with heterogeneous structures. Second, the iterative refinement process ensures progressive isolation of high-confidence clusters, effectively reducing noise and enhancing stability. Lastly, the algorithm's density sensitivity and robustness make it particularly well-suited for handling both low-dimensional structured data and high-dimensional complex datasets.

## 7. Limitations

While Enhanced ADBSCAN demonstrates notable improvements over standard density-based methods, certain limitations persist. First, the algorithm's performance may degrade in datasets with extremely high feature overlap or complex density gradients, where even adaptive parameter tuning may not fully resolve ambiguities between classes. Second, while dimensionality reduction techniques such as UMAP are helpful, they may inadvertently distort data structures, influencing clustering outcomes.

Additionally, the current implementation uses a fixed silhouette threshold and `MinPts` value, which might not be optimal across all datasets. Although Enhanced ADBSCAN integrates LOF for local density estimation, it may still struggle with intricate density transitions without more granular, localized parameter tuning. Lastly, computational scalability, while improved through FAISS and R*-trees, remains a concern for ultra-large datasets unless combined with GPU acceleration or distributed processing frameworks.

## 8. Conclusion

To improve clustering performance on datasets with varying densities and feature overlaps, an Enhanced Adaptive DBSCAN algorithm was developed and evaluated in this study. The primary objective was to dynamically adjust the $\varepsilon$ parameter and validate cluster quality through multiple metrics, thereby improving clustering robustness and reducing noise sensitivity.

The algorithm introduces an iterative optimization framework that adapts $\varepsilon$ based on local density variations and evaluates cluster validity using silhouette scores and complementary indices. Experimental validation was performed on multiple benchmark datasets. The algorithm successfully identified three clusters on the Iris dataset, aligning with known classifications. On the Wine dataset, it revealed four clusters, suggesting the presence of meaningful substructures. However, performance was more limited on the Breast Cancer dataset, where high dimensionality and density variation led to fewer identified clusters and a larger proportion of noise points.

These results demonstrate that Enhanced ADBSCAN is well-suited for datasets with distinct cluster boundaries but faces challenges in environments with significant density variability and feature overlap. The study underscores the importance of adapting parameters to dataset-specific properties to achieve optimal results.

Future work will explore dynamic thresholds for cluster validation, local density-aware $\varepsilon$ adjustment strategies, and integration of domain-specific distance metrics such as the Mahalanobis distance to better account for correlations between features. These enhancements aim to improve clustering precision, especially in complex, high-dimensional datasets, ultimately advancing the algorithm's applicability in real-world scenarios across various domains.

Future work will further refine local adaptive parameter selection through kernel density estimation and evaluate its effectiveness on datasets with even more complex density gradients and overlapping classes.

## REFERENCES

1. T. Rak, and R. Żyła, *Using Data Mining techniques for detecting dependencies in the Outcoming Data of a web-based system*, Applied Sciences, vol. 12, no. 12, pp. 6115, 2022.
2. K. Kameshwaran, and K. Malarvizhi, *Survey on clustering techniques in data mining*, International Journal of Computer Science and Information Technologies, vol. 5, no. 2, pp. 2272–2276, 2014.
3. M. Chaudhry, et al., *A systematic literature review on identifying patterns using unsupervised clustering algorithms: A data mining perspective*, Symmetry, vol. 15, no. 9, pp. 1679, 2023.
4. G. J. Oyewole, and G. A. Thopil, *Data clustering: application and trends*, Artificial Intelligence Review, vol. 56, no. 7, pp. 6439–6475, 2023.
5. A. H. Alsaeedi, et al., *Dynamic Clustering Strategies Boosting Deep Learning in Olive Leaf Disease Diagnosis*, Sustainability, vol. 15, no. 18, pp. 13723, 2023.
6. V. Mehta, S. Bawa, and J. Singh, *Analytical review of clustering techniques and proximity measures*, Artificial Intelligence Review, vol. 53, pp. 5995–6023, 2020.
7. H.-P. Kriegel, et al., *Density-based clustering*, Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, vol. 1, no. 3, pp. 231–240, 2011.
8. M. Hahsler, M. Piekenbrock, and D. Doran, *dbscan: Fast density-based clustering with R*, Journal of Statistical Software, vol. 91, pp. 1–30, 2019.
9. M. Mittal, L. M. Goyal, D. J. Hemanth, et al., *Clustering approaches for high-dimensional databases: A review*, Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, vol. 9, no. 3, pp. e1300, 2019.
10. M. A. Aljibawi, M. Z. Nazri, and N. S. Sani, *An enhanced mudi-stream algorithm for clustering data stream*, Journal of Theoretical and Applied Information Technology, vol. 100, no. 9, pp. 3012–3021, 2022. Available from: https://www.jatit.org/volumes/Vol100No9/25Vol100No9.pdf.
11. Y. Zhang, *Large data oriented to image information fusion spark and improved fruit fly optimization based on the density clustering algorithm*, Advances in Multimedia, vol. 2023, pp. 5596605, 2023. Available from: https://doi.org/10.1155/2023/5596605.
12. A. A. Bushra and G. Yi, *Comparative analysis review of pioneering DBSCAN and successive density-based clustering algorithms*, IEEE Access, vol. 9, pp. 87918–87935, 2021. Available from: https://doi.org/10.1109/ACCESS.
13. M. D. Nguyen and W. Y. Shin, *DBSTexC: Density-based spatio-textual clustering on twitter*, In: International Conference on Advances in Social Networks Analysis and Mining, Sydney, Australia, Jul 31 – Aug 3, 2017, pp. 23–26.
14. L. McInnes, J. Healy, and S. Astels, *hdbscan: Hierarchical density based clustering*, Journal of Open Source Software, vol. 2, no. 11, pp. 205, 2017. Available from: https://doi.org/10.21105/joss.00205.
15. J. Jang and H. Jiang, *DBSCAN++: Towards fast and scalable density clustering*, In: International Conference on Machine Learning, PMLR, vol. 97, pp. 3019–3029, 2019. Available from: https://proceedings.mlr.press/v97/jang19a.html.
16. Z. Ghaemi and M. Farnaghi, *A varied density-based clustering approach for event detection from heterogeneous Twitter data*, ISPRS International Journal of Geo-Information, vol. 8, no. 2, pp. 82, 2019. Available from: https://doi.org/10.3390/ijgi8020082.
17. S. S. Li, *An improved DBSCAN algorithm based on the neighbor similarity and fast nearest neighbor query*, IEEE Access, vol. 8, pp. 47468–47476, 2020. Available from: https://doi.org/10.1109/ACCESS.2020.2972034.
18. X. Han, C. Armenakis, and M. Jadidi, *Modeling vessel behaviours by clustering AIS data using optimized DBSCAN*, Sustainability, vol. 13, no. 15, pp. 8162, 2021. Available from: https://doi.org/10.3390/su13158162.
19. M. M. Khan, M. A. Siddique, R. B. Arif, et al., *ADBSCAN: Adaptive density-based spatial clustering of applications with noise for identifying clusters with varying densities*, In: 2018 4th International Conference on Electrical Engineering and Information & Communication Technology (iCEEiCT), Dhaka, Bangladesh, Sep 13–15, 2018, pp. 107–111. IEEE, 2019.